

ASPETTI COMPUTAZIONALI NEI MODELLI STATISTICI LINEARI

L'analisi di regressione multipla

Si tratta di esaminare i problemi computazionali fondamentali per quanto riguarda l'analisi di regressione multipla.

Il modello è dato da:

$$\underset{(n \times 1)}{\underline{y}} = \underset{(n \times p)}{\underline{X}} \overset{(p \times 1)}{\underline{\beta}} + \underset{(n \times 1)}{\underline{\varepsilon}}$$

Lo stimatore per $\underline{\beta}$ è la soluzione delle equazioni normalizzate:

$$(\underline{X}'\underline{X})\hat{\underline{\beta}} = \underline{X}'\underline{y}$$

vale a dire:

$$\hat{\underline{\beta}} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{y}$$

Problema:

Analizzare gli aspetti computazionali che riguardano:

- Il calcolo dello stimatore $\hat{\beta}$
- la bontà di adattamento del modello di regressione
- una misura della variabilità di $\hat{\beta}$

Il problema dei minimi quadrati e le trasformazioni ortogonali

$\|x\|_2^2 = (x'x)$

(*) Se posso pensare da questo problema ad un problema equivalente utile sendo più comoda trasformazione. Le trasformazioni ortogonali.

I POTESI:

a) $E(\underline{\varepsilon}) = \underline{0}$

b) $E(\underline{\varepsilon}\underline{\varepsilon}') = \sigma^2 \underline{I}_n \Rightarrow E(\varepsilon_i \varepsilon_j) = \begin{cases} \sigma^2 & i=j \\ 0 & i \neq j \end{cases}$

c) X è un insieme di numeri certi

d) X ha rango $p < n$

\swarrow
n° dei parametri
da stimare

\searrow
n° osservazioni

Se trasformo, tramite le
trasformazioni ortogonali, il problema
originale in un problema
equivalente la cui soluzione è
identica all'originale me compoto
vantaggi numerici e computazionali

IL MODELLO DI REGRESSIONE MULTIPLA

Nel modello di regressione multiple la variabile dipendente Y è influenzata da più variabili indipendenti.

NB:

Se per il modello di regressione semplice ovvero per il modello di regressione con due variabili esplicative è ancora possibile calcolare le stime dei parametri senza ricorrere all'uso di un elaboratore nel caso di più di due variabili esplicative l'uso di tale strumento diventa indispensabile.

Obiettivo della regressione multiple
è quello di spiegare la variabile dipen-
dente y in funzione delle n variabili
indipendenti: X_1, \dots, X_n
ovvero di descrivere la dipendenza di
delle n variabili X_i mediante una f-
lineare delle n variabili X_i .

[dal punto di vista geometrico
ciò corrisponde ad un iper-piano
in un iperspazio ad $n+1$ dimensioni
(analoga al fatto che l'equazione
di regressione con una variabile indi-
pendente corrisponde ad una retta
nel piano)]

In generale i punti individuati
non giacciono che su un'unica ped
nell'iperspazio, ma si dispongono
per una serie di elementi pres.

errori di misurazione nelle variabili;
 variabili esplicative non incluse nel
 modello, fattori di non-linearità ecc.

⇒

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_p x_{pi} + \varepsilon_i$$

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

in forma compatta

$$\underline{y} = X \underline{\beta} + \underline{\varepsilon}$$

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_{21} & \dots & x_{p1} \\ 1 & x_{22} & \dots & x_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{2n} & \dots & x_{pn} \end{bmatrix}$$

(n \times p) ^{matrice} ~~matrice~~ ^{di} ~~variabili~~ ^{esplicative}

$$\underline{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$\underline{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

IPOTESI:

a)* $E(\underline{\varepsilon}) = \underline{0}$

b)* $E(\underline{\varepsilon} \underline{\varepsilon}') = \sigma^2 I_n = \text{Var}(\underline{\varepsilon}) = \begin{cases} \sigma^2 & i=j \\ 0 & i \neq j \end{cases}$

c)* X è un insieme di numeri certi

(l'unica fonte di variabilità è data dal vettore $\underline{\varepsilon}$)

d)* X ha rango $p < n \rightarrow n^\circ$ di osservazioni
 n° di parametri da stimare

STIMA DEI PARAMETRI DEL MODELLO

Si applichi il metodo dei minimi quadrati per stimare i parametri.

Sia $\underline{\hat{\beta}} = \{ \hat{\beta}_1, \dots, \hat{\beta}_p \}$

Vettore colonna composto dalle stime

Sotto le ipotesi prima scritte

Lo stimatore dei minimi quadrati
è corretto: $E(\hat{\beta}) = \beta$

è la varianza minima $Var(\hat{\beta}) = \sigma^2 (X'X)^{-1}$

→ TH. DI GAUSS MARKOV

Gli stimatori dei minimi quadrati
sono BLUE

PER HARE IL PRODOTTO:

$$(X'X)^{-1}(X'Y)$$

Si consideri le formule standard:

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

APPROCCIO PER I CALCOLI:

1) CALCOLARE $(X'X)$, $(X'Y)$

2) INVERTIRE $(X'X) \rightarrow (X'X)^{-1}$

3) FORMARE IL PRODOTTO:

$$(X'X)^{-1} (X'Y)$$

APPROCCIO VALIDO SE m e p
SONO PICCOLI.

APPROCCIO FACILE DA IMPLEMENTARE
(eccezione può derivare dall'inversione
della matrice) CHE HA PERÒ
ALCUNI LATI NEGATIVI.

OSS

•) $(X^T X)$

sono matrici derivanti da prodotti interni
di colonne e i prodotti interni
sono fonte di errore.

⇒ meglio un approccio che minimizza
i prodotti interni.

•) Un'altra fonte di errore e di inefficienza
computazionale deriva dal fatto che
vi è un peso in cui si deve effettuare
l'inversione di una matrice.

⇒ INSTABILITÀ