

Dimostrazione formula del calcolo di Beta cappello

Prima parte, gradiente

Nel contesto:

HULTPLE LINEAR REGRESSION

$$X^T \in (X_1, ..., X_p)$$
 potentialmente può influenzare

la dinumica di y

 $Y \in \{(X) \mid E = \beta_0 + \sum_{j=1}^p \beta_j \mid X_j \mid E \text{ linearita}$

Usando Least Squares, quindi puntando alla minimizzazione di:

Valgono le formule:

Essenzialmente abbiamo preso l'espressione e l'abbiamo scritta come matrici e vettori.

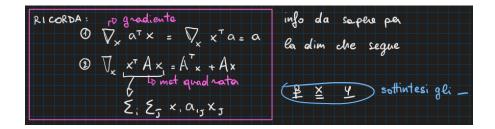
Nella matrice delle x, o matrice di design, le x hanno due pedici, il primo distingue i diversi campionamenti, il secondo indica il numero del regressore.

Dobbiamo trovare $\hat{\beta}$ che minimizza l'RSS:

RSS(
$$\beta$$
) = ξ ; $(\gamma; -(\underline{x}\beta);)^2 = (\underline{\gamma} - \underline{x}\beta)^T(\underline{\gamma} - \underline{x}\beta)$

by calcoli i quadrati

Se ho un vettore e voglio calcolarne i quadrati facciamo riga per colonna, la sommatoria la riscriviamo. Ricordando che:



Dimostriamo la formula:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Per trovare il minimo RSS iniziamo ponendo il gradiente a 0:

$$\nabla_{\mathbf{F}} RSS(\mathbf{F}) = \mathbf{D} = \nabla_{\mathbf{F}} \left[(\mathbf{A} - \mathbf{F} \mathbf{F}), (\mathbf{A} - \mathbf{F} \mathbf{F}) \right]$$

Eseguiamo il prodotto:

=
$$\nabla_{\underline{p}} \left[\underline{Y}^{T} \underline{Y} - \underline{Y}^{T} \underline{\times} \underline{B} - \underline{P}^{T} \underline{\times}^{T} \underline{Y} + \underline{B}^{T} \underline{\times}^{T} \underline{\times} \underline{F} \right]$$

II thas posto odel I

olim:

1×n n×(p+1) (p+1)×1

1×(p+1) /

1×1

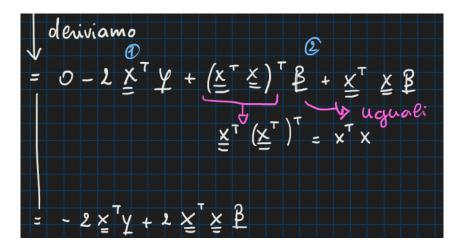
Spras posto odello scalare è

lo scalore stesso

quindi II = I

Sommiamo i due termini rosa:

Svolgiamo il gradiente:



Ricordando che quest'ultima espressione va posta a 0 ricaviamo il valore $\hat{\beta}$ che minimizza l'RSS:

$$\hat{\beta}_{i,j} = (x^T x)^{-1} x^T y \rightarrow LSE$$
stimotore che useremo
least squar

Invertibilità

Nella formula che abbiamo ottenuto c'è l'inversione di una matrice ma non tutte le matrici sono invertibili (se il rango è 0 non si può).

Risulta che la matrice è invertibile se la matrice di design è a rango pieno (full rank).

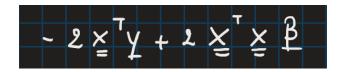


La matrice di design è full rank se n>p, cioè il numero di punti del dataset deve essere maggiore del numero di regressori.

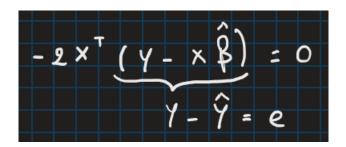
Questa affermazione può essere interpretata come il fatto che ci serve avere abbastanza punti per permettere la stima dei parametri che vogliamo stimare.

Ortogonalità di errore e span generato dalle x

Ripartiamo dall'espressione:



Mettiamo in evidenza -2 e specializziamo l'equazione per $\hat{\beta}$.



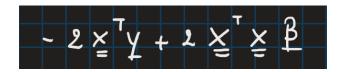


Il fatto che il risultato dell'equazione sia 0 ci dice che l'errore è ortogonale allo SPAN generato dai vettori x (visto che moltiplicato per la trasposta di x fa 0).

Il sottospazio generato dai vettori x si chiama SPAN. L'errore è ortogonale a quello spazio, dove? nel punto che minimizza l'RSS. Questa è una proprietà interessante.

Seconda parte, matrice Hessiana

Ripartiamo dall'espressione:



Abbiamo detto che a noi serve lo stimatore di β che minimizza l'RSS, per ottenere questo non basta porre il gradiente pari a 0 ma bisogna anche fare considerazioni sulla Matrice Hessiana.

Valutiamo la Matrice Hessiana nel punto $\hat{\beta}$.

Nello specifico la Matrice Hessiana deve essere definita positiva, quindi visto che è simmetrica basta che gli autovalori siano tutti positivi.

Hatrice Hessiana
$$H_{RSS}(\hat{\beta}) = \begin{bmatrix} \frac{\partial^2 RSS(\beta)}{\partial \beta \cdot \partial \beta_1} & \frac{\partial^2 RSS(\beta)}{\partial \beta \cdot \partial \beta_2} \\ \frac{\partial^2 RSS(\beta)}{\partial \beta \partial \beta^2} \end{bmatrix}$$

Se effettus tutte be observate

 $H_{RSS} = 2 \times 7 \times 6$

matrice oli varianza co-varianza
quinoli è definita positiva

 $\Rightarrow \hat{\beta}_{LS}$ minimizzatore

La matrice $H_{RSS}=2X^TX$ è una somma di quadrati, è quindi proporzionale alla matrice di varianza-covarianza dei dati x. Si può quindi dimostrare che questa matrice è una stima della matrice di var-covar e sapendo che quest'ultima ha per definizione tutti i valori positivi questo varrà anche per la matrice che stiamo cercando.

Possiamo quindi confermare che il nostro \hat{eta}_{LS} è un punto di minimo.

Ultime considerazioni

Avendo confermato la validità di \hat{eta}_{LS} possiamo sostituire l'espressione dello stimatore nell'espressione.

$$\frac{\hat{Y}}{\hat{Y}} = \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} \hat{Y}_{LS} = \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times}^T Y$$

$$= \underbrace{\times} (\underbrace{\times} (\underbrace{\times}^T \underbrace{\times})^T \underbrace{\times} (\underbrace{\times}^T \underbrace{\times}$$

La matrice H, Hat Matrix, è una matrice che prende i dati y e li trasforma in \hat{y} , il modo in cui effettua questa azione è proiettando la y nello SPAN generato dalle x.

Constatiamo che la varianza del nostro stimatore dipende da σ^2 , questo ha senso perché la varianza dell'errore implica la dispersione più o meno grande delle mie stime.

Hatrice di varianta covarianta di
$$\hat{B}_{L_s}$$
 grazie al Var $(\hat{\beta}_{L_s}) = (x^7 \times)^{-1} G^2$

Hi Gouss- Harkov $(\hat{\beta}_{L_s}) = 0 < -> E[\hat{\beta}_{L_s}] = 2$
 $\hat{\beta}_{L_s}$ Best Linear Unbiesed Est. (BLUE)

In virtù del teorema di Gauss-Markov ha media che coincide con il valore vero, quindi bias nullo. Questo è il miglior stimatore unbiased tra gli stimatori lineari, BLUE. La stima di $\,\sigma^2$ è un problema affrontato nel continuo di 2 - Linear Regression B.