

Generalized Linear Models

Vers. 1.1.2

Gianluca Mastrantonio

gianluca.mastrantonio@polito.it

"AN INTRODUCTION TO GLM (GENERALIZED LINEAR MODELS)"

ANNETTE J. DOBSON
(II EDIZIONE)

I modelli lineari generalizzati I

Cerchiamo di estendere il modello lineare per poter considerare dati che provengano da distribuzioni diverse. Per fare questo consideriamo una famiglia di distribuzioni, chiamata “esponenziale”.

La densità di una variabile aleatoria y_i si dice appartenere alla famiglia esponenziale se può essere scritta come

$$f(y_i; \theta_i, \phi_i) = \exp \left(\frac{y_i \theta_i - b(\theta_i)}{a(\phi_i)} + c(y_i, \phi_i) \right) \quad (1)$$

dove

- θ_i è chiamato parametro *naturale* e
- ϕ_i è il parametro di *dispersione*

- $a(\cdot), b(\cdot), c(\cdot)$
funzioni generiche
- θ_i, ϕ_i sono dipendenti
dalla distribuzione considerata

I modelli lineari generalizzati II

Per le distribuzioni della famiglia esponenziale, media e varianza si calcolano facilmente

Assumiamo che sia possibile cambiare l'ordine di integrazione e derivazione nella seguente formula

$$\frac{\partial}{\partial \theta_i} \int \underbrace{f(y_i; \theta_i, \phi_i) dy_i}_{\text{1}} = \int \frac{\partial f(y_i; \theta_i, \phi_i)}{\partial \theta_i} dy_i \quad \left. \vphantom{\int} \right] \text{ASSUNZIONE}$$

che si può dimostrare essere possibile per la famiglia esponenziale, allora

$$\int \frac{\partial f(y_i; \theta_i, \phi_i)}{\partial \theta_i} dy_i = 0$$

Se lo scriviamo in forma estesa per la famiglia esponenziale abbiamo

$$\int \frac{\partial f(y_i; \theta_i, \phi_i)}{\partial \theta_i} dy_i = \int \overbrace{\frac{y_i - b'(\theta_i)}{a(\phi_i)} f(y_i; \theta_i, \phi_i)}^{\text{DERIVATA (1)}} dy_i =$$

$$\frac{1}{a(\phi_i)} \underbrace{\int y_i f(y_i; \theta_i, \phi_i) dy_i}_{E[y_i] !!} - \frac{b'(\theta_i)}{a(\phi_i)} \underbrace{\int f(y_i; \theta_i, \phi_i) dy_i}_{= 1} = \frac{E(y_i)}{\cancel{a(\phi_i)}} - \frac{b'(\theta_i)}{\cancel{a(\phi_i)}} = 0$$

I modelli lineari generalizzati III

Quindi

$$E(y_i) = \mu_i = b'(\theta_i) \quad (2)$$

che ci dà la formula per la media.

Assumendo di poter cambiare l'ordine di integrazione e derivazione anche per le derivate seconde

$$\frac{\partial^2}{\partial \theta_i^2} \int \overbrace{f(y_i; \theta_i, \phi_i)}^1 dy_i = \int \frac{\partial^2 f(y_i; \theta_i, \phi_i)}{\partial \theta_i^2} dy_i \quad \Bigg] \text{ ASSUNZIONE}$$

che è possibile per la famiglia esponenziale, abbiamo che

$$\int \frac{\partial^2 f(y_i; \theta_i, \phi_i)}{\partial \theta_i^2} dy_i = 0$$

I modelli lineari generalizzati IV

e possiamo scrivere

$$\begin{aligned} \int \frac{\partial^2 f(y_i; \theta_i, \phi_i)}{\partial \theta_i^2} dy_i &= \int \frac{-b''(\theta_i)}{a(\phi_i)} f(y_i; \theta_i, \phi_i) + \frac{(y_i - b'(\theta_i))^2}{a(\phi_i)^2} f(y_i; \theta_i, \phi_i) dy_i = \\ &= -\frac{b''(\theta_i)}{a(\phi_i)} \int \underbrace{f(y_i; \theta_i, \phi_i)}_{=1} dy_i + \frac{1}{a(\phi_i)^2} \int (y_i - \overbrace{b'(\theta_i)}^{\mu_i})^2 f(y_i; \theta_i, \phi_i) dy_i = \\ &= -\frac{b''(\theta_i)}{\cancel{a(\phi_i)}} + \frac{\text{Var}(y_i)}{\cancel{a(\phi_i)^2}} = 0 \end{aligned}$$

$\downarrow \mathbb{E}[Y_i] = b'(\theta_i) !!$
 μ_i

che ci porta ad avere

$$\boxed{\text{Var}(y_i) = b''(\theta_i) a(\phi_i)} \quad (3)$$

vediamo adesso come possiamo passare dai modelli lineari a quello generalizzato.

I modelli lineari generalizzati V

Abbiamo sempre scritto il modello lineare come

$$\left. \begin{array}{l} y_i = \sum_{j=1}^p \beta_j x_{ij} + \epsilon_i \\ \epsilon \sim N(0, \sigma^2) \end{array} \right\} \mu_i$$

x_{ij} = j -esima covariata associata all' i -esima osservazione

con tutte le varie ipotesi. Ma un modo equivalente per scriverlo è il seguente

$$\left. \begin{array}{l} y_i \sim N(\mu_i, \sigma^2) \\ \mu_i = \sum_{j=1}^p \beta_j x_{ij} \end{array} \right\} \text{RAPPRESENTAZIONE GERARCHICA DEL MODELLO LINEARE}$$

dove è chiaro che stiamo modellizzando il parametro media della normale.

I modelli lineari generalizzati VI

I GLM sono una generalizzazione del modello lineare, e si possono scrivere come

$$g: \mathcal{D} \rightarrow \mathbb{R}$$
$$y_i \sim \overbrace{H(\theta_i, \phi_i)}^{\text{FAMIGLIA ESPONENZIALE}}$$
$$\underbrace{g(\mu_i)}_{!!!} = \eta_i = \sum_{j=1}^p \beta_j x_{ij}$$

- H è una distribuzione membro della famiglia esponenziale;
- $f(y_i; \theta_i, \phi_i) = \exp \left(\frac{y_i \theta_i - b(\theta_i)}{a(\phi_i)} + c(y_i, \phi_i) \right)$
- $E(y_i) = \mu_i = b'(\theta_i)$;
- $\text{Var}(y_i) = b''(\theta_i) a(\phi_i)$;
- $\mu_i \in \mathcal{D}$, dove \mathcal{D} può essere \mathbb{R}^+ , oppure $[0, 1]$, etc...;
- $g: \mathcal{D} \rightarrow \mathbb{R}$ è la funzione link;

I modelli lineari generalizzati VII

Tra tutte le funzioni link, quella che soddisfa $\theta_i = g(\mu_i)$ è chiamata **funzione link**

canonica

USEREMO SEMPRE QUESTA !!

vediamo alcuni esempi

Distribuzione normale $y \sim N(\alpha_i, \sigma^2)$ La densità è

$$f(y_i; \alpha_i, \sigma^2) = (2\pi\sigma^2) \exp\left(-\frac{(y_i - \alpha_i)^2}{2\sigma^2}\right) \quad (*)$$

$$= e^{-\log(2\pi\sigma^2) \cdot \frac{1}{2}} = \exp\left(-\frac{1}{2} \log(2\pi\sigma^2)\right)$$

$$(*) \quad \frac{-y_i^2 + 2y_i\alpha_i - \alpha_i^2}{2\sigma^2} = \exp\left(\underbrace{\frac{y_i\alpha_i - \frac{1}{2}\alpha_i^2}{\sigma^2}}_{a(\phi)} - \underbrace{\frac{1}{2} \log(2\pi\sigma^2)}_{c(y_i, \phi)} - \underbrace{\frac{y_i^2}{2\sigma^2}}_{b(\theta_i)}\right)$$

Con

- $\theta_i = \alpha_i$ e $\phi = \sigma^2$;
- $b(\theta_i) = \frac{\theta_i^2}{2}$ $a(\phi_i) = \phi_i$ $c(y_i; \phi_i) = -\frac{1}{2} \log(2\pi\phi_i) - \frac{y_i^2}{2\phi_i}$;
- $E(y_i) = \theta_i = \alpha_i$ $\text{Var}(y_i) = \phi_i = \sigma^2$

$$E[y_i] = b'(\theta) = \theta_i \quad \text{Var}[y_i] = + b''(\theta) a(\phi) = 1 \cdot \sigma^2$$

I modelli lineari generalizzati VIII

- la funzione link canonica è $g(\mu_i) = \mu_i$

Distribuzione Poisson $y_i \sim \text{Pois}(\lambda_i)$ La densità è

$$\begin{aligned} f(y_i; \lambda) &= \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!} & \log(\lambda_i) = \theta \\ &= \exp \left(\underbrace{y_i \log(\lambda_i)}_1 - \underbrace{\lambda_i - \log(y_i!)}_1 \right) \\ &\propto \phi \longrightarrow 1 \end{aligned}$$

- $\theta = \log(\lambda_i)$ e $\phi = \emptyset$;
- $b(\theta_i) = \exp(\theta_i)$ $a(\phi_i) = 1$ $c(y_i; \phi_i) = -\log(y_i!)$;
- $E(y_i) = \exp(\theta_i) = \lambda_i$ $\text{Var}(y_i) = \exp(\theta_i) = \lambda_i$
- la funzione link canonica è $g(\mu_i) = \log(\mu_i)$

I modelli lineari generalizzati IX

Distribuzione Binomiale $y_i \sim \text{Bin}(n_i, \pi_i)$ (n_i si considera noto e non parametro)

$$f(y_i; n_i, \pi_i) = \binom{n_i}{y_i} \cdot \pi_i^{y_i} (1 - \pi_i)^{n_i - y_i}$$

$$= \exp \left(\underbrace{\frac{y_i \log(\pi_i / (1 - \pi_i)) + n_i \log(1 - \pi_i)}{1}}_{\log_{\pi_i}(1 - \pi_i)} - \underbrace{\log \binom{n_i}{y_i}}_{\text{INVERSIÓN}} \right)$$

$\pi_i = f(\theta)$ INVERSIÓN

- $\theta = \log \left(\frac{\pi_i}{1 - \pi_i} \right)$ e $\phi = \emptyset$;

- $b(\theta_i) = n_i \log(1 + \exp(\theta_i))$ $a(\phi_i) = 1$ $c(y_i; \phi_i) = \log \binom{n_i}{y_i}$;

- $E(y_i) = n_i \frac{\exp(\theta_i)}{1 + \exp(\theta_i)} = n_i \pi_i$ $\text{Var}(y_i) = \frac{\exp(\theta_i) n_i}{(1 + \exp(\theta_i))^2} = n_i \pi_i (1 - \pi_i)$

- la funzione link canonica è $g(\mu_i) = \log \left(\frac{\mu_i}{n_i - \mu_i} \right) = \log \left(\frac{\pi_i}{1 - \pi_i} \right)$

LOGIT !!!

Stima dei parametri I

Per ognuno di queste distribuzioni abbiamo un valore medio e siamo interessati a modellizzarlo con una funzione lineare

$$g(\mu_i) = \eta_i = \beta_1 + \sum_{j=2}^p \beta_j x_{ij} = \mathbf{X}_i \boldsymbol{\beta}$$

e come nel modello lineare, vogliamo trovare le stime dei parametri trovando il massimo della verosimiglianza (o log verosimiglianza), che può essere trovato imponendo la derivata uguale a 0.

Indichiamo con $L(\vec{\boldsymbol{\beta}}; \mathbf{y})$ la log-verosimiglianza di $\vec{\boldsymbol{\beta}}$

$$L(\boldsymbol{\beta}; \mathbf{y}) = \sum_{i=1}^n \frac{y_i \theta_i - b(\theta_i)}{a(\phi_i)} + \sum_{i=1}^n c(y_i, \phi_i)$$

e con \mathbf{U} il vettore di derivate rispetto a $\boldsymbol{\beta}$, chiamato anche **statistica score**, con elementi:

$$U_j = \frac{\partial L(\boldsymbol{\beta}; \mathbf{y})}{\partial \beta_j}$$

Stima dei parametri II

Per calcolare le derivate utilizziamo la chain rule:

$$U_j = \frac{\partial L(\beta; y_i)}{\partial \beta_j} = \underbrace{\frac{\partial L(\beta; y_i)}{\partial \theta_i}}_{(1)} \underbrace{\frac{\partial \theta_i}{\partial \mu_i}}_{(2)} \underbrace{\frac{\partial \mu_i}{\partial \eta_i}}_{(3)} \underbrace{\frac{\partial \eta_i}{\partial \beta_j}}_{(4)}$$

dove

$$(1) \quad \frac{\partial L(\beta; y_i)}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{a(\phi_i)} \xrightarrow{b'(\theta_i) = \mu_i} \frac{y_i - \mu_i}{a(\phi_i)}$$

$$(2)^{-1} \quad \frac{\partial \mu_i}{\partial \theta_i} = b''(\theta_i) = \frac{\text{Var}(y_i)}{a(\phi_i)} \rightarrow \mu_i = E[y_i] = b'(\theta)$$

$$\eta_i = \sum_{j=1}^p \beta_j x_{ij} \rightarrow (4) \quad \frac{\partial \eta_i}{\partial \beta_j} = x_{ij}$$

$$\frac{\partial \mu_i}{\partial \theta_i} = b''(\theta)$$

e sfruttando $\text{Var}(y_i) = b'' \cdot a$
 $\rightarrow b'' = \frac{\text{Var}(y_i)}{a(\phi_i)}$

Abbiamo quindi

$$U_j = \sum_{i=1}^n \frac{y_i - \mu_i}{\cancel{a(\phi_i)} \text{Var}(y_i)} x_{ij} \frac{\partial \mu_i}{\partial \eta_i} = \sum_{i=1}^n \frac{(y_i - \mu_i) x_{ij}}{\text{Var}(y_i)} \frac{\partial \mu_i}{\partial \eta_i} \quad (4)$$

$$E[U_j] = \dots$$

sfruttando
 $E[y_i] = \mu_i$

DIPENDENTE DA
 MODELLO E DA $g(\cdot)$!!!

Stima dei parametri III

dove $\frac{\partial \mu_i}{\partial \eta_i}$ dipende dal particolare modello che stiamo usando.

La soluzione $U = 0$ rispetto a β generalmente non si riesce a trovare, e algoritmi numerici vengono usati (Newton-Raphson o Fisher Scoring) **ITERATIVI !!!**

Ipotizziamo di aver trovato la soluzione $\hat{\beta}$, abbiamo adesso bisogno di trovare la sua distribuzione per poter fare inferenza: test, intervalli di confidenza etc.. Visto che $\hat{\beta}$ non si può trovare in forma chiusa, non sappiamo neanche la sua distribuzione, e abbiamo bisogno di un qualche metodo per trovarla (almeno approssimativamente)

A questo scopo, possiamo utilizzare la statistica score U . Indichiamo con **V** e **D** due matrici diagonali aventi come elemento i-esimo $\text{Var}(y_i)$ e $\partial \mu_i / \partial \eta_i$, rispettivamente.

Allora

$$U = X^T D V^{-1} (y - \mu)$$

RISULTATO COMPATTO
D. 4

La statistica U è una variabile aleatoria con media 0 ($E(y_i) = \mu_i$)

Stima dei parametri IV

I METODO

Definiamo \mathbf{J} come la matrice di varianza e covarianze di \mathbf{U} , e chiamiamola **matrice d'informazione**, l'elemento (j, k) di \mathbf{J} si calcola come

$$\text{Cov}(U_j, U_k) = E(U_j U_k) = E \left(\sum_{i=1}^n \frac{(y_i - \mu_i) x_{ij}}{\text{Var}(y_i)} \frac{\partial \mu_i}{\partial \eta_i} \sum_{h=1}^n \frac{(y_h - \mu_h) x_{hk}}{\text{Var}(y_h)} \frac{\partial \mu_h}{\partial \eta_h} \right)$$

Dato che y_i è indipendente da y_j abbiamo che $E((y_i - \mu_i)(y_h - \mu_h)) = 0$, quindi

$$\textcircled{*} \quad J_{jk} = E(U_j U_k) = E \left(\sum_{i=1}^n \frac{(y_i - \mu_i)^2 x_{ij} x_{ik}}{\text{Var}(y_i)^2} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2 \right) = \sum_{i=1}^n \frac{x_{ij} x_{ik}}{\text{Var}(y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2$$

Possiamo allora definire

$$\mathbf{J} = \mathbf{X}^T \mathbf{W} \mathbf{X}$$

⑤

con $\underline{W_{ii} = \frac{1}{\text{Var}(y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2}$

←
RIMANGONO SOLO
I TERMINI IN CUI
 $i = k$ $\textcircled{*}$

Stima dei parametri V

II METODO

Abbiamo anche un diverso modo di calcolare la matrice di informazione. Calcoliamo prima il valore di $E(U_j)$

$$E(U_j) = \int U_j f(\mathbf{y}; \boldsymbol{\beta}, \phi_i) dy = \int \frac{\partial \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_j} f(\mathbf{y}; \boldsymbol{\beta}, \phi_i) dy = 0.$$

$U_j = \frac{\partial \log f}{\partial \beta_j}$

$\partial / \partial \beta_j$

Se differenziamo rispetto a β_k , e scambiamo l'ordine di integrale con derivata, otteniamo

$$0 = \int \frac{\partial^2 \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_j \partial \beta_k} f(\mathbf{y}; \boldsymbol{\beta}, \phi_i) dy + \int \frac{\partial \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_j} \frac{\partial f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_k} dy =$$

usando la derivata del logaritmo, otteniamo

$$\int \frac{\partial^2 \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_j \partial \beta_k} f(\mathbf{y}; \boldsymbol{\beta}, \phi_i) dy + \frac{\partial f}{\partial \beta_k} = \frac{\partial \log f}{\partial \beta_k} \cdot f$$

$$\int \frac{\partial \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_j} \frac{\partial \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_k} f(\mathbf{y}; \boldsymbol{\beta}, \phi_i) dy =$$

Stima dei parametri VI

Si vede facilmente che entrambi gli integrali sono valori attesi, e precisamente

$$\int \frac{\partial^2 \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_j \partial \beta_k} f(\mathbf{y}; \boldsymbol{\beta}, \phi_i) d\mathbf{y} = E \left(\frac{\partial U_j}{\partial \beta_k} \right) \quad \bullet$$

e

$$\int \underbrace{\frac{\partial \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_j}}_{U_j} \underbrace{\frac{\partial \log f(\mathbf{y}; \boldsymbol{\beta}, \phi_i)}{\partial \beta_k}}_{U_k} f(\mathbf{y}; \boldsymbol{\beta}, \phi_i) d\mathbf{y} = E(U_j U_k) \quad \dots$$

Abbiamo allora che

\mathbf{U}_j

\mathbf{U}_k

$$J_{jk} = E(U_j U_k) = -E \left(\frac{\partial U_j}{\partial \beta_k} \right)$$

$$\odot + \odot = \bigcirc$$

$$\Rightarrow \odot = -\odot$$

Oppure, in forma matriciale

$$\boxed{\mathbf{J} = E(\mathbf{U}\mathbf{U}^T) = -E(\mathbf{U}')} \quad \textcircled{6}$$

Le matrici \mathbf{U} e \mathbf{J} sono funzioni di $\boldsymbol{\beta}$, e quindi le indichiamo, in generale, $\mathbf{U}(\mathbf{b})$ se abbiamo bisogno di sapere dove sono calcolate.

Stima dei parametri VII

Queste matrici ci serviranno per fare i test di ipotesi su β . Ma per fare questo abbiamo bisogno della distribuzione di U

Calcoliamo l'approssimazione di Taylor per la log verosimiglianza, calcolata in $\hat{\beta}$

$$\begin{aligned}
 & \bullet L(\beta; y) \approx L(\hat{\beta}; y) + (\beta - \hat{\beta})^T \frac{\partial L(\hat{\beta}; y)}{\partial \beta} + \frac{1}{2} (\beta - \hat{\beta})^T \frac{\partial^2 L(\hat{\beta}; y)}{\partial \beta \partial \beta^T} (\beta - \hat{\beta}) = \\
 & \text{DEFINIZIONE} \quad \text{DI } U \quad L(\hat{\beta}; y) + (\beta - \hat{\beta})^T U(\hat{\beta}) + \frac{1}{2} (\beta - \hat{\beta})^T U(\hat{\beta})' (\beta - \hat{\beta})
 \end{aligned}$$

Approssimiamo U' con $\underbrace{E(U')}_{=0} = -J$ e ricordiamo che $U(\hat{\beta}) = 0$, allora

$$\underbrace{U(\hat{\beta})}_{\text{non è } \hat{\beta} \text{ MLE}} = 0 !!! \quad L(\beta; y) \approx L(\hat{\beta}; y) - \frac{1}{2} (\beta - \hat{\beta})^T J(\hat{\beta}) (\beta - \hat{\beta})$$

Calcoliamo la derivata rispetto a β ottenendo

$$\frac{\partial L(\beta; y)}{\partial \beta} = 0 \quad U(\beta) = -(\beta - \hat{\beta})^T J(\hat{\beta}) = -J(\hat{\beta})(\beta - \hat{\beta})$$

Stima dei parametri VIII

da cui abbiamo che

$$\hat{\beta} = \beta + J^{-1}U \quad !!!$$

$$oss: J^{-1} \cdot U = J^{-\frac{1}{2}} \cdot (J^{-\frac{1}{2}})^T \cdot U$$

Vogliamo trovare la distribuzione di U con il teorema del limite centrale. Ricordiamo che l'elemento j -esimo della statistica score è

$$U_j = \sum_{i=1}^n \frac{(y_i - \mu_i)x_{ij}}{\text{Var}(y_i)} \frac{\partial \mu_i}{\partial \eta_i} \quad ||$$

e

$$J^{-1/2}U = \sum_{i=1}^n a_i \left(\frac{y_i - \mu_i}{\text{Var}(y_i)} \right) \quad \rightarrow \text{Var} = 1$$

è una combinazione lineare di $(y_i - \mu_i)$, con a_i specifici e dove $J^{-1/2}(J^{-1/2})^T = J^{-1}$. Il teorema classico del limite centrale non si riesce ad applicare nella sua versione base, ma possiamo usare la **condizione di Lindeberg** che ci assicura che (sotto opportune

Stima dei parametri IX

condizioni) si può applicare il TLC anche in questo caso, dato che $a_i(y_i - \mu_i)$ è una variabile aleatoria con media 0 e varianza a_i^2 , con $\sum_{i=1}^n a_i^2 = 1$ perchè $\text{Var}(\mathbf{J}^{-1/2} \mathbf{U}) = \mathbf{I}$

!!!

Abbiamo allora che (approssimativamente, per $n \Rightarrow \infty$)

$$\textcircled{*} \rightarrow \boxed{\hat{\beta} \sim N(\beta, \mathbf{J}^{-1})} \textcircled{7}$$

$$\hat{\beta} = \beta + \mathbf{J}^{-\frac{1}{2}} \underbrace{(\mathbf{J}^{-\frac{1}{2}}) \mathbf{U}}_{\substack{N(\vec{0}, \mathbf{I}) \\ \text{quindi } \textcircled{*}}}$$

Notate come la stima di β è più precisa se aumenta la curvatura della funzione di verosimiglianza ($-\mathbf{J}$)

$$\checkmark \quad \mathbf{C} = (0, 0, 0 \dots 1, 0, 0, 0)$$

$$\mathbf{C} \hat{\beta} = \hat{\beta}_J$$

$$\mathbf{C} \beta = \beta_J$$

$$\mathbf{C} \mathbf{J}^{-1} = \mathbf{J}_{JJ}^{-1}$$

Stima dei parametri X

Come nel caso lineare, possiamo usare la distribuzione di $\hat{\beta}$ per fare test d'ipotesi, intervalli di confidenza, etc ...

Per esempio, potremmo testare

$$H_0 : \beta_j = b_j \quad H_1 : \beta_j \neq b_j$$

con la statistica

$$z = (\hat{\beta}_j - b_j) \sqrt{J_{jj}} \sim N(0, 1) \quad \bullet$$

se si assume nota la varianza, oppure

$$z = (\hat{\beta}_j - b_j) \sqrt{\hat{J}_{jj}} \sim N(0, 1) \quad \bullet$$

se viene stimata, dove $\hat{\mathbf{J}} = \mathbf{X}^T \hat{\mathbf{W}} \mathbf{X}$ e \hat{W}_{ii} è una stima di $W_{ii} = \frac{1}{\text{Var}(y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2$.

Nel secondo caso dovremmo usare una T di student, ma la distribuzione di $\hat{\beta}$ è approssimativamente normale, come la T .

Stima dei parametri XI

Possiamo testare un'ipotesi **sull'intero modello**

$$\underline{H_0 : \beta = b} \quad \underline{H_1 : \beta \neq b}$$

possono essere testate utilizzando la statistica di Wald

$$\underline{(\beta - b)^T J (\beta - b) \sim \chi_p^2}$$

sfruttando il fatto
che $(\beta - b) J^{-\frac{1}{2}} \sim N(\vec{0}, I)$

o nella forma

$$\underline{U^T J^{-1} U \sim \chi_p^2} \quad \leftarrow$$

visto che $U \sim N(0, J)$ $\Rightarrow U^T J^{-\frac{1}{2}} \sim N(\vec{0}, I)$

Stima dei parametri XII

Attenzione: Sebbene abbiamo derivato solo rispetto a β la log-verosimiglianza, delle volte dobbiamo stimare anche il parametro di dispersione (per esempio σ^2 nella regressione o la dispersione della binomiale negativa). Il parametro da trovare è ϕ .

I modelli vengono scelti in modo tale che la matrice di informazione ha elementi che connettono i β e ϕ asintoticamente pari a 0. Rendendo i parametri indipendenti e possono essere stimati indipendentemente.

Stima dei parametri XIII

La Devianza:

La verosimiglianza può essere usata per fare dei test sul modello. Una possibile idea è di confrontare il modello in esame con il modello saturato (o saturo), i.e. il modello con il maggior numero di parametri possibile (1 per ogni osservazione univoca).

Indichiamo con $L(\mathbf{y}; \hat{\beta}_{\max})$ la log-verosimiglianza per il modello saturo e come sempre $L(\mathbf{y}; \hat{\beta})$ quella del modello in esame. Siamo interessati al rapporto

$$\text{RAPPORTO} \\ \text{FRA} \\ \text{VEROSIMIGLIANZE} \quad \lambda = \frac{\exp(L(\hat{\beta}_{\max}; \mathbf{y}))}{\exp(L(\hat{\beta}; \mathbf{y}))}$$

o più precisamente al suo logaritmo moltiplicato per due, che chiamiamo **Devianza**

- $D = 2 \log(\lambda) = 2 \left(L(\hat{\beta}_{\max}; \mathbf{y}) - L(\hat{\beta}; \mathbf{y}) \right)$

Stima dei parametri XIV

Ricordiamo che l'approssimazione di Taylor della log-verosimiglianza ci dice che

$$\underline{L(\beta; \mathbf{y}) = L(\hat{\beta}; \mathbf{y}) + (\beta - \hat{\beta})^T \cancel{\mathbf{U}(\hat{\beta})} - \frac{1}{2}(\beta - \hat{\beta})^T \mathbf{J}(\hat{\beta})(\beta - \hat{\beta})}$$

(approssimativamente) Se $\hat{\beta}$ è lo stimatore di massima verosimiglianza, allora possiamo scrivere

$$L(\beta; \mathbf{y}) = L(\hat{\beta}; \mathbf{y}) - \frac{1}{2}(\beta - \hat{\beta})^T \mathbf{J}(\hat{\beta})(\beta - \hat{\beta})$$

e dato che $\hat{\beta} \sim N(\beta, \mathbf{J}^{-1})$ allora

$$2 \left(L(\hat{\beta}; \mathbf{y}) - L(\beta; \mathbf{y}) \right) = (\beta - \hat{\beta})^T \mathbf{J}(\hat{\beta})(\beta - \hat{\beta}) \sim \chi_p^2$$

Il risultato vale per ogni β di lunghezza p .

Stima dei parametri XV

Aggiungiamo e eliminiamo il valore “vero” dei β , nel modello saturo e nel modello sotto esame, avendo

$$\bullet \quad D = \underbrace{2 \left(L(\hat{\beta}_{max}; \mathbf{y}) - L(\hat{\beta}; \mathbf{y}) \right)}_{(*)} = \underbrace{2 \left(L(\hat{\beta}_{max}; \mathbf{y}) - L(\beta_{max}; \mathbf{y}) \right)}_{(*)} - \underbrace{2 \left(L(\hat{\beta}; \mathbf{y}) - L(\beta; \mathbf{y}) \right)}_{(**)} + \underbrace{2 \left(L(\beta_{max}; \mathbf{y}) - L(\beta; \mathbf{y}) \right)}_{\star \text{ COSTANTI !!!}}$$

$$(*) \bullet \quad \underline{2 \left(L(\hat{\beta}_{max}; \mathbf{y}) - L(\beta_{max}; \mathbf{y}) \right) \sim \chi_n^2, \text{ dove } n \text{ è la lunghezza di } \beta_{max};}$$

$$(**) \bullet \quad \underline{2 \left(L(\hat{\beta}; \mathbf{y}) - L(\beta; \mathbf{y}) \right) \sim \chi_p^2, \text{ dove } p \text{ è la lunghezza di } \beta;}$$

\star \bullet $\nu = 2 \left(L(\beta_{max}; \mathbf{y}) - L(\beta; \mathbf{y}) \right) \geq 0$ è una costante che è uguale a 0, se i due modelli spiegano i dati nello stesso modo. Deve essere maggiore di zero perchè il saturo contiene il modello da testare (è un possibile caso particolare)

Stima dei parametri XVI

La differenza di due chi quadrato non è nota in generale, ma usando il teorema di Wilks, si può vedere che è

WILK'S THEOREM

$$\boxed{D \sim \chi^2_{n-p, \nu}} \cong D = X + \nu$$

con $X \sim \chi^2_{n-p}$

Relazione simile vale per qualsiasi devianza, non solo calcolata con il saturo

La devianza può essere usata per confrontare modelli *annidati*. nello specifico ipotizziamo di voler testare

$$\boxed{H_0 : \beta = \beta_0 \quad H_1 : \beta = \beta_1}$$

dove β_0 ha q elementi (modello M_0) e β_1 ha p elementi (modello M_1), con $q < p < n$.
Possiamo calcolare

piccolo

grande

$$\Delta D = \underbrace{D_0 - D_1}_{\geq D_1} = 2 \left(\cancel{L(\hat{\beta}_{\max}; \mathbf{y})} - L(\hat{\beta}_0; \mathbf{y}) \right) - 2 \left(\cancel{L(\hat{\beta}_{\max}; \mathbf{y})} - L(\hat{\beta}_1; \mathbf{y}) \right) =$$

Stima dei parametri XVII

e abbiamo quindi che

$$\Delta D = 2 \left(L(\hat{\beta}_1; \mathbf{y}) - L(\hat{\beta}_0; \mathbf{y}) \right) \sim \chi_{p-q, v_0 - v_1}^2$$

dove

- $\nu_0 = 2 \left(L(\beta_{max}; \mathbf{y}) - L(\beta_0; \mathbf{y}) \right);$
- $\nu_1 = 2 \left(L(\beta_{max}; \mathbf{y}) - L(\beta_1; \mathbf{y}) \right);$

sotto H_0 : $L(\beta_0; \mathbf{y}) = L(\beta_1; \mathbf{y})$
che i dati sono
spiegati ugualmente
da M_0 e M_1

Se entrambi i modelli sono equivalenti in termine di spiegazione dei dati, abbiamo

$$\nu_0 = \nu_1.$$

e $\Delta D \sim \chi_{p-q}^2$.

Se M_0 spiega in maniera peggiore i dati, rispetto a M_1 , abbiamo che

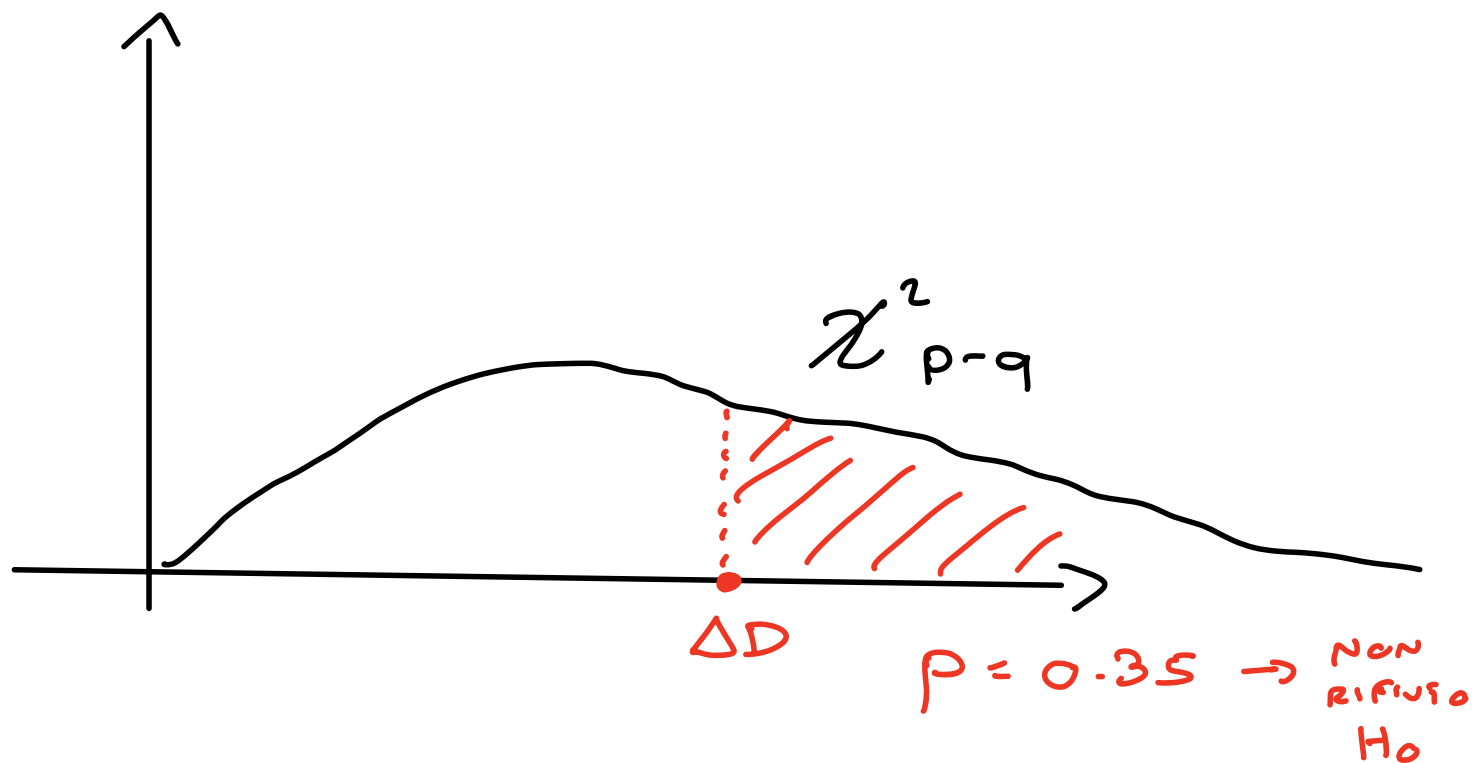
$$\nu_0 > \nu_1$$

Stima dei parametri XVIII

$$\text{e } \Delta D \sim \chi^2_{p-q, \nu_0 - \nu_1}.$$

Possiamo confrontare ΔD con un χ^2_{p-q} .

Se il valore ΔD si trova nella regione critica $\int_D^\infty f^*(x)dx < \alpha$, dove $f^*(x)$ è la densità di una χ^2_{p-q} , rifiutiamo H_0 , altrimenti, per il rasoio di Occam, scegliamo il modello più parsimonioso.



Stima dei parametri XIX

Prendiamo il caso normale.

Possiamo calcolare la devianza per un **modello normale**, con lo scopo di confrontare modelli nested. Nel modello saturo, il vettore dei parametri è della stessa dimensione delle osservazioni, e abbiamo che $\mathbf{y} = \mathbf{\hat{y}}$, quindi

$$\underline{L(\hat{\beta}_{max}; \mathbf{y}) = -\frac{n}{2} \log(2\pi\sigma^2)}$$

$(2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(\frac{-(y_i - \mu_i)^2}{2\sigma^2}\right)$
 $\mu_i = y_i$

mentre per il modello di interesse è

$$L(\hat{\beta}; \mathbf{y}) = -\frac{\sum_{i=1}^n (y_i - \mathbf{X}_i \hat{\beta})^2}{2\sigma^2} - \frac{n}{2} \log(2\pi\sigma^2)$$

e

$$D = 2 \left(L(\hat{\beta}_{max}; \mathbf{y}) - L(\hat{\beta}; \mathbf{y}) \right) = \frac{\sum_{i=1}^n (y_i - \mathbf{X}_i \hat{\beta})^2}{\sigma^2} \sim \chi_{n-p}^2$$

Sfortunatamente, questa non può essere usata direttamente dato che dipende da σ^2 , che non conosciamo.

Stima dei parametri XX

Indichiamo con

- $D_0 = \frac{\sum_{i=1}^n (y_i - \mathbf{X}_i \cdot \hat{\boldsymbol{\beta}}_0)^2}{\sigma^2} \sim \chi_{n-q, \nu_0}^2$
- $D_1 = \frac{\sum_{i=1}^n (y_i - \mathbf{X}_i \cdot \hat{\boldsymbol{\beta}}_1)^2}{\sigma^2} \sim \chi_{n-p, \nu_1}^2$

Assumendo che M_1 fitta bene i dati, e quindi la χ^2 ha parametro di non centralità pari a 0, possiamo utilizzare la seguente statistica

$$F = \frac{\overbrace{D_0 - D_1}^{\Delta D}}{\underbrace{p - q}} / \frac{\overbrace{D_1}^{D_1}}{n - p} \sim F_{p-q, n-p, \nu_0 - \nu_1} \quad \star$$

Quindi, se assumiamo che M_1 e M_0 spiegano allo stesso modo i dati, $\nu_0 = \nu_1$ e possiamo confrontare F con una $F_{p-q, n-p}$. Valori elevati di F suggeriscono che M_1 spieghi meglio di M_0 .

Stima dei parametri XXI

Rappresentazione grafica del test di Wald e rapporto di verosimiglianze

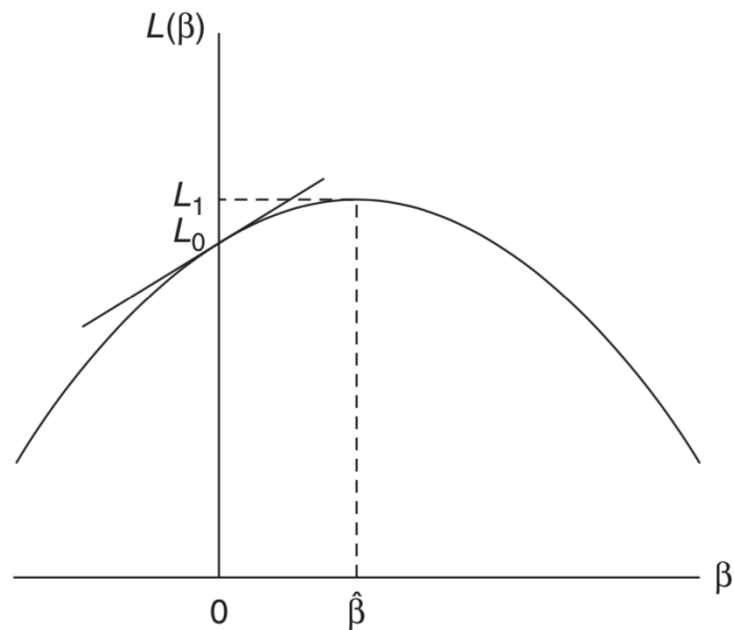


Figure 4.1 Log-likelihood function and information used in likelihood-ratio, score, and Wald tests of $H_0: \beta = 0$.

Stima dei parametri XXII

Concludiamo con dei residui

- Residui di Pearson

$$e_i = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{\text{Var}}(y_i)}}$$

- deviance residuals

$$d_i = \text{sign}(y_i - \hat{\mu}_i) \sqrt{2 \left(L_i(\hat{\beta}_{\max}; \mathbf{y}) - L_i(\hat{\beta}); \mathbf{y} \right)}$$

dove L_i è il contributo dato alla log-ver. dalla i-esima osservazione

