

# Statistical Models

Kevin Michael Frick

3 gennaio 2020

## 1 Formulario

### 1.1 Stima di valori attesi e intervalli di predizione

$$E[\hat{y}_h] = X^T \beta \quad (1)$$

$$\text{Var}[\hat{y}_h] = \sigma^2 x_h^T (X^T X)^{-1} x_h \quad (2)$$

$$s^2[\hat{y}_h] = \text{MSE} \cdot x_h^T (X^T X)^{-1} x_h \quad (3)$$

$$\epsilon_i \sim N \implies \hat{y}_h \sim N \quad (4)$$

$$\epsilon_i \sim N \implies \frac{\hat{y}_h - x_h^T \beta}{\sqrt{s^2[\hat{y}_h]}} \sim t(n-p) \quad (5)$$

$$E[y_{h(\text{new})} - \hat{y}_h] = 0 \quad (6)$$

$$\text{Var}[y_{h(\text{new})} - \hat{y}_h] = \sigma^2 (1 + x_h^T (X^T X)^{-1} x_h) \quad (7)$$

$$s^2[y_{h(\text{new})} - \hat{y}_h] = \text{MSE} \cdot (1 + x_h^T (X^T X)^{-1} x_h) \quad (8)$$

$$\frac{y_{h(\text{new})} - \hat{y}_h}{\sqrt{s^2[y_{h(\text{new})} - \hat{y}_h]}} \sim t(n-p) \quad (9)$$

### 1.2 Variabili standardizzate

$$y_{is} = \frac{y_i - \bar{y}}{\sqrt{s_y^2}} \quad (10)$$

$$x_{ks} = \frac{x_k - \bar{x}_k}{\sqrt{s_k^2}} \quad (11)$$

$$X_{RS} = X_{RC} D_R^{-1} \quad (12)$$

$$D_R = \text{diag}_k \{ \sqrt{s_k^2} \} \quad (13)$$

$$b_{0s} = 0 \quad (14)$$

$$R_{xx} = D_R^{-1} S_{xx} D_R^{-1} \quad (15)$$

$$y_s = \frac{M_0 y}{\sqrt{s_y^2}} \quad (16)$$

$$r_{yx} = \frac{D_R^{-1} s_{yx}}{\sqrt{s_y^2}} \quad (17)$$

$$b_{Rs} = \frac{D_R b_R}{\sqrt{s_y^2}} = R_{xx}^{-1} r_{yx} \quad (18)$$

$$b_{kR} = \sqrt{\frac{s_y^2}{s_k^2}} b_{ks} \quad (19)$$

### 1.3 Diagnostiche di regressione

$$h_{ii} = x_i^T (X^T X)^{-1} x_i \quad (20)$$

$$e \sim MVN(0, \sigma^2 M) \quad (21)$$

$$r_i = \frac{e_i}{\sqrt{\text{MSE} \cdot (1 - h_{ii})}} \quad (22)$$

$$d_i = y_i - \hat{y}_{i(i)} \quad (23)$$

$$t_i = \frac{d_i}{\sqrt{s^2[d_i]}} = e_i \sqrt{\frac{n-p-1}{\text{SSE}(1-h_{ii}) - e_i^2}} \quad (24)$$

$$y_i \text{ non outlier} \implies t_i \sim t(n-1-p) \quad (25)$$

$$\text{Pr}[\text{Err. t. 1 su famiglia}] \approx (1-\alpha)^n \quad (26)$$

$$\text{Pr}[\text{Err. t. 1 su famiglia (corr. Bonf.)}] = 1 - (1 - \frac{\alpha}{n})^n \approx \alpha \quad (27)$$

$$\sum_i h_{ii} = p \quad (28)$$

$$\bar{h} = \frac{p}{n} \quad (29)$$

$$n \gg 3p \implies h_{ii} \geq 2p/n \text{ (o } 3p/n \text{) valore grande} \quad (30)$$

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(i)})^2}{p \cdot \text{MSE}} \quad (31)$$

$$D_i = \frac{r_i^2 h_{ii}}{p(1-h_{ii})} \quad (32)$$

$$D_i > F(1/2; n, p) \implies y_i \text{ influente} \quad (33)$$

### 1.4 Regressori categorici

$$b_0 = \bar{y}_{reg} \quad (34)$$

$$b_k = \bar{y}_k - \bar{y}_{reg} \quad (35)$$

### 1.5 Selezione del modello

$$R_a^2 = 1 - \left( \frac{\text{SSE}}{\text{SSTOT}} \cdot \frac{n-1}{n-p} \right) = 1 - \frac{\text{MSE}}{s_y^2} \quad (36)$$

$$\text{AIC} = n \log\left(\frac{\text{SSE}}{n}\right) + 2p \quad (37)$$

$$\text{SBC} = n \log\left(\frac{\text{SSE}}{n}\right) + p \log(n) \quad (38)$$

### 1.6 Regressione logistica semplice

$$L(\beta_0, \beta_1) = \prod_i \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right)^{y_i} \left( \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right)^{1-y_i} \quad (39)$$

$$l(\beta_0, \beta_1) = \log L(\beta_0, \beta_1) \quad (40)$$

$$\hat{\pi}_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \quad (41)$$

## 1.7 Regressione logistica multipla

$$E[b] \rightarrow \beta \quad (42)$$

$$Var[b] \rightarrow -E^{-1}\left[\frac{\partial^2}{\partial\beta\partial\beta^T} \log L(\beta)\right] \quad (43)$$

$$b \sim MVN_p \quad (44)$$

$$s[b] = (-G)^{-1} \quad (45)$$

$$G = \frac{\partial^2}{\partial\beta\partial\beta^T} \log L(\beta)|_{\beta=b} \quad (46)$$

$$b_k \in_{\alpha} \beta_k \pm Z(1 - \alpha/2)\sqrt{s^2[b_k]} \quad (47)$$

$$\frac{b_k - \beta_k}{\sqrt{s^2[b_k]}} \rightarrow N(0, 1) \quad (48)$$

$$G^2 = 2 \log \frac{L(F)}{L(R)} \geq 0 \approx \chi^2(p - q)|H_0 \quad (49)$$

$$G^2 \approx N^2(0, 1) = \frac{b_k^2}{s^2[b_k]} \quad (50)$$

## 2 Ipotesi

### 2.1 Approccio di regressione all'analisi della varianza

1.  $y_{ij} = \beta_0 + \sum_k \beta_k x_{ijk} + \epsilon_{ij}$ ;
2.  $x_{ijk}$  fissi  $\implies n$  fissi;
3.  $\epsilon_{ij} \sim N(0, \sigma)$  iid.

### 2.2 Regressione logistica semplice

1.  $y_i \sim Be(\pi_i)$  iid;
2.  $\pi_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$ ;  $\pi'_i = \log\left(\frac{\pi_i}{1 - \pi_i}\right)$ ;
3.  $x_i$  fissi.

### 2.3 Regressione logistica multipla

1.  $y_i \sim Be(\pi_i)$  iid;
2.  $\pi_i = \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}$ ;  $\pi'_i = \log\left(\frac{\pi_i}{1 - \pi_i}\right)$ ;
3.  $x_i$  fissi.

**Disclaimer:** Questo documento pu contenere errori e imprecisioni che potrebbero danneggiare sistemi informatici, terminare relazioni e rapporti di lavoro, liberare le vesciche dei gatti sulla moquette e causare un conflitto termonucleare globale. Procedere con cautela.

Questo documento rilasciato sotto licenza CC-BY-SA 4.0. 