

REGRESSIONE LINEARE

Note Title

04/12/2017

$$D_n := \sup_{t \in \mathbb{R}} \left| \frac{\#\{i: F_0(x_i) \leq \bar{F}_0(t)\}}{n} - \bar{F}_0(t) \right|$$

$$\mathbb{P}(D_n \geq d) = \mathbb{P} \left(\sup_{y \in (0,1)} \left| \frac{\#\{i: U_i \leq y\}}{n} - y \right| \geq d \right)$$

dove U_1, \dots, U_n sono v.o. i.i.d. con $\mathbb{P}_{U_i} = U(0,1)$ purché $F_0: \mathbb{R} \rightarrow [0,1]$ sia continua.

X v.o. con legge $F = 0$ $U = F \circ X$ ha dist. densità $U(0,1)$
 $\mathbb{P}_X = f(x) dx$ $F(t) = \int_{-\infty}^t f(x) dx$
 $\psi: \mathbb{R} \rightarrow \mathbb{R}$ di Borel nonnegativa

$$\int_{\mathbb{R}} \psi(u) \mathbb{P}_U(du) = \int_{\mathbb{R}} \psi(u) \mathbb{P}_{F \circ X}(du) = \int_{\mathbb{R}} \psi(F(x)) \mathbb{P}_X(dx) = \int_{\mathbb{R}} \psi(F(x)) f(x) dx$$

rispetto alla misura di Lebesgue

$$u = F(x) \quad du = F'(x) dx = 0 \text{ per p.o. } x \in \mathbb{R}$$

$$F'(x) = f(x) \quad \Rightarrow \quad du = f(x) dx$$

$$\begin{array}{ll} x \rightarrow -\infty & u \rightarrow 0 \\ x \rightarrow +\infty & u \rightarrow 1 \end{array}$$

$$\Rightarrow \int_{\mathbb{R}} \psi(u) \mathbb{P}_U(du) = \int_0^1 \psi(u) du = \int_{\mathbb{R}} \psi(u) \mathbb{1}_{(0,1)}(u) du$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(D_n \sqrt{n} \leq t) = \begin{cases} 1 - 2 \sum_{j=1}^{\infty} (-1)^{j-1} \exp(-2j^2 t^2) & t > 0 \\ 0 & t \leq 0 \end{cases}$$

Riguardo Ho se $d_n > \varepsilon$

$$\alpha = \mathbb{P}(D_n > \varepsilon) = \mathbb{P}(D_n \sqrt{n} > \varepsilon \sqrt{n}) = 1 - \mathbb{P}(D_n \sqrt{n} \leq \varepsilon \sqrt{n})$$

$$= 2 \sum_{j=1}^{\infty} (-1)^{j-1} \exp(-2j^2 \epsilon^2 n) =$$

$$= 2 \exp(-2\epsilon^2 n) + 2 \sum_{j=2}^{\infty} (-1)^{j-1} \exp(-2j^2 \epsilon^2 n)$$

serie a segni alterni

$$\alpha = 2 \exp(-2\epsilon^2 n)$$

$$\exp(-2\epsilon^2 n) = \frac{\alpha}{2} \quad -2\epsilon^2 n = \log \frac{\alpha}{2} \quad \alpha < \frac{1}{2}$$

$$\epsilon^2 = \frac{-1}{2n} \log \left(\frac{\alpha}{2} \right) = \frac{1}{2n} \log \left(\frac{2}{\alpha} \right)$$

$$\epsilon = \sqrt{\frac{1}{2n} \log \left(\frac{2}{\alpha} \right)}$$

RETTE DI REGRESSIONE

Suppongo di fare un esperimento in cui si può controllare direttamente un dato di ingresso che indico con x

$$x_1 \quad \dots \quad x_n$$

Ad ogni prova ho un risultato $y_1 \dots y_n$

(y_i la risposta in corrispondenza dell'input x_i)

$$y_i = ax_i + b \quad y_i = ax_i + b + \epsilon_i$$

Vedo la risposta dell'esperimento come una v.e. y_i di valore atteso $ax_i + b$

$$\forall i=1 \dots n \quad \sum_{i=1}^n |y_i - (ax_i + b)|^2 \rightarrow \min$$

$$\sum_{i=1}^n (y_i - (Ax_i + B))^2 \rightarrow \min$$

$$A = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{(Y_1 - Y_n)(x_1 - x_n)}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$B = \bar{Y} - A\bar{x} = b(Y_1 - Y_n)(x_1 - x_n)$$

$$Y - \bar{Y} = A(x - \bar{x}) \quad \text{si dice STIMA DELLA REGRESSIONE}$$

$$A = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i Y_i - \bar{x} Y_i - x_i \bar{Y} + \bar{x} \bar{Y}) =$$

$$= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \left(\sum_{i=1}^n x_i Y_i - \bar{x} \sum_{i=1}^n Y_i - n\bar{x}\bar{Y} + n\bar{x}\bar{Y} \right)$$

$$A = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) Y_i$$

$$B = \bar{Y} - A\bar{x} = \frac{1}{n} \sum_{i=1}^n Y_i - \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) Y_i$$

$$= \sum_{i=1}^n \left(\frac{1}{n} - \frac{x_i - \bar{x}}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) Y_i$$

$$E[Y_i] = ax_i + b$$

Y_1, \dots, Y_n indipendenti.

In generale supponiamo $\mathbb{P}_{Y_i} = \mathcal{N}(ax_i + b, \sigma^2)$

$$\mathbb{P}_{Y_i - ax_i + b} = \mathcal{N}(0, \sigma^2)$$

$$A = \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) Y_i \quad B = \sum_{i=1}^n \left(\frac{1}{n} - \frac{(x_i - \bar{x})\bar{x}}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) Y_i$$

Poiché Y_1, \dots, Y_n sono indipendenti e gaussiane lo che
 A e B sono ancora gaussiane.

$$\begin{aligned}
 E[A] &= \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) E[Y_i] = \\
 &= \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \underbrace{(2x_i + b - 2\bar{x} + 2\bar{x})}_{2(x_i - \bar{x}) + 2\bar{x} + b} \\
 &= \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n \left(2(x_i - \bar{x})^2 + (2\bar{x} + b)(x_i - \bar{x}) \right) = \\
 &= \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \left(2 \sum_{i=1}^n (x_i - \bar{x})^2 + (2\bar{x} + b) \sum_{i=1}^n (x_i - \bar{x}) \right) \\
 &= 2 \quad \left(\text{poiché } \sum_{i=1}^n (x_i - \bar{x}) = n\bar{x} - n\bar{x} = 0 \right)
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}[A] &= \text{Var} \left[\frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) Y_i \right] = \\
 &= \frac{1}{\left(\sum_{j=1}^n (x_j - \bar{x})^2 \right)^2} \sum_{i=1}^n (x_i - \bar{x})^2 \underbrace{\text{Var}[Y_i]}_{= \sigma^2 \quad \forall i=1, \dots, n}
 \end{aligned}$$

$$= \frac{1}{\left(\sum_{j=1}^n (x_j - \bar{x})^2 \right)^2} \sigma^2 \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$B = \sum_{i=1}^n \left(\frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) Y_i \quad E[Y_i] = 2x_i + b$$

$$\begin{aligned}
 E[B] &= \sum_{i=1}^n \left(\frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) \underbrace{E[Y_i]}_{= 2x_i + b} = \\
 &= \sum_{i=1}^n \frac{1}{n} (2x_i + b) - \frac{\bar{x}}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) (2(x_i - \bar{x}) + 2\bar{x} + b)
 \end{aligned}$$

$$= 2\bar{x} + \cancel{\frac{1}{n}b} - \frac{\bar{x}}{\sum_{j=1}^n (x_j - \bar{x})^2} \left(2 \sum_{i=1}^n (x_i - \bar{x})^2 + (\cancel{2\bar{x} + b}) \sum_{i=1}^n (x_i - \bar{x}) \right)$$

$n\bar{x} - n\bar{x} = 0$

$$= \cancel{2\bar{x}} + b - \cancel{2\bar{x}} = b$$

$$\text{Var}[B] = \sum_{i=1}^n \left(\frac{1}{n} - \frac{(x_i - \bar{x})\bar{x}}{\sum_{j=1}^n (x_j - \bar{x})^2} \right)^2 \text{Var}[Y_i]$$

$= \sigma^2 \quad \forall i=1, \dots, n$

$$= \sigma^2 \sum_{i=1}^n \left(\frac{1}{n^2} + \frac{\bar{x}^2 (x_i - \bar{x})^2}{\left(\sum_{j=1}^n (x_j - \bar{x})^2 \right)^2} - \frac{2\bar{x} (x_i - \bar{x})}{n \sum_{j=1}^n (x_j - \bar{x})^2} \right)$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\left(\sum_{j=1}^n (x_j - \bar{x})^2 \right)^2} \sum_{i=1}^n (x_i - \bar{x})^2 - \frac{2\bar{x} \sum_{i=1}^n (x_i - \bar{x})}{n \sum_{j=1}^n (x_j - \bar{x})^2} \right)$$

$n\bar{x} - n\bar{x} = 0$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) =$$

$$= \sigma^2 \frac{n\bar{x}^2 + \sum_{j=1}^n (x_j^2 - 2\bar{x}x_j + \bar{x}^2)}{n \sum_{j=1}^n (x_j - \bar{x})^2} =$$

$$= \sigma^2 \frac{\cancel{n\bar{x}^2} + \sum_{j=1}^n x_j^2 - \cancel{2\bar{x}n\bar{x}} + \cancel{n\bar{x}^2}}{n \sum_{j=1}^n (x_j - \bar{x})^2} = \frac{\sigma^2 \sum_{j=1}^n x_j^2}{n \sum_{j=1}^n (x_j - \bar{x})^2}$$

$$\sum_{i=1}^n |Y_i - Ax_i + B|^2 = S_R \quad \text{residuo}$$

$$\frac{S_R}{\sigma^2} \text{ he distribucione } \chi^2_{n-2}$$

A, B e S_R sono v.o. independenti.

$$\mathbb{E} \left[\frac{S_R}{n-2} \right] = \mathbb{E} \left[\frac{S_R}{\sigma^2} \frac{\sigma^2}{n-2} \right] = \frac{\sigma^2}{n-2} \mathbb{E} \left[\frac{S_R}{\sigma^2} \right] =$$

$$= \frac{\sigma^2}{\cancel{n-2}} (\cancel{n-2}) = \sigma^2$$

$$\sum_{i=1}^n (x_i - \bar{x})^2 = S_{xx}$$

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = S_{xy}$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = S_{yy}$$

$$A = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

$$B = \bar{y} - A\bar{x}$$

$$S_R = \sum_{i=1}^n (y_i - Ax_i - B)^2 = \sum_{i=1}^n (y_i - Ax_i - \bar{y} + A\bar{x})^2$$

$$= \sum_{i=1}^n \left((y_i - \bar{y}) - A(x_i - \bar{x}) \right)^2 = \sum_{i=1}^n \left((y_i - \bar{y}) - \frac{S_{xy}}{S_{xx}} (x_i - \bar{x}) \right)^2$$

$$S_R = \sum_{i=1}^n \left((y_i - \bar{y})^2 + \frac{S_{xy}^2}{S_{xx}^2} (x_i - \bar{x})^2 - \frac{2S_{xy}}{S_{xx}} (x_i - \bar{x})(y_i - \bar{y}) \right)$$

$$= S_{yy} + \frac{S_{xy}^2}{S_{xx}} \cancel{S_{xx}} - \frac{2S_{xy}}{S_{xx}} S_{xy} =$$

$$= S_{yy} - \frac{S_{xy}^2}{S_{xx}} = \frac{S_{xx}S_{yy} - S_{xy}^2}{S_{xx}}$$

— 0 —

$$H_0: \alpha = \bar{\alpha}$$

$$H_A: \alpha \neq \bar{\alpha}$$

$$P_A = N\left(\alpha, \frac{\sigma^2}{S_{xx}}\right)$$

$$P_{S_R} = \chi_{n-2}^2$$

$$Z = \frac{A - \alpha}{\sigma / \sqrt{S_{xx}}} = \frac{(A - \alpha) \sqrt{S_{xx}}}{\sigma} \quad P_Z = N(0, 1)$$

$$T = \frac{Z \sqrt{n-2}}{\sqrt{S_R / \sigma^2}} = \frac{Z \sqrt{n-2}}{\sqrt{S_R}} \sigma \Rightarrow$$

$$T = \frac{(A - \bar{a}) \sqrt{S_{xx}} \sqrt{n-2}}{\sqrt{S_R}} \quad \mathbb{P}_T = t(n-2)$$

So da $E[A] = a - \bar{a}$ SSE $A - \bar{a}$ ha valore atteso nullo SSE $E[T] = 0$

A e S_R sono indipendenti. $= 0$

$$E[T] = E\left[\frac{(A - \bar{a}) \sqrt{S_{xx}} \sqrt{n-2}}{\sqrt{S_R}}\right] =$$

$$\sqrt{S_{xx}} \sqrt{n-2} E\left[\underbrace{(A - \bar{a})}_{=0} \underbrace{S_R^{-1/2}}_{=0}\right] =$$

$$= \sqrt{S_{xx}} \sqrt{n-2} E[A - \bar{a}] E[S_R^{-1/2}] = 0 \quad \text{SSE}$$

$E[A - \bar{a}] = 0$
 $E[A] = a$

Avetto H_0 SSE

$$\left| \frac{A(x_n - x_n, y_n - y_n) - \bar{a}}{\sqrt{S_R(x_n - x_n, y_n - y_n)}} \right| < \varepsilon$$

$$\alpha = \mathbb{P}\left(\frac{|A - \bar{a}| \sqrt{S_{xx}} \sqrt{n-2}}{\sqrt{S_R}} > \varepsilon \mid E[A] = \bar{a}\right) = \mathbb{P}(|T_{n-2}| > \varepsilon)$$

$$= \mathbb{P}(T_{n-2} > \varepsilon) + \mathbb{P}(T_{n-2} < -\varepsilon) =$$

$$= 1 - F_{T_{n-2}}(\varepsilon) + F_{T_{n-2}}(-\varepsilon) = 2(1 - F_{T_{n-2}}(\varepsilon))$$

$$F_{T_{n-2}}(\varepsilon) = 1 - \frac{\alpha}{2} \quad \varepsilon = t_{n-2, 1 - \frac{\alpha}{2}}$$

INTERVALLO DI CONFIDENZA PER IL PARAMETRO α

$$\frac{(A - \alpha) \sqrt{n-2} \sqrt{S_{xx}}}{\sqrt{S_R}} \text{ ha distribuzione } t(n-2)$$

$$1 - \alpha = P \left(\frac{|A - \alpha| \sqrt{n-2} \sqrt{S_{xx}}}{\sqrt{S_R}} < t_{n-2, 1 - \frac{\alpha}{2}} \right) =$$

$$= P \left(|A - \alpha| < \frac{\sqrt{S_R} t_{n-2, 1 - \frac{\alpha}{2}}}{\sqrt{n-2} \sqrt{S_{xx}}} \right) =$$

$$= P \left(A - \frac{\sqrt{S_R} t_{n-2, 1 - \frac{\alpha}{2}}}{\sqrt{(n-2) S_{xx}}} < \alpha < A + \frac{\sqrt{S_R} t_{n-2, 1 - \frac{\alpha}{2}}}{\sqrt{(n-2) S_{xx}}} \right)$$

$$H_0: b = E[B] = \bar{b}$$

$$H_A: b = E[B] \neq \bar{b}$$

$$P_B = N \left(\bar{b}, \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n S_{xx}} \right)$$

$$Z = \frac{(B - \bar{b}) \sqrt{S_{xx} n}}{\sigma \sqrt{\sum x_i^2}} \text{ ha distribuzione } N(0, 1)$$

$$\frac{S_R}{\sigma^2} \text{ ha distribuzione } \chi^2_{n-2}$$

$$T := \frac{Z \sqrt{n-2}}{\sqrt{S_R / \sigma^2}} = \frac{(B - \bar{b}) \sqrt{S_{xx} n (n-2)}}{\cancel{\sigma \sqrt{\sum x_i^2}} \sqrt{S_R}} \text{ ha distribuzione } t(n-2)$$

$$E \left[\frac{(B - \bar{b}) \sqrt{S_{xx} n (n-2)}}{\sqrt{\sum x_i^2} \sqrt{S_R}} \right] = 0 \text{ SSE } \bar{b} = E[B]$$

$$\text{Accetto } H_0 \text{ se } \frac{|B - \bar{b}| \sqrt{S_{xx} n (n-2)}}{\sqrt{\sum x_i^2} \sqrt{S_R}} < \epsilon$$

$$\alpha = P\left(\frac{|B - \bar{b}| \sqrt{S_{xx} n(n-2)}}{\sqrt{\sum x_i^2} \sqrt{S_R}} > \varepsilon \mid \bar{b} = E[B]\right)$$

$$= P(|T_{n-2}| > \varepsilon) = 2 \left(1 - F_{T_{n-2}}(\varepsilon)\right)$$

=> ottengo livello di significat. info. e scegliendo $\varepsilon = t_{n-2, 1-\frac{\alpha}{2}}$

$$\left(B - \frac{\sqrt{\sum x_i^2} \sqrt{S_R}}{\sqrt{S_{xx} n(n-2)}} t_{n-2, 1-\frac{\alpha}{2}}, B + \frac{\sqrt{\sum x_i^2} \sqrt{S_R}}{\sqrt{S_{xx} n(n-2)}} t_{n-2, 1-\frac{\alpha}{2}} \right)$$

— 0 —

$$x_0 \quad x_n$$

$$x_0 \quad E[Ax_0 + B] = x_0 E[A] + E[B] = ax_0 + b$$

$$Ax_0 + B = Ax_0 + \bar{Y} - A\bar{x} = \bar{Y} + A(x_0 - \bar{x}) =$$

$$= \frac{1}{n} \sum_{i=1}^n Y_i + (x_0 - \bar{x}) \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_{xx}} Y_i$$

$$= \sum_{i=1}^n \left(\frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{S_{xx}} \right) Y_i$$

$$\text{Var}[Ax_0 + B] = \sum_{i=1}^n \left(\frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{S_{xx}} \right)^2 \overbrace{\text{Var}[Y_i]}^{= \sigma^2} =$$

$$= \sigma^2 \sum_{i=1}^n \left(\frac{1}{n^2} + \frac{(x_0 - \bar{x})^2 (x_i - \bar{x})^2}{S_{xx}^2} + \frac{2(x_0 - \bar{x})}{n S_{xx}} (x_i - \bar{x}) \right) =$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} + \frac{2(x_0 - \bar{x})}{n S_{xx}} \underbrace{\sum_{i=1}^n (x_i - \bar{x})}_{n\bar{x} - n\bar{x}} \right)$$

$$\text{Var}[Ax_0 + B] = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

$$Z = \frac{(Ax_0 + B) - (ax_0 + b)}{\sigma \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \quad P_Z = N(0,1)$$

$\frac{S_R}{\sigma^2}$ ha distribuzione χ_{n-2}^2

$A, B \in S_R$ sono indipendenti $\Rightarrow Z \in \frac{S_R}{\sigma^2}$ sono indipendenti.

$\Rightarrow \frac{Z \sqrt{n-2}}{\sqrt{S_R/\sigma^2}}$ ha distribuzione $t(n-2)$ cioè

$$\frac{(Ax_0 + B) - (ax_0 + b)}{\sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \cdot \frac{\sqrt{n-2}}{\sqrt{S_R}} \text{ ha distribuzione } t(n-2)$$

$$P \left(\frac{|(Ax_0 + B) - (ax_0 + b)|}{\sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \cdot \frac{\sqrt{n-2}}{\sqrt{S_R}} < t_{n-2, 1-\frac{\alpha}{2}} \right) = 1 - \alpha$$

$$Ax_0 + B \pm t_{n-2, 1-\frac{\alpha}{2}} \frac{\sqrt{S_R} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}}{\sqrt{n-2}}$$

