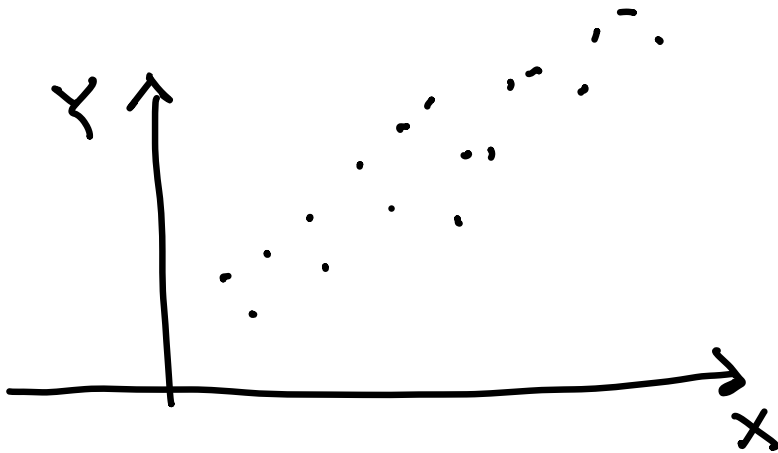


# RETTE DI REGRESSIONE LINEARE (RETTE AI MINIMI QUADRATI, LEAST SQUARES)

DATI  
 $(x_1, y_1)$   $(x_2, y_2)$  —  $(x_N, y_N)$



Sembra esserci una  
relazione lineare  
approssimata tra le  
 $x_i$  e le  $y_i$ , per  $i=1, \dots, N$

$$\boxed{Y = mX + q}$$

INCOGNITE  $m = ?$   $q = ?$

$$\begin{cases} Y_1 = mX_1 + q \\ Y_2 = mX_2 + q \end{cases}$$

$\vdots$

$$Y_N = mX_N + q$$

Ho un sistema lineare  
con solo 2 incognite e

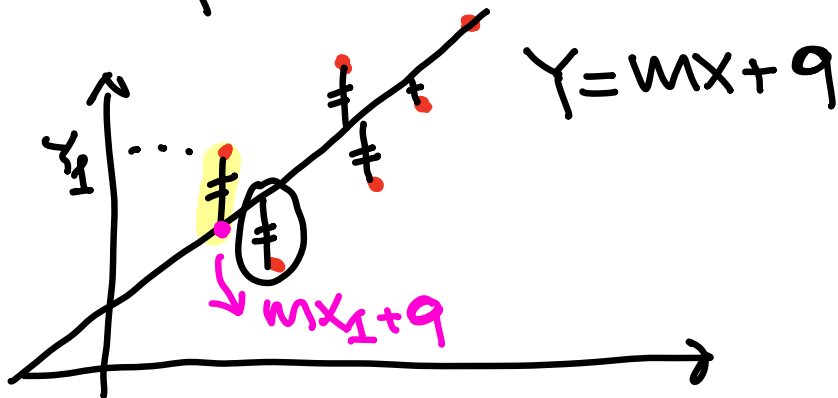
$$N \gg 2$$

equazioni.

$\xrightarrow{\text{#}}$   
dato più grande

In generale tale sistema  
non ammette soluzione.  
Fa eccezione il caso  
in cui le  $x_i$  e le  $y_i$

sono perfettamente allineate.



IDEA Cerchiamo la retta  
(il valore di  $m$  e quello  
di  $q$ ) che rende minima  
la seguente quantità

$$L(m, q) = \sum_{i=1}^N (y_i - mx_i - q)^2$$

$L$  si dice funzione di costo  
(cost function, loss function)

Per cercare il minimo  
si pone uguale a zero  
la derivata di  $L$  rispetto

ad  $m$  e la derivata di  $L$  rispetto a  $q$  :

$$\begin{cases} \frac{\partial L}{\partial m} = 0 \\ \frac{\partial L}{\partial q} = 0 \end{cases}$$

Noi daremo il valore di  $m$  e  $q$  utilizzando dei concetti di statistica.

DEF Si dice valor medio

di  $x_1 - x_N$  la

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

Analogamente

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

Nel nostro esempio (Altezza padre/figlio)

$$\bar{x} \approx 179.85$$

$$\bar{y} \approx 145.2$$

DEF Si definisce varianza di  $x_1, \dots, x_N$

$$\underbrace{s_x^2}_{\text{NOTAZIONE}} = \underbrace{\left(\frac{1}{N}\right)}_{\text{medio}} \sum_{i=1}^N \underbrace{(x_i - \bar{x})^2}_{\text{scarto}} \quad \textcircled{2} \quad \text{quadrato}$$

$s_x^2$  è uno scarto quadratico medio

La varianza fornisce informazioni su come le  $x_i$  sono distribuite intorno ad  $\bar{x}$ .  
Più  $s_x^2$  è piccola e più i dati sono vicini ad  $\bar{x}$ .

Analogamente

$$s_y^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$$

Nel nostro esempio  $S_x^2 \approx 42.63$

DEF Si dice covarianza di

$x_1, \dots, x_N$  ed  $y_1, \dots, y_N$

$$\underbrace{S_{xy}}_{\text{Notazione}} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

Per il nostro esempio

$$S_{xy} \approx 37.63$$

TEOREMA la retta ai minimi quadrati per i dati  $(x_i, y_i)$ ,

$i=1, \dots, N$ , è la retta di equazione

$$Y = mX + q$$

con  $m = \frac{S_{xy}}{S_x^2}$  e

$$q = \bar{y} - m\bar{x}$$

Oss  $(\bar{x}, \bar{y})$  appartiene  
alla retta di regressione  
lineare

$$A = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix}$$

$$\bar{A}^T = \begin{pmatrix} 1 & \text{---} & 1 \\ x_1 & \text{---} & x_N \end{pmatrix}$$

$$(\bar{A}^T A) \begin{pmatrix} m \\ q \end{pmatrix} = \bar{A}^T \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$$