

Momento del secondo ordine

$$E[XY] = \int_{-\infty}^{+\infty} xy f_{XY}(x, y) dx dy$$

Se non si conosce f_{XY} è possibile stimare il momento congiunto del secondo ordine dai dati. Date n osservazioni delle variabili X e Y , $x = (x_1, x_2, \dots, x_n)$ e $y = (y_1, y_2, \dots, y_n)$ possiamo scrivere

$$E[XY] = \frac{1}{n} \sum_{k=1}^n x_k y_k \quad (1)$$

Matlab

```
x=[x1 x2 x3 x4 ... xn];
y=[y1 y2 y3 y4 ... yn];
```

$$E[xy] = x * y' / n;$$

Covarianza

$$C_{X,Y} = E[(X - E[X])(Y - E[Y])]$$

con $E[X] = \eta_X$ e $E[Y] = \eta_Y$ valori medi di X e Y .

Se non si conoscono le distribuzioni marginali e congiunte si possono stimare i valori medi e la covarianza dai dati.

$$E[X] = \frac{1}{n} \sum_{k=1}^n x_k \quad \text{e} \quad E[Y] = \frac{1}{n} \sum_{k=1}^n y_k$$

Il momento congiunto centrale delle variabili rispetto al valor medio è dato da

$$C_{X,Y} = \frac{1}{n} \sum_{k=1}^n (x_k - \eta_X)(y_k - \eta_Y) \quad (2)$$

Si dimostra che la stima non polarizzata della covarianza delle variabili X e Y è data da

$$C_{X,Y} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \eta_X)(y_k - \eta_Y) \quad (3)$$

Matlab

```
x=[x1 x2 x3 x4 ... xn];
y=[y1 y2 y3 y4 ... yn];
```

Il momento congiunto centrale in (2) si stima nel modo seguente

$$c_xy = \text{cov}(x, y, 1);$$

La stima della covarianza in (3) si stima nel modo seguente

$c_{xy} = \text{cov}(x, y, 0)$ o $c_{xy} = \text{cov}(x, y)$

N.B. Nel seguito faremo riferimento alla notazione in (2)

Coefficiente di Correlazione

$$\rho_{XY} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[XY] - E[X]E[Y]}{\sqrt{E[(X - \eta_X)^2] E[(Y - \eta_Y)^2]}}$$

dove la deviazione standard σ_X è la radice quadrata della varianza

$$\sigma_X^2 = E[(X - \eta_X)^2]$$

Se non si conoscono le distribuzioni marginali e congiunte si possono le deviazioni standard dai dati

Il momento centrale del secondo ordine è dato da

$$\sigma_X^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \eta_X)^2 \quad (4)$$

Si dimostra che la stima non polarizzata della varianza è data da

$$\sigma_X^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \eta_X)^2 \quad (5)$$

N.B. Nel seguito faremo riferimento alla notazione in (4)

Matlab

```
x=[x1 x2 x3 x4 ... xn];
y=[y1 y2 y3 y4 ... yn];
```

Tramite il comando $C = \text{corrcoef}(x, y)$ viene stimata una matrice C di dimensione 2×2 così formata

$$C = \begin{pmatrix} \rho_{XX} & \rho_{YX} \\ \rho_{XY} & \rho_{YY} \end{pmatrix}$$

N.B. L'elemento (i, j) è calcolato come $\frac{\text{cov}(\text{Variabile}_i, \text{Variabile}_j)}{\text{std}(\text{Variabile}_i) * \text{std}(\text{Variabile}_j)}$ dove $\text{Variabile}_i = X$
 $\text{Variabile}_j = Y$

o equivalentemente $\frac{\text{cov}(\text{Variabile}_i, \text{Variabile}_j)}{\text{std}(\text{Variabile}_i) * \text{std}(\text{Variabile}_j)}$

Alcuni esempi

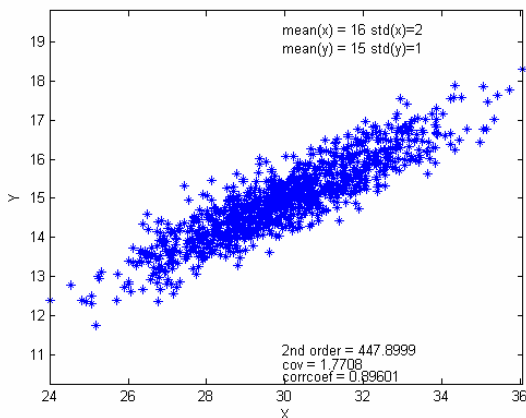
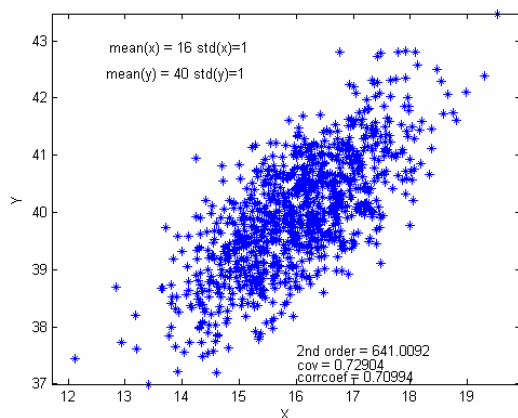
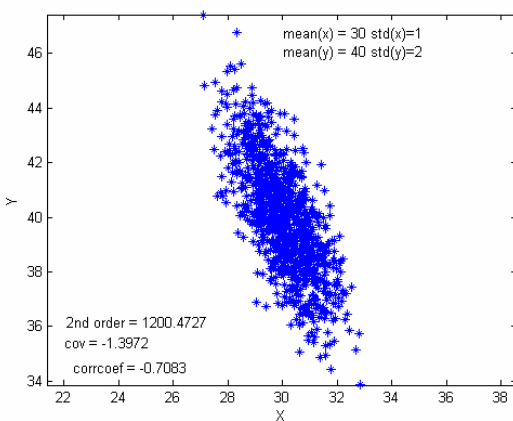
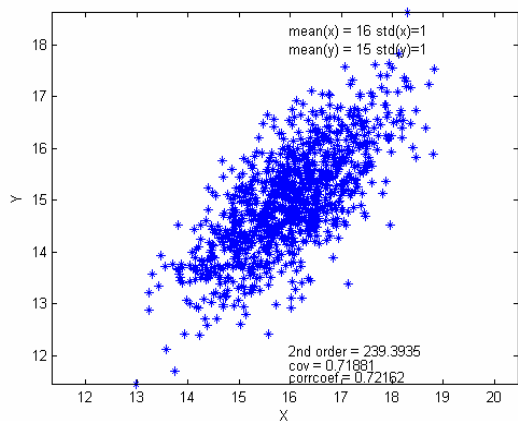
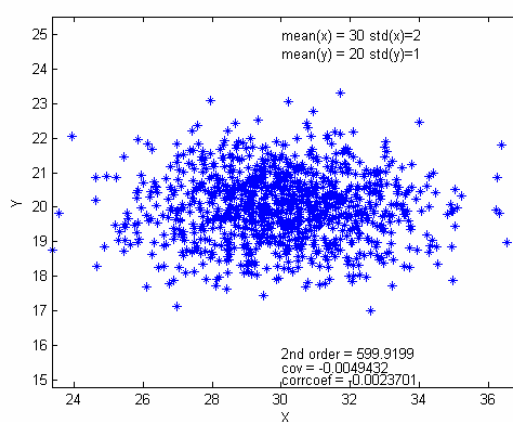
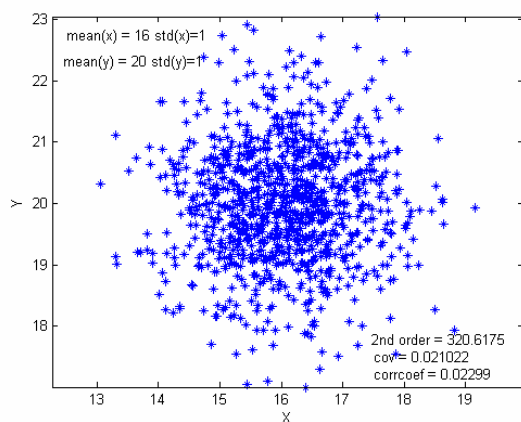
Nei grafici seguenti sono mostrati gli scatter plot relativi a variabili aleatorie X e Y a diverso valore medio, deviazione standard e coefficiente di correlazione.

Una volta generati i due vettori

$$x=[x_1 \ x_2 \ x_3 \ x_4 \ \dots \ x_n];$$

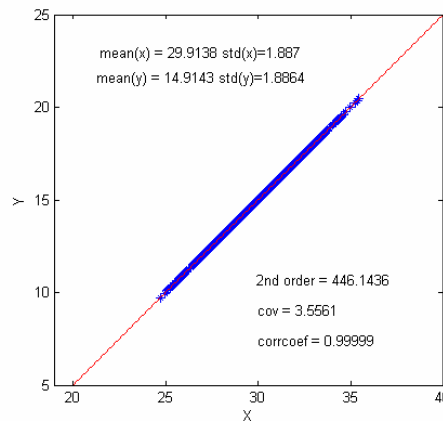
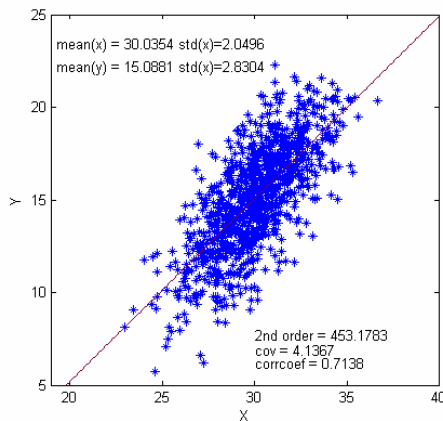
$$y=[y_1 \ y_2 \ y_3 \ y_4 \ \dots \ y_n];$$

ogni punto del piano ha coordinate (x_i, y_i) . In Matlab un grafico di questo tipo si può ottenere disegnando un simbolo, ad esempio un '*' in tali punti, eseguendo il comando `plot(x,y,'*')`



Si deve notare come i valori medi influenzino la posizione dei dati, mentre il coefficiente di correlazione insieme alle deviazioni standard siano legate al tipo di distribuzione spaziale dei dati.

Si deve precisare che il coefficiente di correlazione non è la pendenza della retta sulla quale si distribuiscono i dati. Si vedano a tale proposito i seguenti grafici. La retta può essere stimata attraverso la regressioni lineare. Vedremo in seguito la relazione che li lega.



Regressione

Nel modello di regressione lineare si assume una relazione di tipo lineare tra il valore medio della variabile dipendente Y da quello indipendente per cui

$$E(Y | X) = \eta_{Y|X} = a + b * x$$

Il modello si scrive come

$$y = a + b * x + \varepsilon$$

Dove ε è l'errore gaussiano con varianza σ^2 e valore medio nullo.

Dati due vettori delle osservazioni x e y, per gli elemento i -esimi dei vettori si può scrivere:

$$y_i = a + b * x_i + \varepsilon_i$$

per cui dati

$$x = [x_1 \ x_2 \ x_3 \ x_4 \ \dots \ x_n];$$

$$y = [y_1 \ y_2 \ y_3 \ y_4 \ \dots \ y_n];$$

si hanno n relazioni

$$y_1 = a + b * x_1 + \varepsilon_1$$

$$y_2 = a + b * x_2 + \varepsilon_2$$

.

.

.

$$y_n = a + b * x_n + \varepsilon_n$$

Usando la notazione matriciale si può scrivere:

$$y = a + bx + \varepsilon$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \quad \text{con } X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \text{ la matrice nota del modello}$$

e $\beta = [a \ b]'$ la matrice incognita dei coefficienti

In Matlab tale matrice può essere stimata come

$$\text{beta} = X \backslash y$$

gli elementi di beta sono tali da minimizzare l'errore quadratico medio $\sum_{k=1}^n (y_k - (a + b * x_k))^2$

Relazione tra coefficiente di correlazione e pendenza della retta di regressione

Date la variabile dipendente y e la variabile indipendente x il coefficiente di correlazione ρ_{xy} è legato alla pendenza della retta dalla seguente relazione

$$b = \frac{\text{cov}(x, y)}{\text{var}(x)} = \rho_{xy} \frac{\text{std}(y)}{\text{std}(x)}$$