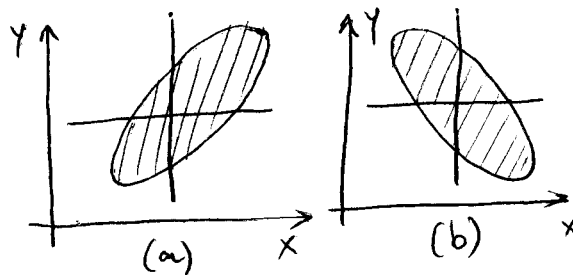
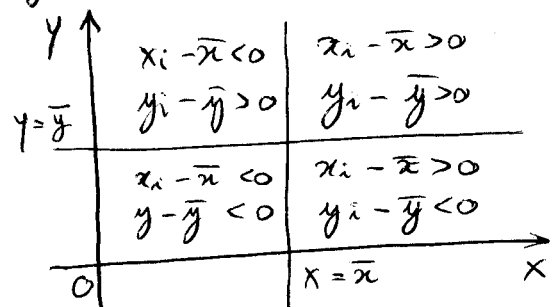


LA DIPENDENZA TRA DUE CARATTERI DI TIPO QUANTITATIVO VIENE CHIAMATA CORRELAZIONE, PER DISTINGUERLA DALLA DIPENDENZA TRA DUE CARATTERI QUALITATIVI, CHIAMATA CONNESSIONE.

SIANO X ED Y DUE VARIABILI STATISTICHE DI MEDIE \bar{x} E \bar{y} , RILEVATE SU UN COLLETTIVO DI N UNITA' STATISTICHE. SIANO x_1, x_2, \dots, x_N I VALORI OSSERVATI DI X E y_1, y_2, \dots, y_N I VALORI OSSERVATI DI Y . SI CHIAMA COVARIANZA DI X E Y , E SI INDICA CON σ_{xy} , IL VALORE COSI' DEFINITO:

$$\sigma_{xy} = \frac{\sum_{k=1}^N (x_k - \bar{x})(y_k - \bar{y})}{N}$$

SI OSSERVI COME QUESA DEFINIZIONE RICORDI QUELLA DI BARICENTRO DI UN SISTEMA DI CORPI, IN PRATICA SI TROVA IL "BARICENTRO" DI UNA DISTRIBUZIONE DI VALORI STATISTICI. SE LA COVARIANZA E' POSITIVA, VOL DIRE CHE LA MAGGIOR PARTE DEI PRODOTTI $(x_i - \bar{x})(y_i - \bar{y})$ SONO POSITIVI, E QUINDI LA MAGGIOR PARTE DEI PUNTI DI COORDINATE (x_i, y_i) DEVE CADERE ALL'INTERNO DEL QUADRANTE DEL PIANO XOY A DESTRA DI $y = \bar{y}$ E SOPRA $x = \bar{x}$, COME SI VEDE NELLO SCHEMA QUI A DESTRA, OPPURE IN QUELLO A SINISTRA DI $y = \bar{y}$ E SOTTO $x = \bar{x}$, PER CUI LA DISTRIBUZIONE DEI VALORI E' LA (a). SE INVECE LA COVARIANZA E' NEGATIVA, LA MAGGIORANZA DEI PRODOTTI $(x_i - \bar{x})(y_i - \bar{y})$ E' NEGATIVA, E QUINDI LA MAGGIOR PARTE DEI PUNTI DI COORDINATE (x_i, y_i) DEVE CADERE INTERAMENTE AGLI ALTRI DUE QUADRANTI, PER CUI LA DISTRIBUZIONE DEI VALORI E' SIMILE ALLA (b). ORA, LA (a) E' UNA DISTRIBUZIONE DI TIPO LINEARE CRESCENTE, LA (b) E' UNA DISTRIBUZIONE DI TIPO LINEARE DECRESCENTE. SE INVECE LA COVARIANZA E' NULLA, I PUNTI SONO SPARPAGLIATI SENZA ALCUNA REGOLARITA', OPPURE SONO DISPOSTI SECONDO RELAZIONI MOLTO LONTANE DA QUELLE LINEARI, COME AD ESEMPIO UNA RELAZIONE QUADRATICA.



LA COVARIANZA PUO' ESSERE CALCOLATA ANCHE CON LA FORMULA $\sigma_{xy} = \frac{\sum_{i=1}^N x_i y_i}{N} - \bar{x} \bar{y}$

UNA VOLTA APPURATA UNA CORRELAZIONE TRA DUE VARIABILI STATISTICHE X E Y , OCCORRE CAPIRE IL GRADO DI TALE CORRELAZIONE, COME SI FA CON IL TEST DEL CHI QUADRO; L'IDEA E' QUELLA DI METTERE A RAPPORTO LA COVARIANZA CON IL SUO VALORE MASSIMO. SE σ_x E σ_y SONO LE DENAZIONI STANDARD DI X ED Y , LA COVARIANZA PUO' ASSUMERE VALORI COMPRESI IN QUESTO INTERVALLO:

$$-\sigma_x \sigma_y \leq \sigma_{xy} \leq +\sigma_x \sigma_y$$

SI CHIAMA ALLORA COEFFICIENTE DI CORRELAZIONE LINEARE DI DUE VARIABILI X E Y IL COEFFICIENTE COSI' DEFINITO: (→)

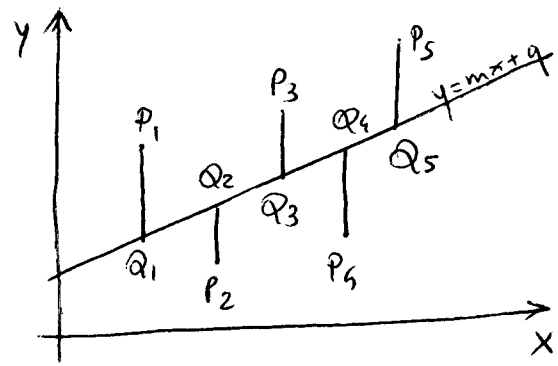
(→)

$$\rho = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

2/3

RISULTA $-1 \leq \rho \leq 1$. IL SEGNO DEL COEFFICIENTE DI CORRELAZIONE LINEARE È LO STESSO DELLA COVARIANZA; $\rho > 0$ INDICA UNA RELAZIONE LINEARE CRESCENTE, MENTRE $\rho < 0$ NE INDICA UNA DECRESCENTE. L'INDICE DI CORRELAZIONE LINEARE È UGUALE A ± 1 SE E SOLO SE LA CORRELAZIONE È PERFETTA = LINEARE; TANTO PIÙ ρ SI AVVICINA A ± 1 , TANTO PIÙ LA CORRELAZIONE È LINEARE. SE $\rho = 0$, TRA X E Y NON VI È ALCUNA CORRELAZIONE LINEARE. COME PERO' DETERMINARE L'ESATTA FUNZIONE LINEARE CHE MEGLIO INTERPRETA TALE LEGAME?

CONSIDERIAMO UNA FUNZIONE LINEARE DI EQUAZIONE $y = mx + q$ E, PER OGNI PUNTO $P(x_i, y_i)$ CHE RAPPRESENTA I DATI, CONSIDERIAMO IL CORRESPONDENTE PUNTO $Q_i(x_i, y_i')$ DI ASCISSA x_i APPARTENENTE ALLA RETTA CHE COSTITUISCE IL GRAFICO DI $y = mx + q$, A FIANCO SI VEDE LA COSTRUZIONE DEI SEGMENTI $P_i Q_i, P_2 Q_2, \dots, P_n Q_n$ CON $N = 5$. IN GENERAZIONE $P_i Q_i = |y_i - y_i'|$ ELEVANO AL QUADRATO TALI LUNGHEZZE E LE SOMMIAMO:



$$\sum_{i=1}^N (y_i - y_i')^2$$

QUESTA SOMMA ESISTE CON UN UNICO NUMERO UNA SINGOLA DISTANZA TRA I DATI y_i OSSERVATI E I VALORI TEORICI y_i' CALCOLATI SUL GRAFICO DELLA RETTA. COME FUNZIONE LINEARE CHE MEGLIO APPROSSIMA I DATI, CONVIENE SCEGLIERE QUELLA PER CUI QUESTA SOMMA RISULTA MINIMA. TALE RETTA PRENDE IL NOME DI RETTA DI REGRESSIONE.

DATTE DUE VARIABILI X E Y , DI VALORI MEDI RISPETTIVAMENTE \bar{x} E \bar{y} , LA RETTA DI REGRESSIONE CHE ESPRIME Y IN FUNZIONE DI X È LA RETTA CHE PASSA PER IL PUNTO DI COORDINATE (\bar{x}, \bar{y}) E CHE HA COME COEFFICIENTE ANGOLARE IL COSIDDETTO COEFFICIENTE DI REGRESSIONE, DEFINITO DA $m = \frac{\sigma_{xy}}{\sigma_x^2}$. L'EQUAZIONE È ADORA $y - \bar{y} = m(x - \bar{x})$. LA RETTA DI REGRESSIONE È DETTA ANCHE RETTA DEI MINIMI QUADRATI.

VEDIAMO ORA UN ESEMPIO PRATICO DI TUTTO QUESTO. SU CINQUE PAZIENTI È STATA RILEVATA L'ETÀ (X) E IL TASSO DI COLESTEROLO NEL SANGUE:

ETÀ	25	33	40	55	72
COLESTEROLO	125	130	137	146	170

DETERMINIAMO ANZITUTTO IL COEFFICIENTE DI CORRELAZIONE LINEARE, RICORRENDO ALLA COVARIANZA σ_{xy} E DALLE DEVIAZIONI STANDARD σ_x DELL'ETÀ E σ_y DEL TASSO DI COLESTEROLO. IN MODO DA AGEVOLARE I CALCOLI, CHE SONO TANTI E COMPLESSI, CONVIENE ORGANIZZARE I DATI IN UNA TABELLA COME QUESTA:

(→)

(→)

x_i	y_i	$x_i y_i$	x_i^2	y_i^2
25	125	3125	625	15625
33	130	4290	1089	16900
40	137	5480	1600	18769
55	146	8030	3025	21316
72	170	12240	5184	28900
$\sum x_i = 225$	$\sum y_i = 708$	$\sum x_i y_i = 33165$	$\sum x_i^2 = 11523$	$\sum y_i^2 = 101510$

A QUESTO PUNTO OTTERREMO:

$$\bar{x} = \frac{\sum x_i}{5} = 45$$

$$\bar{y} = \frac{\sum y_i}{5} = 141,6$$

$$\sigma_x^2 = \frac{\sum x_i^2}{5} - \bar{x}^2 = \frac{11523}{5} - 45^2 = 279,6$$

$$\sigma_y^2 = \frac{\sum y_i^2}{5} - \bar{y}^2 = \frac{101510}{5} - 141,6^2 = 251,44$$

$$\sigma_{xy} = \frac{\sum x_i y_i}{5} - \bar{x} \bar{y} = \frac{33165}{5} - 45 \cdot 141,6 = 261$$

CONCLUDIAMO PERCIÒ CHE IL COEFFICIENTE DI CORRELAZIONE LINEARE NEL NOSTRO CASO VALE:

$$\rho = \frac{\sigma_{xy}}{\sigma_x \cdot \sigma_y} = \frac{261}{\sqrt{279,6} \cdot \sqrt{251,44}} = 0,984$$

PERCIÒ LA RELAZIONE SI PUÒ CONSIDERARE LINEARE AL 98,4%. CERTAMENTE LA RETTA DI REGRESSIONE È UN MODELLO CHE INTERPRETA MOLTO BENE LA CORRELAZIONE TRA X E Y. RICORDIAMO CHE ρ VIENE CHIAMATO ANCHE INDICE DI BRAVASS - PEARSON, ESSENDO STATO INTRODOTTO DAL MATEMATICO FREGUENSE BRAVASS (1811-1863) E KARL PEARSON (1857-1936). SCRIVIAMO ALLORA NEL NOSTRO CASO L'EQUAZIONE DELLA RETTA DI REGRESSIONE. IL SUO COEFFICIENTE ANGOLARE È DATO DALLA FORMULA SOPRA RIPORTATA:

$$m = \frac{\sigma_{xy}}{\sigma_x^2} = \frac{261}{279,6} = 0,933$$

L'EQUAZIONE DELLA RETTA DI REGRESSIONE PERCIÒ È:

$$y - 141,6 = 0,933(x - 45)$$

CIÒ È:

$$y = 0,933x + 99,594$$

AD ESEMPIO, SE IN ESSA INTERESSO $x = 55$, OTTIENGO $y = 150,909$, UN VALORE MOLTO VICINO ALL' $y = 146$ OSSERVATO (LO SCARTO È APPENA DEL 3%)