

# Γραμμική συσχέτιση τυχαίων μεταβλητών

## Ανεξάρτητες τυχαίες μεταβλητές (independent random variables)

Δύο τυχαίες μεταβλητές  $X, Y$  λέγονται ανεξάρτητες αν η μία δεν επηρεάζει την πιθανότητα πραγματοποίησης της άλλης. Ο ορισμός είναι ο ακόλουθος

$$P(X = a, Y = b) = P(X = a) \cdot P(Y = b),$$

για κάθε δύο αριθμούς  $a, b$ .

Ας θεωρήσουμε την ρήψη 2 νομισμάτων και τις μεταβλητές  $X, Y$  οι οποίες ορίζονται ως εξής

	K	Γ
$X(\omega)$	0	1
$Y(\omega)$	0	1

Οι  $X, Y$  είναι ανεξάρτητες εφόσον η ρήψη του κάθε νομίσματος δεν επηρεάζει την ρήψη του άλλου. Αυτό μπορούμε να το πετύχουμε ρίχνοντας τα νομίσματα σε διαφορετικά μέρη του δωματίου.

$$P(X = 0, Y = 1) = \frac{1}{4}$$

$$P(X = 0) \cdot P(Y = 1) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

Μια σημαντική ιδιότητα των ανεξάρτητων τυχαίων μεταβλητών είναι η ακόλουθη

$$\mathbf{E}(X \cdot Y) = \mathbf{E}X \cdot \mathbf{E}Y$$

Η μέση τιμή των  $\mathbf{X}$ ,  $\mathbf{Y}$  στο προηγούμενο παράδειγμα είναι  $1/2$  και συνεπώς η μέση τιμή του γινομένου τους θα είναι

$$\mathbf{E}(X \cdot Y) = 1/4.$$

Ποια είναι η μέση τιμή του γινομένου των  $\mathbf{X}$ ,  $\mathbf{Y}$  κατά την ρήψη 2 τίμιων ζαριών, όπου  $\mathbf{X}$  είναι η ένδειξη του πρώτου και  $\mathbf{Y}$  η ένδειξη του δεύτερου ζαριού;

## Συνδιακύμανση (covariance)

Για δύο τυχαίες μεταβλητές  $X, Y$  ορίζουμε την συνδιακύμανση τους  $\text{cov}(X, Y)$  ως εξής:

$$\text{cov}(X, Y) = \mathbf{E}(X - \mu_X) \cdot (Y - \mu_Y),$$

$$\mu_X = \mathbf{E}X \text{ και } \mu_Y = \mathbf{E}Y.$$

Κάνοντας πράξεις παίρνουμε

$$\text{cov}(X, Y) = \mathbf{E}(X \cdot Y) - \mathbf{E}X \cdot \mathbf{E}Y.$$

(\*) Να θυμίσουμε ότι ο τύπος της διακύμανσης είναι

$$\text{Var}(X) = \mathbf{E}X^2 - (\mathbf{E}X)^2.$$

Η συνδιακύμανση είναι ας πούμε μια ποσότητα που εκφράζει κατά πόσο «συσχετίζονται» οι δύο τυχαίες μεταβλητές και μπορεί να πάρει τόσο θετικές όσο και αρνητικές τιμές.

Αν οι  $X, Y$  είναι ανεξάρτητες τότε η συνδιακύμανση τους είναι μηδέν,

$$\text{cov}(X, Y) = \mathbf{E}X \cdot \mathbf{E}Y - \mathbf{E}X \cdot \mathbf{E}Y = 0$$

**Παράδειγμα.** Ας θεωρήσουμε το πείραμα της ρήψης 2 τίμιων νομισμάτων. Θα θεωρήσουμε τις τυχαίες μεταβλητές  $X, Y, Z$  τις οποίες ορίζουμε στο παρακάτω πίνακάκι

	ΚΚ	ΚΓ	ΓΚ	ΓΓ
X	1	1	0	0
Y	1	0	1	0
Z	2	1	1	0

Οι  $X, Y$  είναι ανεξάρτητες γιατί η  $X$  αφορά το αποτέλεσμα του «πρώτου» νομίσματος ενώ η  $Y$  το αποτέλεσμα του «δεύτερου».

$$\left. \begin{array}{l} \mathbf{EX} = \mathbf{EY} = 1/2 \\ \mathbf{E}(X \cdot Y) = \mathbf{EX} \cdot \mathbf{EY} \end{array} \right\} \Rightarrow \text{cov}(X, Y) = 0$$

Οι  $X, Z$  όμως δεν είναι ανεξάρτητες, αν η  $Z$  έχει την τιμή 2, τότε η  $X$  έχει αναγκαστικά την τιμή 1. Για να υπολογίσουμε την  $\text{cov}(X, Z)$  κάνουμε πράξεις,

$$\left. \begin{array}{l} \mathbf{E}(X \cdot Z) = 2 \cdot P(X \cdot Z = 2) + 1 \cdot P(X \cdot Z = 1) + 0 \cdot P(X \cdot Z = 0) \\ \quad = 2 \cdot \frac{1}{4} + 1 \cdot \frac{1}{4} = \frac{3}{4} \\ \mathbf{EZ} = 2 \cdot P(Z = 2) + 1 \cdot P(Z = 1) + 0 \cdot P(Z = 0) = 1 \\ \mathbf{EX} = 1/2 \end{array} \right\} \Rightarrow \text{cov}(X, Z) = \frac{1}{4}$$

## Συντελεστής (γραμμικής) συσχέτισης (correlation coefficient)

Δεν είναι όμως η συνδιακύμανση η ποσότητα που χρησιμοποιούμε για να εκφράσουμε κατά πόσο συσχετίζονται δύο τυχαίες μεταβλητές  $X, Y$  αλλά ο συντελεστής συσχέτισης  $\rho(X, Y)$  που ορίζεται ως εξής,

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma(X) \cdot \sigma(Y)}$$

με την προϋπόθεση οι τυπικές αποκλίσεις  $\sigma(X)$  και  $\sigma(Y)$  να μην είναι μηδέν.

Ο συντελεστής συσχέτισης παίρνει τιμές από -1 έως 1, δηλαδή ισχύει

$$-1 \leq \rho(X, Y) \leq 1$$

$$\rho(X, Y) = 0 \rightarrow X, Y \text{ ασυσχέτιστες}$$

$$\rho(X, Y) > 0 \rightarrow X, Y \text{ θετικά συσχετισμένες}$$

$$\rho(X, Y) < 0 \rightarrow X, Y \text{ αρνητικά συσχετισμένες}$$

Στην πράξη όταν δύο μεταβλητές είναι θετικά συσχετισμένες και η μία από τις δύο είναι μεγάλη, αυτό ενισχύει την πιθανότητα και η άλλη να είναι μεγάλη. Όταν είναι αρνητικά συσχετισμένες και η μία από τις δύο είναι μεγάλη, τότε αυτό ενισχύει την πιθανότητα η άλλη να είναι μικρή. Όταν είναι ασυσχέτιστες και γνωρίζουμε κάτι για την  $X$  δεν μπορούμε να πούμε κάτι για την  $Y$ .

$X, Y$  ανεξάρτητες συνεπάγει  $X, Y$  ασυσχέτιστες (το αντίστροφο δεν ισχύει)

Όσο πιο κοντά είμαστε στο 1 ή -1 τόσο μεγαλώνει και η σιγουριά μας όταν πιθανολογούμε για την τιμή της  $Y$  όταν γνωρίζουμε την τιμή της  $X$ .

Στην ακραία περίπτωση όπου  $\rho(\mathbf{X}, \mathbf{Y})=1$  ή  $-1$ , τότε οι  $\mathbf{X}, \mathbf{Y}$  συνδέονται με μία σχέση της μορφής  $\mathbf{Y}=\alpha\mathbf{X}+\beta$ , δηλαδή αν γνωρίζουμε την τιμή της  $\mathbf{X}$  γνωρίζουμε με πιθανότητα  $1$  την τιμή της  $\mathbf{Y}$ .

**Παράδειγμα.** Ας πάμε πάλι στο πείραμα της ρήψης 2 τίμιων νομισμάτων και τις τυχαίες μεταβλητές  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$  τις οποίες ορίσαμε πριν

	ΚΚ	ΚΓ	ΓΚ	ΓΓ
X	1	1	0	0
Y	1	0	1	0
Z	2	1	1	0

Η συνδιακύμανση των  $\mathbf{X}, \mathbf{Y}$  είναι όπως είδαμε μηδέν και συνεπώς ο συντελεστής συσχέτισης θα είναι και αυτός μηδέν,

$$\text{cov}(\mathbf{X}, \mathbf{Y}) = 0 \Rightarrow \rho(\mathbf{X}, \mathbf{Y}) = \frac{0}{\sigma(\mathbf{X}) \cdot \sigma(\mathbf{Y})} = 0$$

Οι τυπικές αποκλίσεις των  $\mathbf{X}, \mathbf{Y}$  είναι  $1/2$ .

Ο συντελεστής συσχέτισης των  $\mathbf{X}, \mathbf{Z}$  θα είναι

$$\text{cov}(\mathbf{X}, \mathbf{Z}) = \frac{1}{4} \Rightarrow \rho(\mathbf{X}, \mathbf{Z}) = \frac{1/4}{\sigma(\mathbf{X}) \cdot \sigma(\mathbf{Z})} = \frac{1/4}{(1/2) \cdot (1/\sqrt{2})} = \frac{\sqrt{2}}{2}$$

Αναμενόμενο να έχουμε θετική συσχέτιση. Για «μεγάλες» τιμές της  $\mathbf{X}$  αναμένουμε «μεγάλες» τιμές της  $\mathbf{Z}$ . Θεωρήστε την μεταβλητή  $\mathbf{Z}'$  η οποία μετράει τα «Γράμματα». Τώρα ποια είναι η εκτίμηση σας για το πρόσημο της  $\rho(\mathbf{X}, \mathbf{Z}')$ ;

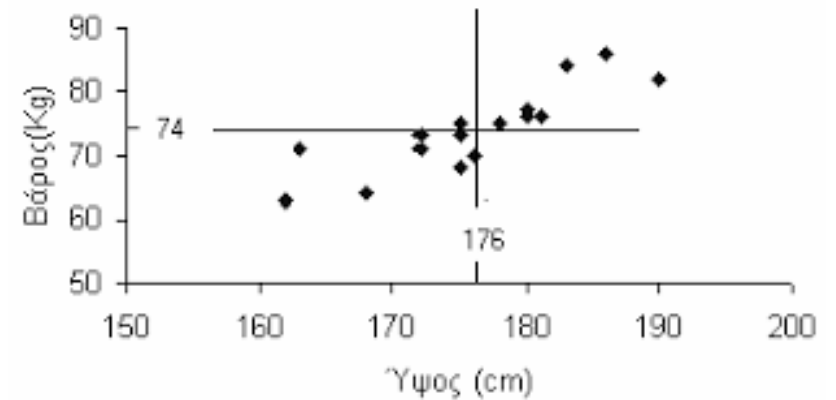
## Διάγραμμα διασποράς (scatter plot)

Ας υποθέσουμε ότι εκτελούμε ένα πείραμα πολλές φορές και παίρνουμε ζεύγη τιμών  $(x,y)$  δύο τυχαίων μεταβλητών  $\mathbf{X},\mathbf{Y}$ , οπότε προκύπτει ένα πλήθος  $n$  δειγμάτων

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n).$$

Φανταστείτε για παράδειγμα ότι μετράτε το βάρος και το ύψος των Ελλήνων. Τα δείγματα αυτά που είναι ζευγάρια αριθμών μπορούμε να τα αποτυπώσουμε στο επίπεδο και να πάρουμε το διάγραμμα διασποράς, όπως φαίνεται στο διπλανό σχήμα. Σε κάθε σημείο αντιστοιχεί και ένα δείγμα.

Όσο μεγαλύτερο είναι το  $n$  τόσο περισσότερα πράγματα μας λέει το διάγραμμα διασποράς για την κατανομή του ζεύγους  $(\mathbf{X},\mathbf{Y})$ .

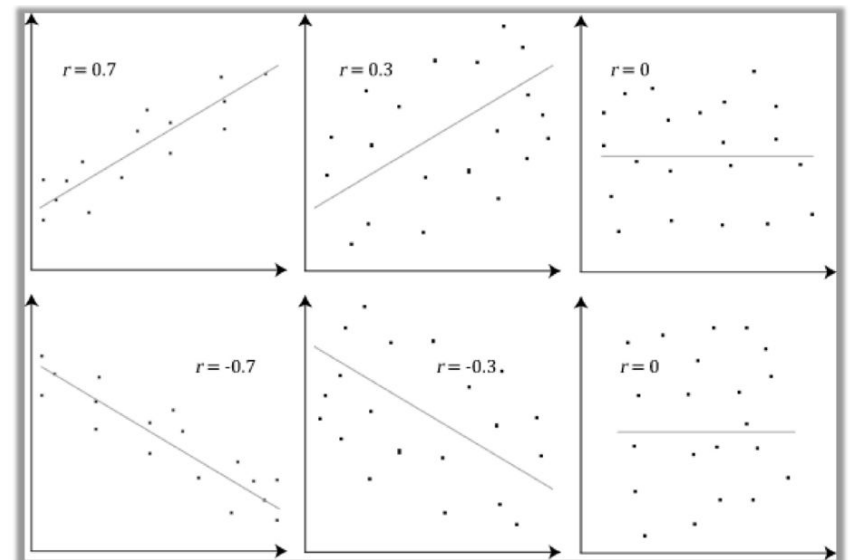


Μπορούμε από το διάγραμμα διασποράς να βγάλουμε συμπεράσματα για την συσχέτιση των  $\mathbf{X},\mathbf{Y}$ ; Η απάντηση είναι συνήθως ναι, αρκεί να έχουμε αρκετά μεγάλο πλήθος δειγμάτων.

Ο συντελεστής συσχέτισης μας λέει κατά πόσον οι δύο μεταβλητές σχετίζονται μέσω μιας σχέσης της μορφής

$$Y = \alpha \cdot X + \beta$$

και όπως γνωρίζουμε δύο μεταβλητές που συνδέονται με μια τέτοια σχέση σχηματίζουν μία ευθεία γραμμή στο επίπεδο.



Μπορούμε να εκτιμήσουμε τον συντελεστή συσχέτισης δύο μεταβλητών  $\mathbf{X}, \mathbf{Y}$  από τα δείγματα που έχουμε συλλέξει αντικαθιστώντας στον τύπο

$$\rho(\mathbf{X}, \mathbf{Y}) = \frac{\text{cov}(\mathbf{X}, \mathbf{Y})}{\sigma(\mathbf{X}) \cdot \sigma(\mathbf{Y})}$$

την δειγματική συνδιακύμανση  $\mathbf{s}(\mathbf{X}, \mathbf{Y})$  και τις δειγματικές τυπικές αποκλίσεις  $\mathbf{s}(\mathbf{X}), \mathbf{s}(\mathbf{Y})$ , οπότε θα πάρουμε τον τύπο του δειγματικού συντελεστή συσχέτισης  $\mathbf{r}(\mathbf{X}, \mathbf{Y})$ ,

$$\mathbf{r}(\mathbf{X}, \mathbf{Y}) = \frac{\mathbf{s}(\mathbf{X}, \mathbf{Y})}{\mathbf{s}(\mathbf{X}) \cdot \mathbf{s}(\mathbf{Y})}.$$

Οι τύποι είναι οι ακόλουθοι,

$$\mathbf{s}(\mathbf{X}, \mathbf{Y}) = \frac{1}{\nu - 1} [(x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + \dots + (x_\nu - \bar{x})(y_\nu - \bar{y})]$$

$$\mathbf{s}(\mathbf{X}) = \sqrt{\frac{1}{\nu - 1} [(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_\nu - \bar{x})^2]}$$

$$\mathbf{s}(\mathbf{Y}) = \sqrt{\frac{1}{\nu - 1} [(y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + \dots + (y_\nu - \bar{y})^2]}$$

και μετά από πράξεις και απλοποιήσεις παίρνουμε τους τύπους

$$\mathbf{r}(\mathbf{X}, \mathbf{Y}) = \frac{\sum_{i=1}^{\nu} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{\nu} (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^{\nu} (y_i - \bar{y})^2}} = \frac{\sum_{i=1}^{\nu} x_i y_i - \nu \bar{x} \cdot \bar{y}}{\sqrt{\sum_{i=1}^{\nu} x_i^2 - \nu \bar{x}^2} \cdot \sqrt{\sum_{i=1}^{\nu} y_i^2 - \nu \bar{y}^2}}$$

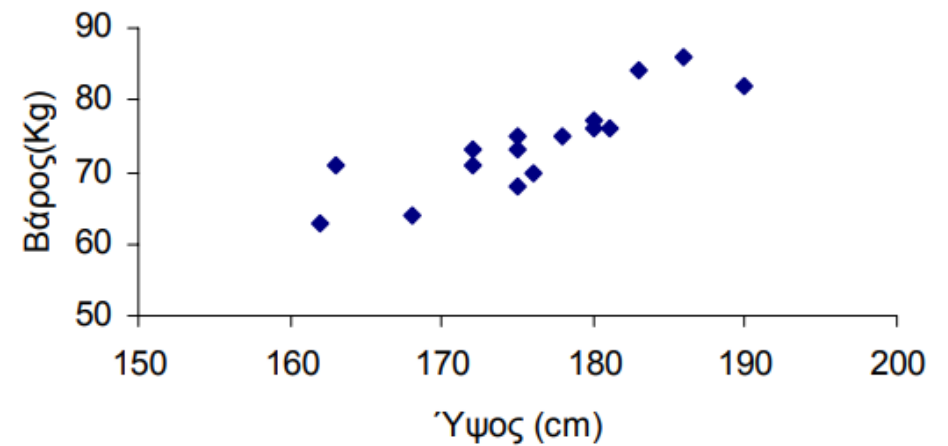


Για παράδειγμα θα πάρουμε τα διπλανά δείγματα από 16 εργάτες μίας βιομηχανίας.

	1	2	3	4	5	6	7	8
Ύψος (cm)	183	162	172	181	180	168	176	180
Βάρος (Kg)	84	63	71	76	77	64	70	76

	9	10	11	12	13	14	15	16
Ύψος (cm)	190	175	178	175	186	172	175	163
Βάρος (Kg)	82	68	75	73	86	73	75	71

Το διάγραμμα διασποράς φαίνεται στο διπλανό σχήμα και εκτιμούμε αμέσως μία θετική συσχέτιση μεταξύ του ύψους και του βάρους, όπως άλλωστε ήταν και αναμενόμενο.



Μπορείτε να κάνετε τις πράξεις και να τον υπολογίσετε ή να χρησιμοποιήσετε ένα πρόγραμμα επεξεργασίας δεδομένων, όπως είναι το PSCP ή το SPSS.

Πολλές φορές όταν δεν έχουμε έναν υπολογιστή ίσως χρειαστεί να κάνουμε πράξεις φτιάχνοντας ένα τέτοιο πίνακάκι

$x_i$	$y_i$	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\sum x_i$	$\sum y_i$			$\sum (x_i - \bar{x})^2$	$\sum (y_i - \bar{y})^2$	$\sum (x_i - \bar{x})(y_i - \bar{y})$

και εφαρμόζοντας τον τύπο

$$r(X, Y) = \frac{\sum_{i=1}^v (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^v (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^v (y_i - \bar{y})^2}}.$$

Ή μπορούμε να φτιάξουμε το πίνακάκι

$x_i$	$y_i$	$x_i^2$	$y_i^2$	$x_i y_i$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\sum x_i$	$\sum y_i$	$\sum x_i^2$	$\sum y_i^2$	$\sum x_i y_i$

και να εφαρμόσουμε τον τύπο

$$r(X, Y) = \frac{\sum_{i=1}^v x_i y_i - v \bar{x} \cdot \bar{y}}{\sqrt{\sum_{i=1}^v x_i^2 - v \bar{x}^2} \cdot \sqrt{\sum_{i=1}^v y_i^2 - v \bar{y}^2}}$$

Παράδειγμα 1.

$x_i$	$y_i$	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x}) \cdot (y_i - \bar{y})$
1	4	-2	4	4	16	-8
2	2	-1	2	1	4	-2
3	0	0	0	0	0	0
4	-2	1	-2	1	4	-2
5	-4	2	-4	4	16	-8
$\sum x_i = 15$	$\sum y_i = 0$			$\sum (x_i - \bar{x})^2 = 10$	$\sum (y_i - \bar{y})^2 = 40$	$\sum (x_i - \bar{x}) \cdot (y_i - \bar{y}) = -20$

$$\bar{x} = 3, \bar{y} = 0$$

$$r = \frac{\sum_{i=1}^5 (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^5 (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^5 (y_i - \bar{y})^2}} = \frac{-20}{\sqrt{10} \cdot \sqrt{40}} = -1$$

Παράδειγμα 2.

$x_i$	$y_i$	$x_i^2$	$y_i^2$	$x_i \cdot y_i$
1	-2	1	4	-2
3	0	9	0	0
5	1	25	1	5
7	3	49	9	21
9	5	81	25	45
10	6	100	36	60
12	8	144	64	96
13	10	169	100	130
$\sum x_i = 60$	$\sum y_i = 31$	$\sum x_i^2 = 578$	$\sum y_i^2 = 239$	$\sum x_i \cdot y_i = 355$

$$\bar{x} = 7,5 \text{ και } \bar{y} = 3,9$$

$$r = \frac{s_{xy}}{s_x \cdot s_y} = \frac{\sum_{i=1}^n x_i y_i - n \cdot \bar{x} \cdot \bar{y}}{\sqrt{\sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2} \cdot \sqrt{\sum_{i=1}^n y_i^2 - n \cdot \bar{y}^2}} = \frac{355 - 8 \cdot 7,5 \cdot 3,9}{\sqrt{578 - 8 \cdot 7,5^2} \cdot \sqrt{239 - 8 \cdot 3,9^2}} = 0,99$$