

## Προβλήματα

Φροντιστήριο #1

Συσχέτιση δύο μεταβλητών

Θέλουμε να μελετήσουμε τη σχέση μεταξύ δύο μεταβλητών (ποσοτήτων, χαρακτηριστικών)  $X$  και  $Y$ . Ένας πρώτος απλός τρόπος για να διαπιστωθεί αν υπάρχει συσχέτιση μεταξύ των δύο μεταβλητών είναι η κατασκευή ενός διαγράμματος διασποράς (scatter plot). Το διάγραμμα διασποράς καθορίζει (προσδιορίζει) τη πιθανή σχέση μεταξύ των μεταβλητών  $X$  και  $Y$ . Διάφορες μορφές διαγραμμάτων διασποράς εμφανίζονται στα σχήματα 8.1-8.3 του βιβλίου (I. Χατζηνιάς).

Μία ποσοτική μέτρηση της έντασης της συσχέτισης μεταξύ δύο μεταβλητών  $X$  και  $Y$  μπορεί να γίνει με τη χρήση του συντελεστή συσχέτισης  $r = r(X, Y)$ . Αν υποθέσουμε ότι για τις τιμές  $x_1, x_2, \dots, x_n$  που δίνουμε στη μεταβλητή  $X$  καταγράφουμε τις αντίστοιχες τιμές  $y_1, y_2, \dots, y_n$  που λαμβάνει η μεταβλητή  $Y$ . Τότε, για τις μετρήσεις  $x_1, x_2, \dots, x_n$  της  $X$  και  $y_1, y_2, \dots, y_n$  της  $Y$  ο (δειγματικός) συντελεστής συσχέτισης  $r = r(X, Y)$  ορίζεται ως εξής:

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

$$\text{όπως } \bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad (\text{αριθμητικός μέσος των μετρήσεων (τιμών) της } X)$$

-2-

$$\text{και } \bar{y} = \frac{\sum_{i=1}^n y_i}{n} \quad (\text{ο αριθμητικός μέσος των μετρήσεων (τιμών) της } Y)$$

Όπως,

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

$$= \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \bar{y} - \sum_{i=1}^n \bar{x} y_i + \sum_{i=1}^n \bar{x} \bar{y}$$

$$= \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} - \bar{x} n \bar{y} + n \bar{x} \bar{y}$$

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

και

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

$$= \sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i \bar{x} + \sum_{i=1}^n \bar{x}^2$$

$$= \sum_{i=1}^n x_i^2 - 2 \bar{x} \sum_{i=1}^n x_i + n \bar{x}^2 = \sum_{i=1}^n x_i^2 - 2 \bar{x} n \bar{x} + n \bar{x}^2$$

$$= \sum_{i=1}^n x_i^2 - 2n \bar{x}^2 + n \bar{x}^2 = \sum_{i=1}^n x_i^2 - n \bar{x}^2$$

$$= \sum_{i=1}^n x_i^2 - n \left( \frac{\sum_{i=1}^n x_i}{n} \right)^2 = \sum_{i=1}^n x_i^2 - \frac{\left( \sum_{i=1}^n x_i \right)^2}{n}$$

$$= \frac{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}{n}$$

Ομοίως

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - n\bar{y}^2 = \frac{n \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i\right)^2}{n} \quad -3-$$

Μπορούμε να γράψουμε:

$$r = r(X, Y) = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sqrt{\left(n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2\right) \left(n \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i\right)^2\right)}}$$

Αποδεικνύεται ότι:  $-1 \leq r(X, Y) \leq +1$ .

- Τιμές του  $r$  που τείνουν στην τιμή  $-1$ , υποδηλώνουν έντονη αρνητική γραμμική συσχέτιση μεταξύ των  $X$  και  $Y$ .

- Τιμές του  $r$  που τείνουν στην τιμή  $+1$ , υποδηλώνουν έντονη θετική γραμμική συσχέτιση μεταξύ των  $X$  και  $Y$ .

- Τιμές του  $r$  που τείνουν στην τιμή  $0$ , υποδηλώνουν μη-ύπαρξη γραμμικής συσχέτισης μεταξύ των  $X$  και  $Y$ .

-  $0$  συντελεστής γραμμικής συσχέτισης  $r$  καλείται δείκτης συντελεστής γραμμικής συσχέτισης διότι υπολογίζεται από (το δείγμα) τις μετρήσεις των μεταβλητών  $X$  και  $Y$ .

- Δεν εκφράζεται σε κάποια μονάδα μέτρησης, είναι ένας καθαρός αριθμός.

Έλεγχος στατιστικής σημαντικότητας του συντελεστή -4-  
συσχέτισης.

Θέλουμε να ελέγξουμε αν η γραμμική συσχέτιση μεταξύ των  $X$  και  $Y$  είναι στατιστικά σημαντική. Βασίζόμαστε στο δείγμα των μετρήσεων των μεταβλητών  $X$  και  $Y$ . Θέλουμε να διερευνήσουμε αν ο συντελεστής γραμμικής συσχέτισης του πληθυσμού των μετρήσεων των μεταβλητών  $X$  και  $Y$ ,  $\rho$  είναι διάφορος του μηδενός. Οι προς έλεγχο υποθέσεις είναι:

$$H_0: \rho = 0 \text{ (Δεν υπάρχει γραμμική συσχέτιση μεταξύ των } X \text{ και } Y)$$

$$H_1: \rho \neq 0 \text{ (Υπάρχει στατιστικά σημαντική γραμμική συσχέτιση μεταξύ των } X \text{ και } Y)$$

Ο έλεγχος υποθέσεων πραγματοποιείται με το κριτήριο  $t$ . Χρησιμοποιούμε τη ελεγχοσυνάρτηση:

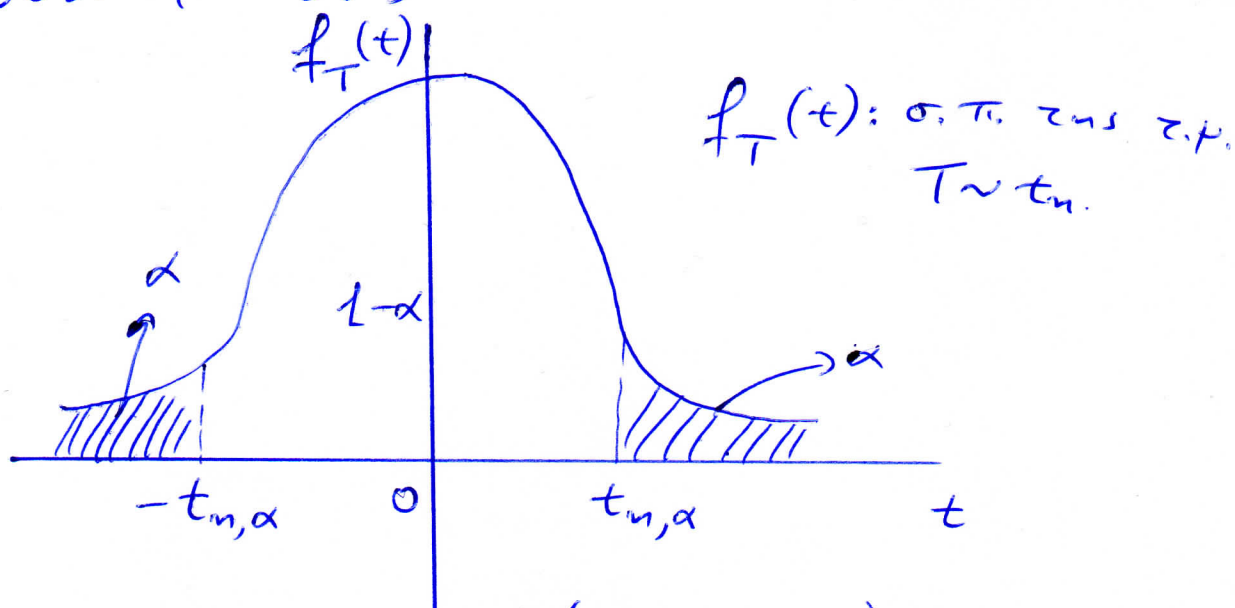
$$T = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

Όταν ισχύει η  $H_0$ ,  
(υπό την  $H_0$ )  
 $T \sim t_{n-2}$

Απαραίτητη υπόθεση: Οι μεταβλητές  $X$  και  $Y$  να κατανοούνται τουλάχιστον οριακά σύμφωνα με την κανονική κατανομή. Το κριτήριο  $t$  δεν είναι ιδιαίτερα ευαίσθητο στην υπόθεση της κανονικότητας.

ΥΠΕΝΟΤΗΜΙΣΗ: Κατανομή  $t_{n-2}$ -student με βαθμούς ελευθερίας. Η γ.π. της σ.π. της  $T \sim t_{n-2}$  έχει κωδωνοειδή μορφή, είναι συμμετρική ως προς τον

κατακόρυφο άξονα στο 0 και έχει πιο παχιές ουρές -5-  
 από την  $N(0,1)$  (είναι πιο πεπλατυσμένη). Όσο το  $n$   
 μεγαλώνει προσεγγίζεται πολύ καλά από την  $N(0,1)$



Αν  $T \sim t_n$  τότε  $t_{n, \alpha} : P(T > t_{n, \alpha}) = \alpha$

Επίσης, λόγω συμμετρίας  $t_{n, 1-\alpha} = -t_{n, \alpha}$  ☒

Απόφαση: Απορρίπτουμε την  $H_0$  σε εσο  $\alpha$  αν

$$|T| > t_{n-2, \alpha/2}$$

Προσοχή!! Μία χαμηλή τιμή του  $r$  δεν σημαίνει ότι  
 πάντα η σχέση μεταξύ των  $X$  και  $Y$  είναι ασθενής.

Οι μεταβλητές  $X$  και  $Y$  μπορεί να συσχετίζονται  
 έντονα αλλά η σχέση να είναι καρπυλόγραφη.

Παραδείγματα

#1 Σε 10 παιδιά ηλικίας 10-15 ετών μετρήθηκε το  
 βάρος  $Y$  σε Kg και το ύψος  $X$  σε cm.

$X$	130	148	154	145	162	158	152	144	138	164
$Y$	34	42	45	41	46	49	43	46	36	50

Να εξετασθεί αν το βάρος είναι ή όχι ανεξάρτητο του

ύψους.

Λύση

Μετά από υπολογισμούς έχουμε:  $\bar{x} = 249,5$

$\bar{y} = 43,2$       $\sum_{i=1}^{10} (x_i - \bar{x})^2 = 10,7^2$ ,      $\sum_{i=1}^{10} (y_i - \bar{y})^2 = 5,18^2$

$\sum_{i=1}^{10} (x_i - \bar{x})(y_i - \bar{y}) = 49,88$

Ελέγχουμε τις υποθέσεις:  $H_0: \rho = 0$

vs  
 $H_1: \rho \neq 0$

όπου  $\rho$  είναι ο συντελεστής συσχέτισης για τον πληθυσμό των μετρήσεων των μεταβλητών  $X$  και  $Y$ . Υποθέτουμε ότι οι μεταβλητές  $X$  και  $Y$  κατανέμονται (έστω οριακά) κανονικά.

Χρησιμοποιούμε τη σ.σ. ελέγχου:

$$T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$
 όπου  $n$ : μέγεθος δείγματος μετρήσεων  
(εδώ  $n=10$ )

και  $r$  ο δείγματός συντελεστής συσχέτισης των  $X$  και  $Y$ . Από τα δεδομένα, έχουμε:

$$r = \frac{49,88}{\sqrt{10,7^2 \cdot 5,18^2}} = 0,9$$

Ε συνεπώς 
$$T = \frac{0,9\sqrt{8}}{\sqrt{1-0,9^2}} \approx 5,84$$

Η  $H_0$  απορρίπτεται σε ε.σ.σ.  $\alpha = 5\%$ , αν

$|T| > t_{n-2, \alpha/2} = t_{8, 0,025} = 2,306$  (από πίνακες) της  $t_\beta$

Συμπέρασμα: η  $H_0$  απορρίπτεται σε ε.σ.σ.  $\alpha = 5\%$  -7-

Άρα το βάρος δεν φαίνεται να είναι ανεξάρτητο του ύψους. Δηλαδή τα δεδομένα παρέχουν ισχυρές ενδείξεις ότι το βάρος συσχετίζεται στατιστικά σημαντικά με το ύψος.  $\boxtimes$

#2 Ένα ζ.δ.  $n=6$  ζευγών μετρήσεων των μεταβλητών  $X$  και  $Y$  που προέρχεται από δύο κανονικούς πληθυσμούς δίνει συντελεστή γραμμικής συσχέτισης  $r=0.874$ . Υπάρχει ή όχι στατιστικά σημαντική διαμετρική συσχέτιση μεταξύ των μεταβλητών  $X$  και  $Y$  στον πληθυσμό σε ε.σ.σ.  $\alpha = 0.01$ ;

Λύση Έχουμε  $n=6, r=0.874$

Ελέγχουμε σε ε.σ.σ.  $\alpha = 0.01 = 1\%$  των

$$H_0: \rho = 0$$

vs

$$H_1: \rho > 0$$

Χρησιμοποιούμε τη σ.σ.

ελέγχου:

$$T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \stackrel{H_0}{\sim} t_{n-2} = t_4$$

$$\text{Όπως } T = \frac{0.874 \cdot 2}{\sqrt{1-0.874^2}} = \frac{1.748}{0.48} \approx 3.64$$

Απόφαση: η  $H_0$  απορρίπτεται σε ε.σ.σ.  $\alpha = 1\%$

$$\text{αν } T > t_{n-2, \alpha} = t_{4, 0.01} = 3.747$$

Η ανισότητα δεν ισχύει.

Συνεπώς η Ηο δεν απορρίπτεται. Άρα μια μη-  
στατιστικά σημαντική θετική συσχέτιση φαίνεται  
ότι υπάρχει μεταξύ των μεταβλητών Χ και Υ.

#3 Οι μηνιαίες πωλήσεις μπίρας (Υ) σε εκ. βαρέλια  
και η μέση μηνιαία θερμοκρασία (Χ) σε βαθμούς με-  
λσίου καταγράφηκαν από το Διευθυντή πωλήσεων  
ενός μεγάλου εργοστασίου για το χρονικό διάστημα  
Δεκ-Σεπτ. και τα παρακάτω δεδομένα συγκεντρώ-  
θηκαν:

$$\sum x_i = 115, \quad \sum y_i = 32,7, \quad \sum x_i^2 = 1585, \quad \sum y_i^2 = 109,07$$

$$\sum x_i y_i = 395,4$$

Να ελεγχθεί στατιστικά αν ισχύει ο ισχυρισμός ότι  
καθώς αυξάνει η θερμοκρασία αυξάνει και η κατα-  
νάλωση μπίρας σε επίπεδο σημαντικότητας  $\alpha = 0.05$ .

Λύση Θα ελέγξουμε τις υποθέσεις:

$H_0: \rho = 0$  (δεν υπάρχει γραμμική συσχέτιση μεταξύ  
της θερμοκρασίας και της κατανάλωσης  
μπίρας)

$H_1: \rho > 0$  (υπάρχει θετική γραμμική συσχέτιση  
μεταξύ της θερμοκρασίας και της  
κατανάλωσης μπίρας).

Έχουμε  $n = 10$ . Από τα παραπάνω δεδομένα, υπολο-  
γίζουμε το συντελεστή γραμμικής συσχέτισης  $r$  του  
δείγματος.



$$r = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sqrt{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \sqrt{n \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i\right)^2}}$$

$$= \frac{10 \cdot (395,4) - 115 \cdot (32,27)}{\sqrt{10(1585) - 115^2} \sqrt{10(109,07) - 32,7^2}} = 0,816$$

Απορρίπτει η  $H_0$  αν  $T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} > t_{n-2, \alpha}$   
για ε.σ.σ.  $\alpha$

Οπως  $T = \frac{0,816 \sqrt{8}}{\sqrt{1-0,816^2}} = 3,9924$

Από τους πίνακες της  $t_8$   
 $t_{8,0.05} = 1,86$

Παρασπούμε ότι  $T > t_{8,0.05}$  άρα τα δεδομένα παρέχουν ισχυρές ενδείξεις ότι οι δύο μεταβλητές  $X$  και  $Y$  συσχετίζονται θετικά σε ε.σ.σ.  $\alpha = 5\%$ .  
Άρα, ισχύει ο ισχυρισμός σε  $\alpha = 5\%$  ότι καθώς αυξάνει η θερμοκρασία αυξάνει και η κατανάλωση μπίρας.

