

Διάλεξη 3^η

Παλινδρόμηση και Ανάλυση Διακύμανσης 22/10/2019

Συμπερασματολογία (τεστ και Διαστήματα εμπιστοσύνης)

ους παραμέτρους β_0 και β_1 της α.χ.π $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

$i = 1, 2, \dots, n$

Έχει νόημα να ενδιαφερόμαστε για έλεγχο υποθέσεων ουσ παραμέτρων β_0 και β_1 ;

πχ Έχει νόημα να ενδιαφερώ για έλεγχο της $H_0: \beta_1 = \beta_1^*$, β_1^* γινωσ?
 ΝΑΙ

Έχει νόημα να κατασκευάσω τεστ για έλεγχο $H_0: \beta_1 = \beta_1^*$ γιατί μου δίνει τη δυνατότητα να ελέγξω συγκεκριμένη μεταβλητή της Y για μοναδιαία μεταβολή της X .

Ακόμη, αν $\beta_1^* = 0$ τότε η αποδοχή της $H_0: \beta_1 = 0$ σημαίνει ότι δεν υφίσταται γραμμική σχέση μεταξύ των X και Y .

Πώς θα κατασκευάσω στατιστικό τεστ για έλεγχο της $H_0: \beta_1 = \beta_1^*$ ή $H_0: \beta_0 = \beta_0^*$?

Τεχνική κατασκευής τεστ του Wald
ή κατασκευάζουμε τεστ σημαντικότητας.

Έστω τυχαίο δείγμα W_1, \dots, W_n από $N(\mu, \sigma^2)$

Έστω ο έλεγχος $H_0: \mu = \mu_0$ έναντι $H_1: \mu \neq \mu_0$.

Ⓘ σ^2 γνωστή (Z-TEST)

Το στατιστικό τεστ για τον έλεγχο της $H_0: \mu = \mu_0$ έχει ΣΣΤ την $Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$ με κατανομή $N(0,1)$

υπό την H_0 και κ.π $|Z| \geq Z_{\alpha/2}$

Ⓜ σ^2 άγνωστο (t-TEST)

Το στατιστικό τεστ για έλεγχο την $H_0: \mu = \mu_0$ έχει ΣΣΤ την $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$ με κατανομή t_{n-1} υπό την H_0 και κ.π $|t| \geq t_{n-1, \alpha/2}$.

* Παρατηρώντας την Z και t η ΣΣΤ έχει τη μορφή:

Εκτιμητής της	Η παράμετρος
παραμέτρου που	που εμφανίζεται
εμφανίζεται στην H_0	στην H_0

• ΣΣΤ = $\frac{\text{Τυπική Απόκλιση του Εκτιμητή}}$

- Χρειάζεται την κατανομή της ΣΣΤ υπό την H_0
- Χρειάζεται η μορφή της κρίσιμης περιοχής (περιοχή απόρριψης της H_0)
- Υπολογισμός του κρίσιμου σημείου $\alpha = P(\text{Απορρ } H_0 / H_0 \text{ αληθής})$.

Υποθέτω ότι οι υποθέσεις για τα σφάλματα είναι ισχύουν
Κατασκευή τεστ για έλεγχο της $H_0: \beta_1 = \beta_1^*$ (β_1^* γνωστό)
έναντι $H_a: \beta_1 \neq \beta_1^*$.

Ξεκινώ από εκτίμησή της β_1 : $\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right)$

Υπό την $H_0: \beta_1 = \beta_1^*$, $\hat{\beta}_1 \sim N\left(\beta_1^*, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right)$

$$\frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\frac{\sigma^2}{\sum (x_i - \bar{x})^2}}} \sim N(0,1) \text{ υπό την } H_0$$

ΠΡΟΒΛΗΜΑ

Το σ^2 είναι
 άγνωστο

Ισχύει ότι $\frac{SS_{res}}{\sigma^2} \sim \chi^2_{n-2}$

Θεωρώ το:

επειδή $\hat{\beta}_1$ ανεξ SS_{res} .

$$t_{n-2} \equiv \frac{N(0,1) \text{ υπό } H_0}{\sqrt{\frac{\chi^2_{n-2}}{n-2}}} \sim t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\frac{SS_{res}}{\sigma^2} / n-2}}$$

$$\Rightarrow t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\frac{SS_{res}/(n-2)}{\sum (x_i - \bar{x})^2}}} = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\frac{MS_{res}}{\sum (x_i - \bar{x})^2}}}$$

ή

$$t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_1)}} \sim t_{n-2} \text{ υπό } H_0$$

Η ΣΣΤ είναι $t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\hat{\text{Var}}(\hat{\beta}_1)}} \sim t_{n-2}$ υπό H_0

$$\text{και } \hat{\text{Var}}(\hat{\beta}_1) = \frac{MS_{\text{res}}}{\sum (x_i - \bar{x})^2}$$

Μικρές τιμές του $t \Rightarrow$ Δέχονται την H_0

Για να απορρ την H_0 θα πάρω μεγάλες τιμές του t .

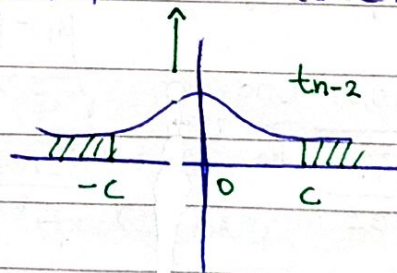
Μορφή κρίσιμης περιοχής. Μεγάλες τιμές της t .

Επειδή το t παίρνει και αρνητικές τιμές

$$\text{μορφή της κ.π είναι } |t| \geq c \quad (\Leftrightarrow \begin{matrix} t \geq c \\ \text{ή} \\ t \leq -c \end{matrix})$$

Υπολογισμός κρίσιμου σημείου c

$$\begin{aligned} \text{Πάντα } \alpha &= P(\text{απορ } H_0 / H_0 \text{ αληθής}) = P(|t| \geq c / H_0 \text{ αληθής}) \\ &= P(t \geq c \text{ ή } t \leq -c / H_0 \text{ αληθής}) = 2P(t \geq c / H_0 \text{ αληθής}) \end{aligned}$$



$$\begin{aligned} &= 2P(t \geq c / t \sim t_{n-2}) = 2P(t_{n-2} \geq c) \\ &\Rightarrow P(t_{n-2} \geq c) = \alpha/2 \end{aligned}$$

$$\text{αρα } c = t_{n-2, \alpha/2}$$

Για τον έλεγχο της $H_0: \beta_1 = \beta_1^* \vee H_1: \beta_1 \neq \beta_1^*$ η ΣΣΤ είναι:

$$t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\hat{\text{Var}}(\hat{\beta}_1)}} \text{ με κατανομή } t_{n-2} \text{ υπό την } H_0$$

και $kn |t| \geq t_{n-2, \alpha/2}$

$$\text{όπου } \widehat{\text{Var}}(\hat{\beta}_1) = \frac{MS_{\text{res}}}{\sum (x_i - \bar{x})^2}$$

Αντίστοιχα για την β_0 .

$$\text{όμως } \widehat{\text{Var}}(\hat{\beta}_0) = \frac{\sum x_i^2}{n \sum (x_i - \bar{x})^2} \cdot MS_{\text{res}} \quad \left(\begin{array}{l} \text{Η μόνη διαφορά σε} \\ \text{σχέση με το } \beta_1 \end{array} \right)$$

Κατασκευή Διαστήματος εμπιστοσύμης (Δ.Ε) με βαθμό εμπιστοσύμης (β.ε) $100(1-\alpha)\%$ για την β_1 .

Διάστημα εμπιστοσύμης

Έστω τ.δ W_1, \dots, W_n από πληθυσμό με κατανομή

$$f(x, \theta) \quad (\text{π.χ } \text{εκθ } f(x, \theta) = \theta e^{-\theta x}, x > 0)$$

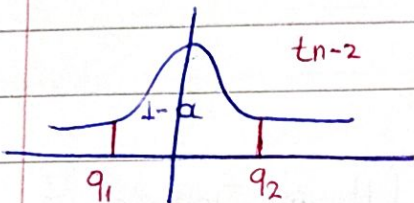
(*) Ορισμός: Ονομάζουμε Δ.Ε για την παράμετρο θ με β.ε $100(1-\alpha)\%$ κάθε διάστημα (L, U) όπου $L = L(W_1, \dots, W_n)$ και $U = U(W_1, \dots, W_n)$ τ.ω $P(L < \theta < U) = 1 - \alpha$. **ΜΕΘΟΔΟΣ ΑΝΤΙΣΤΡΕΠΤΗΣ ΠΡΟΣΟΧΤΑΣ.**

$$\text{π.χ } \Delta.Ε \text{ για το } \mu \sim N(\mu, \sigma^2) \begin{cases} \rightarrow \left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right) \quad \underline{\sigma^2 \text{ γνωστό}} \\ \rightarrow \left(\bar{x} - t_{n-1, \frac{\alpha}{2}} \frac{s}{\sqrt{n}}, \bar{x} + t_{n-1, \frac{\alpha}{2}} \frac{s}{\sqrt{n}} \right) \quad \underline{\sigma^2 \text{ άγνωστο}} \end{cases}$$

Βασίζονται στον $\hat{\beta}_i$ και ειδικότερα στην $\frac{\hat{\beta}_i - \beta_i}{\sqrt{\hat{\text{Var}}(\hat{\beta}_i)}} \sim t_{n-2}$

είναι αντιστρέφτη ποσότητα.

Αφού η t έχει t_{n-2} , υπάρχουν q_1, q_2 με $q_1 < q_2$ τ.ω
 $1-\alpha = P(q_1 < t < q_2) = P(q_1 < \frac{\hat{\beta}_i - \beta_i}{\sqrt{\hat{\text{Var}}(\hat{\beta}_i)}} < q_2)$



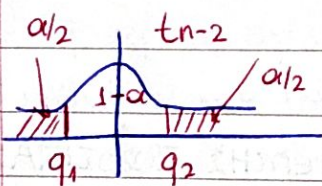
λύω προς β_i προσεγγιστικά.

$$= P(\hat{\beta}_i - q_2 \sqrt{\hat{\text{Var}}(\hat{\beta}_i)} < \beta_i < \hat{\beta}_i - q_1 \sqrt{\hat{\text{Var}}(\hat{\beta}_i)})$$

Από τον ορισμό του Δ.Ε. (*) ένα Δ.Ε για την β_i με β.ε $\pm 100(1-\alpha)\%$ είναι το $(\hat{\beta}_i - q_2 \sqrt{\hat{\text{Var}}(\hat{\beta}_i)}, \hat{\beta}_i - q_1 \sqrt{\hat{\text{Var}}(\hat{\beta}_i)})$

Επιλέγω τα q_1 και q_2 ώστε να κόβουν ίσες ουρές από την t_{n-2} κατανομή.

$$\text{δηλ } P(t_{n-2} > q_2) = P(t_{n-2} < q_1) = \frac{\alpha}{2}$$



$$\text{Από την } P(t_{n-2} > q_2) = \frac{\alpha}{2} \Rightarrow q_2 = t_{n-2, \alpha/2}$$

$$\text{Από την } P(t_{n-2} < q_1) = \frac{\alpha}{2}$$

$$\Rightarrow 1 - P(t_{n-2} > q_1) = 1 - \frac{\alpha}{2}$$

$$\Rightarrow q_1 = t_{n-2, 1-\alpha/2}$$

Συμμετρως το $\pm 100(1-\alpha)\%$ Δ.Ε για β_i είναι

$$(\hat{\beta}_i - t_{n-2, \alpha/2} \sqrt{\hat{\text{Var}}(\hat{\beta}_i)}, \hat{\beta}_i - t_{n-2, 1-\alpha/2} \sqrt{\hat{\text{Var}}(\hat{\beta}_i)})$$

$$\text{με } \hat{\text{Var}}(\hat{\beta}_i) = \text{MSres} / \sum (x_i - \bar{x})^2$$

Το $100(1-\alpha)\%$ ΔΕ για το β_0 είναι:

$$\left(\hat{\beta}_0 - t_{n-2, \frac{\alpha}{2}} \sqrt{\hat{\text{Var}}(\hat{\beta}_0)}, \hat{\beta}_0 - t_{n-2, 1-\frac{\alpha}{2}} \sqrt{\hat{\text{Var}}(\hat{\beta}_0)} \right)$$

$$\mu\epsilon \hat{\text{Var}}(\hat{\beta}_0) = \frac{\sum x_i^2}{n \sum (x_i - \bar{x})^2} \text{MSres}$$

✓ εκτίμηση - προσέγγιση του χ_k .

Έστω η πρόβλεψη $\hat{Y}_k = \hat{\beta}_0 + \hat{\beta}_1 X_k$

Ποια η $\text{Var}(\hat{Y}_k)$?

$$\hat{Y}_k = \hat{\beta}_0 + \hat{\beta}_1 X_k$$

$$= \bar{Y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 X_k = \bar{Y} + \hat{\beta}_1 (X_k - \bar{X})$$

$$\text{Var}(\hat{Y}_k) = \text{Var} \left[\bar{Y} + \hat{\beta}_1 (X_k - \bar{X}) \right]$$

συνάρτηση
ποσότητας

Γνωρίζω ότι:

$$\text{Var} \left(\sum_{i=1}^n a_i W_i \right) = \sum_{i=1}^n a_i^2 \text{Var}(W_i) + \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n a_i a_j \text{Cov}(W_i, W_j)$$

Αν W_i ασυσχετιση τότε

$$\text{Var} \left(\sum_{i=1}^n a_i W_i \right) = \sum_{i=1}^n a_i^2 \text{Var}(W_i)$$

$$\text{Άρα } \text{Var}(\hat{Y}_k) = \text{Var}(\bar{Y}) + (X_k - \bar{X})^2 \text{Var}(\hat{\beta}_1) + (X_k - \bar{X}) \text{Cov}(\bar{Y}, \hat{\beta}_1) \quad (1)$$

$$\text{όπου } \text{Cov}(\bar{Y}, \hat{\beta}_1) = \text{Cov} \left(\frac{1}{n} \sum Y_i, \sum \frac{X_i - \bar{X}}{\sum (X_i - \bar{X})^2} Y_i \right)$$

Αν W_1, \dots, W_n ασυσχ. τότε $\text{Cov}(\sum a_i W_i, \sum \beta_i W_i)$

$$= \sum_{i=1}^n a_i \beta_i \text{Var}(W_i)$$

$$\begin{aligned}
 \left. \begin{aligned}
 \text{Cov}(\bar{Y}, \hat{\beta}_1) &= \sum_{i=1}^n \frac{1}{n} \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})} \text{Var}(y_i) \\
 &= \frac{1}{n} \sum_{i=1}^n \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2} \sigma^2 \\
 &= \frac{\sigma^2}{n} \frac{1}{\sum (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \\
 &= 0 \quad (2)
 \end{aligned} \right\}
 \end{aligned}$$

Από (1) ή (2)

$$\text{Var}(\hat{Y}_k) = \text{Var}(\bar{Y}) + (x_k - \bar{x})^2 \text{Var}(\hat{\beta}_1)$$

• Η διακύμανση μικραίνει όταν x_k πολύ κοντά στο \bar{x} .

ΠΑΡΑΤΗΡΗΣΗ.

Όταν χρησιμοποιούμε το μοντέλο της αληθιάς για προβλέψεις πρέπει να το χρησιμοποιούμε για τιμές της ανεξάρτητης μεταβλητής X κοντά στο \bar{x} γιατί τότε $\text{Var}(\hat{Y}_k)$ γίνεται μικρή και άρα η πρόβλεψη είναι αξιόπιστη

• Το F-τεστ για την Παλινδρόμηση ή Το F-τεστ για τον έλεγχο της $H_0: \beta_1 = 0$.

Γνωρίζουμε ότι $E(MS_{res}) = \sigma^2$

Να βρούμε $E(MS_{reg}) = ?$

$$MS_{reg} = \frac{SS_{reg}}{1} = SS_{reg} = \hat{\beta}_1^2 \sum (x_i - \bar{x})^2$$

$$E(MS_{reg}) = E(\hat{\beta}_1^2 \sum (x_i - \bar{x})^2) = \sum (x_i - \bar{x})^2 E(\hat{\beta}_1^2)$$

$$\underline{\underline{\text{Var}(w) = E(w^2) - (E(w))^2}} \quad \sum (x_i - \bar{x})^2 \left[\text{Var}(\hat{\beta}_1) + (E\hat{\beta}_1)^2 \right]$$

$$= \sum (x_i - \bar{x})^2 \left[\frac{\sigma^2}{\sum (x_i - \bar{x})^2} + \beta_1^2 \right] = \sigma^2 + \beta_1^2 \sum (x_i - \bar{x})^2$$

Συμμεντρωτικά : $E(MS_{res}) = \sigma^2$

$$E(MS_{reg}) = \sigma^2 + \beta_1^2 \sum (x_i - \bar{x})^2$$

Παρατηρώ ότι υπό την $H_0: \beta_1 = 0$

$$E(MS_{res}) = E(MS_{reg}) = \sigma^2.$$

Αρα αν $\beta_1 = 0$ τότε $MS_{res} \approx MS_{reg}$
 (Αν $A \Rightarrow B$ τότε $\sim B \Rightarrow \sim A$ (\sim : αίρνηση))

Αν MS_{res} είναι πολύ διαφορετικό από το MS_{reg}
 τότε $\beta_1 \neq 0$, δηλ. τότε απορρ. $H_0: \beta_1 = 0$.

ΑΡΑ ένα στατιστικό τεστ για τον έλεγχο της $H_0: \beta_1 = 0$

μπορεί να βασιστεί στη σύγκριση του MS_{res} με το MS_{reg} .

Ειδικότερα : Αν MS_{res} πολύ διαφορετικό από MS_{reg}
 τότε απορρ. την $H_0: \beta_1 = 0$.

Αποδείξαμε $\frac{SS_{res}}{\sigma^2} \sim \chi^2_{n-2}$

Μήπως και $\frac{SS_{reg}}{\sigma^2} \sim \chi^2_{??}$

Είναι $SS_{reg} = \hat{\beta}_1^2 \sum (x_i - \bar{x})^2$

$$\text{Ισχύει } \hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum(x_i - \bar{x})^2}\right)$$

$$\text{και υπό την } H_0: \beta_1 = 0, \hat{\beta}_1 \sim N\left(0, \frac{\sigma^2}{\sum(x_i - \bar{x})^2}\right)$$

υπό την $H_0: \beta_1 = 0$

$$\frac{\frac{\hat{\beta}_1}{\sigma}}{\sqrt{\sum(x_i - \bar{x})^2}} \sim N(0, 1)$$

\Rightarrow υπό την $H_0: \beta_1 = 0$

$$\eta \frac{\sqrt{\sum(x_i - \bar{x})^2} \hat{\beta}_1}{\sigma} \sim N(0, 1) \rightarrow \frac{\sum(x_i - \bar{x})^2 \hat{\beta}_1^2}{\sigma^2} \sim N(0, 1)^2 \equiv \chi_1^2$$

$$\Rightarrow \text{υπό την } H_0: \beta_1 = 0 \text{ το } \frac{SS_{\text{reg}}}{\sigma^2} \sim \chi_1^2$$

Θεωρώ της ποσότητα

$$F = \frac{MS_{\text{reg}}}{MS_{\text{res}}} = \frac{SS_{\text{reg}}/1}{SS_{\text{res}}/(n-2)} \quad \text{Διαχωρίζω με } \sigma^2 \text{ αριθμικαί παρανομή.}$$

\Rightarrow

$$F = \frac{SS_{\text{reg}}/\sigma^2}{SS_{\text{res}}/\sigma^2(n-2)} = \frac{\chi_1^2/1}{\chi_{n-2}^2/(n-2)} \equiv F_{1, n-2}$$

Συμπερασματικά: Υπό την $H_0: \beta_1 = 0$ το F -πληθυσ

$$F = \frac{MS_{\text{reg}}}{MS_{\text{res}}} \sim F_{1, n-2}$$

• Μορφή και Μεγάλες τιμές του F (γιατί τότε το MS_{reg} πολύ διαφορετικό MS_{res}).

δηλ $F \geq c$.

ο Υπολογισμός κ.σ. C:

$$\begin{aligned}\alpha &= P(\text{απορ } H_0 / H_0 \text{ αληθής}) \\ &= P(F > c / F \sim F_{1, n-2}) = P(F_{1, n-2} > c)\end{aligned}$$

ορισμός
Εκατοσναιών
σημείων
της $F_{1, n-2}$

$C = F_{1, n-2, \alpha}$

Πα τον έλεγχο της $H_0: \beta_1 = 0$
η ΣΣΤ είναι:

$$F = \frac{MS_{\text{reg}}}{MS_{\text{res}}} \text{ με κατανομή } F_{1, n-2}$$

Υπό την $H_0: \beta_1 = 0$ και κ.π. $F > F_{1, n-2, \alpha}$.

ΠΑΡΑΤΗΡΗΣΗ.

Το F τεστ για τον έλεγχο της $H_0: \beta_1 = 0$ είναι ισοδύναμο με το t-τεστ διότι:

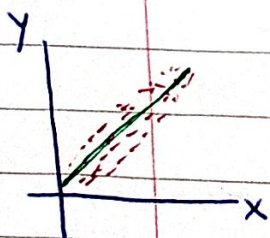
i) $F = t^2$

ii) $F_{1, n-2, \alpha} = t^2_{n-1, \alpha/2}$

Συντελεστής:

Συσχέτισης - Pearson.

Διάγραμμα Διασποράς: (δ, δ) το Μαθηματικό ανάλογο \Rightarrow του δ, δ



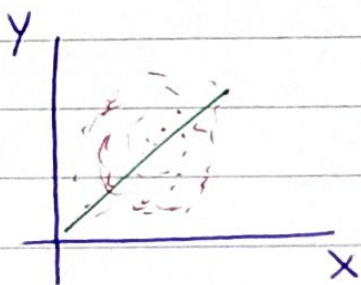
$$r_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Ιδιότητες του $r_{x,y}$.

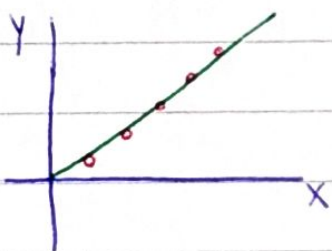
① $r_{x,y}$ καθαρός αριθμός

② $-1 \leq r_{x,y} \leq 1$

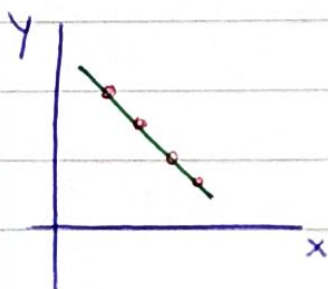
Ιδιαίτερο ενδιαφέρον παρουσιάζουν οι τιμές $0, -1, 1$.



$r_{x,y} = 0$ ← αντιστοιχεί στην μη ύπαρξη γραμμικής σχέσης.



$r_{x,y} = 1$ ← πλήρης θετική γραμμική σχέση.



$r_{x,y} = -1$ ← απόλυτη αρνητική γραμμική σχέση

ΠΑΡΑΤΗΡΗΣΗ

Ο συντελεστής Pearson εκφράζει μόνο γραμμική σχέση.

δηλ το $r_{x,y} = 0$ σημαίνει ότι ~~∃~~ γραμμική σχέση

Μπορεί όμως να ~~∃~~ άλλου είδους σχέση.

ρ_X :	X	-3	-1	0	1	β
	Y	9.54	9.95	10	9.95	9.54

$r_{x,y} = 0$ Άρα ~~∃~~ γραμμική σχέση μεταξύ X, Y
 όμως $X^2 + Y^2 = 100$ (το λαιτάνι)

