

Στατιστική συμπερασματολογία (ΤΕΣΤ, δ.ε.) για τις β_0, β_1 , παραμέτρους του μοντέλου της α.χ.π.
 θεωρώ το μοντέλο της α.χ.π. $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$, $i=1, \dots, n$ και υποθέτω ότι οι υποθέσεις για
 τα σφάλματα ικανοποιούνται.

Αναγκαιότητα ΤΕΣΤ για β_0, β_1

Έχει νοήμα $H_0: \beta_1 = \beta_1^*$ (β_1^* γνωστό) jii
 $H_0: \beta_1 = 0$

ΝΑΙ γιατί $H_0: \beta_1 = \beta_1^*$ ελέγχει αν η μεταβλητότητα της Y για μοναδιαία
 μεταβολή της X είναι ίση με β_1^*
 ή γιατί αν ισχύει η $H_0: \beta_1 = 0$ τότε \nexists σχέση μεταξύ Y και X .

ΤΕΣΤ για έλεγχο $H_0: \beta_1 = \beta_1^*$ (β_1^* γνωστό) έναντι $H_a: \beta_1 \neq \beta_1^*$

ΤΕΣΤ κατά Wald: Έστω τ.δ. W_1, \dots, W_n από $N(\mu, \sigma^2)$, $H_0: \mu = \mu_0$ (μ_0 γνωστό)
 $H_a: \mu \neq \mu_0$
 σ^2 γνωστό $\rightarrow Z$ ΤΕΣΤ, $Z = \frac{\bar{W} - \mu_0}{\sigma/\sqrt{n}}$

σ^2 άγνωστο $\rightarrow t$ -test, $t = \frac{\bar{W} - \mu_0}{S/\sqrt{n}}$, $S^2 = \frac{1}{n-1} \sum_{i=1}^n (W_i - \bar{W})^2$

$\Sigma\Sigma.T = \frac{\text{Ευτιμητής} - E(\text{Ευτιμητή})}{\text{Τυπική απόκλιση (ευτιμητή)}}$ (τουλάχιστον για μεγάλα δείγματα)

ⓐ Κατανομή της $\Sigma\Sigma.T.$ υπό την H_0

ⓑ Κρίσιμο σημείο $\leftarrow \alpha = P(\text{απορ. } H_0 \mid H_0 \text{ αληθεία}) = P(\text{σφάλμα τύπου I})$

Αφού δώσω ΤΕΣΤ για β_1 στηρίζομαι σε ένα ευτιμητή της β_1 , την $\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum (X_i - \bar{X})^2}\right)$

Κατασκευή στατιστικού τεστ:

Υπό την $H_0: \beta_1 = \beta_1^*$, $\hat{\beta}_1 \sim N\left(\beta_1^*, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right)$ έναντι $H_a: \beta_1 \neq \beta_1^*$ (β_1^* γνωστό)
 (Αναλόγως για την $H_0: \beta_0 = \beta_0^*$)

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

$t_v \approx \frac{N(0, 1)}{\sqrt{\chi_v^2/v}}$ οπότε, εφόσον $\frac{SS_{res}}{\sigma^2} \sim \chi_{n-2}^2$ θα έχουμε:

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

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

Το $t = \frac{N(0, 1)}{\sqrt{\frac{SS_{res}}{\sigma^2} / (n-2)}} = \frac{N(0, 1)}{\sqrt{\chi_{n-2}^2 / (n-2)}}$ ανεξ. t_{n-2} $\hat{\beta}_1$ ανεξάρτητο από SS_{res}

Συνοψίζοντας μέχρι εδώ η ΣΣΤ:

Υπό την $H_0: \beta_1 = \beta_1^*$ η $t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\frac{MS_{res}}{\sum (x_i - \bar{x})^2}}} \sim t_{n-2}$ ή $t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\text{Var}(\hat{\beta}_1)}}$ καθώς $\sigma^2 = MS_{res}$

Κριτική περιοχή:

Αν το t έχει μεγάλες τιμές τότε $\hat{\beta}_1$ πολύ διαφορετικό από β_1^* , τότε και το β_1 πολύ διαφορετικό από το β_1^* , άρα απορ. H_0 .

Τέτοιες μεγάλες τιμές του t οδηγούν σε απόρριψη της $H_0: \beta_1 = \beta_1^*$

Άρα η κρίσιμη περιοχή αποτελείται από μεγάλες τιμές του t , δηλ. $|t| \geq c$

Άρα, μορφή κ.π. $|t| \geq c$.

♦ Άρκει να προσδιορίσω το κρίσιμο σημείο c .

Πάντα από $P(\text{σφάλμα τύπου I}) = P(\text{αναρ. } H_0 | \text{ } H_0 \text{ αληθεύει}) = \alpha = P(|t| \geq c | t \sim t_{n-2})$

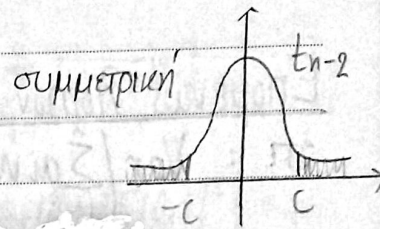
$$= P(|t_{n-2}| \geq c) = P(t_{n-2} \geq c \text{ ή } t_{n-2} \leq -c) = 2P(t_{n-2} \geq c) \quad \leftarrow C = t_{n-2, \frac{\alpha}{2}}$$

Παρατηρώντας ότι $\sqrt{\frac{MS_{res}}{\sum (x_i - \bar{x})^2}} = \sqrt{\frac{\hat{\beta}^2}{\sum (x_i - \bar{x})^2}} = \sqrt{\text{Var}(\hat{\beta}_1)}$ \hookrightarrow ένας ευρηματίας

Οπότε: για τον έλεγχο της $H_0: \beta_1 = \beta_1^*$ έναντι $H_a: \beta_1 \neq \beta_1^*$ η ΣΣΤ είναι $C = t_{n-2, \frac{\alpha}{2}}$

$t = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\text{Var}(\hat{\beta}_1)}}$ με κατανομή t_{n-2} υπό H_0 και κ.π. μεγέθους α , την $|t| \geq t_{n-2, \frac{\alpha}{2}}$

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



Κατασκευή δ.ε. για την παράμετρο β_1 :

Επειδή $\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} \sim t_{n-2}$ κατασκευή δ.ε. για την παράμετρο β_1 είναι αντιστρέφτη, δηλ. q_1, q_2 ώστε:

$$1 - \alpha = P(q_1 < t < q_2) = P\left(q_1 < \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} < q_2\right) = P\left(\hat{\beta}_1 - q_2 \sqrt{\text{Var}(\hat{\beta}_1)} \leq \beta_1 \leq \hat{\beta}_1 - q_1 \sqrt{\text{Var}(\hat{\beta}_1)}\right)$$

- Αντιστρέφτη ποσότητα

 - ① Συναρτηση της παραμέτρου
 - ② Συναρτηση των δεδομένων
 - ③ Η κατανομή της είναι ανεξάρτητη από την παράμετρο.

Άρα, $100(1-\alpha)\%$ δ.ε. για την β_1 είναι $(\hat{\beta}_1 - q_2 \sqrt{\text{Var}(\hat{\beta}_1)}, \hat{\beta}_1 - q_1 \sqrt{\text{Var}(\hat{\beta}_1)})$

Το δ.ε. ίσων ουρών προκύπτει για $q_2 = -q_1$ και $q_1 = -t_{n-2, \frac{\alpha}{2}}$ $\left. \begin{matrix} q_2 = t_{n-2, \frac{\alpha}{2}} \end{matrix} \right\}$ λόγω συμμετρίας

Οπότε Το $100(1-\alpha)\%$ δ.ε. για β_1 είναι $(\hat{\beta}_1 - t_{n-2, \frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\beta}_1)}, \hat{\beta}_1 + t_{n-2, \frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\beta}_1)})$

Ανάλογα Για τον έλεγχο $H_0: \beta_0 = \beta_0^*$ (β_0^* γνωστό) έναντι $H_a: \beta_0 \neq \beta_0^*$ η ΣΣΤ είναι

$t = \frac{\hat{\beta}_0 - \beta_0^*}{\sqrt{\text{Var}(\hat{\beta}_0)}}$, με κατανομή t_{n-2} , υπό H_0 και κ.π. μεγέθους α την

$|t| \geq t_{n-2, \frac{\alpha}{2}}$ όπου $\text{Var}(\hat{\beta}_0) = \frac{\sum (x_i)^2}{n \sum (x_i - \bar{x})^2} MS_{res}$ και $100(1-\alpha)\%$

είναι το $(\hat{\beta}_0 - t_{n-2, \frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\beta}_0)}, \hat{\beta}_0 + t_{n-2, \frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\beta}_0)})$

Διακύμανση της πρόβλεψης

Το μοντέλο της α.χ.π. χρησιμοποιείται για προβλέψεις.

Ειδικότερα το ευσταθές μοντέλο χρησιμοποιείται για προβλέψεις.

Έτσι η πρόβλεψη της Y όταν $X = X_k$ είναι $\hat{Y}_k = \hat{\beta}_0 + \hat{\beta}_1 X_k$

ΕΡΩΤΗΜΑ:

Όμως το θέμα είναι, πόσο καλή είναι η πρόβλεψη;

Η πρόβλεψη \hat{Y}_k είναι καλή όταν $\text{Var}(\hat{Y}_k)$ είναι μικρή.

► Υπολογισμός $\text{Var}(\hat{Y}_k)$

$$\hat{Y}_k = \hat{\beta}_0 + \hat{\beta}_1 X_k = \bar{Y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 X_k = \bar{Y} - \hat{\beta}_1 (\bar{X} - X_k)$$

Επομένως, λοιπόν: $\text{Var}(\hat{Y}_k) = \text{Var}(\bar{Y} - \hat{\beta}_1 (\bar{X} - X_k))$ για να το υπολογίσω θυμάμαι

$$\text{Var}\left(\sum_{i=1}^n a_i w_i\right) = \sum_{i=1}^n a_i^2 \text{Var}(w_i) + \sum_{i \neq j} a_i a_j \text{Cov}(w_i, w_j) \text{ για } i, j = 1, \dots, n$$

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

Ξέρω ότι η συνδιακύμανση είναι: $\text{Cov}(\bar{Y}, \hat{\beta}_1) = \text{Cov}\left(\frac{1}{n} \sum Y_i, \frac{\sum (X_i - \bar{X}) Y_i}{\sum (X_i - \bar{X})^2}\right) =$

$$\left\{ \text{Cov}\left(\sum_{i=1}^n a_i w_i, \sum_{j=1}^n b_j z_j\right) = \sum_{i=1}^n \sum_{j=1}^n a_i b_j \text{Cov}(w_i, z_j) \right\} = \text{Cov}\left(\sum_{i=1}^n \frac{1}{n} Y_i, \frac{\sum (X_i - \bar{X}) Y_i}{\sum (X_i - \bar{X})^2}\right) \quad (2)$$

Αρα, τελικά, από (2), (3) η συνδιακύμανση είναι: $\begin{cases} w_i = z_i = Y_i \\ \text{Cov}(Y_i, Y_j) = 0, Y_i \text{ ανεξάρτητα } (i, j) \end{cases}$

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

$$= \sum_{i=1}^n \frac{1}{n} \frac{(X_i - \bar{X}) \sigma^2}{\sum (X_i - \bar{X})^2} = \frac{\sigma^2}{n} \frac{\sum (X_i - \bar{X})}{\sum (X_i - \bar{X})^2} = \frac{\sigma^2}{n} \frac{\sum (X_i - \bar{X})}{(\sum (X_i - \bar{X})^2)} \quad (4)$$

$$(\sum (X_i - \bar{X}) = \sum X_i - n\bar{X} = 0)$$

δηλ. από (1), (4)

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

Για να είναι η πρόβλεψη Y_k όσο πιο αξιόπιστη μπορεί πρέπει $\text{Var}(\hat{Y}_k)$ όσο δυνατόν ελάχιστο

Αρα, τι πρέπει να συμβαίνει για να είναι η $\text{Var}(\hat{Y}_k)$ όσο πιο μικρή γίνεται,

δηλ. αρκεί το X_k να είναι κοντά στο \bar{X} ή το $(X_k - \bar{X})^2$ όσο το δυνατόν ελάχιστο

Συμπέρασμα: Όταν το ευσταθές μοντέλο $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ χρησιμοποιείται για προβλέψεις θα πρέπει η πρόβλεψη να γίνεται για τιμή του X κοντά στο σύνολο των δεδομένων, για την X αν

στο \bar{X}

Το F-τεστ για έλεγχο παλινδρόμησης ή

το F-τεστ για έλεγχο της $H_0: \beta_1 = 0$ έναντι $H_a: \beta_1 \neq 0$.

Για τον έλεγχο αυτό κατασκευάσαμε t-τεστ ($H_0: \beta_1 = \beta_1^*$, β_1 γνωστό)

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

Ισχύει: $E(MS_{reg}) = \sigma^2 + \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$

Πράγματι, $E(MS_{reg}) = E\left[\frac{SS_{reg}}{1}\right] = E(SS_{reg}) = E\left[\hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2\right] =$
 $= \sum_{i=1}^n (x_i - \bar{x})^2 E(\hat{\beta}_1^2) = \sum_{i=1}^n (x_i - \bar{x})^2 [Var(\hat{\beta}_1) + \{E(\hat{\beta}_1)\}^2] =$
 $= \sum_{i=1}^n (x_i - \bar{x})^2 \left[\frac{\sigma^2}{\sum (x_i - \bar{x})^2} + \beta_1^2 \right] \Rightarrow$

$$\Rightarrow E[MS_{reg}] = \sigma^2 + \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

Υπό την $H_0: \beta_1 = 0$ τότε $E(MS_{res}) = E(MS_{reg})$ ή $MS_{res} \approx MS_{reg}$ } \Rightarrow
Προτασιακός Λογισμός: Αν $A \Rightarrow B$ τότε $\sim B \Rightarrow \sim A$

Αν MS_{res} πολύ διαφορετικό από MS_{reg} , τότε η $H_0: \beta_1 = 0$ δεν ισχύει, άρα απορρίπτεται.

Συμπερασματικά:

- Άρα, ένα τεστ για τον έλεγχο της $H_0: \beta_1 = 0$ μπορεί να στηριχτεί στη σύγκριση των MS_{reg} με το MS_{res} .

Γνωρίζοντας, ότι $\frac{SS_{res}}{\sigma^2} \sim \chi_{n-2}^2$, συγκρίνω τα MS_{res} και MS_{reg} με βάση το

πηλίκο τους: $\frac{MS_{reg}}{MS_{res}} = \frac{SS_{reg}/1}{SS_{res}/(n-2)}$ ή μπορώ να ελέγξω την διαφορά τους αν είναι κοντά στο μηδέν.

► Απομένει να βρω κατανομή για SS_{reg} .

Το $SS_{reg} \sim \chi_1^2$ υπό $H_0: \beta_1 = 0$ αν οι υποθέσεις για τα σφάλματα ικανοποιούνται

$$SS_{reg} = \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

Υπό τις υποθέσεις για σφάλματα: $\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right)$ υπό την $H_0: \beta_1 = 0$.

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

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

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

Θεωρώ τη στατιστική συνάρτηση: $F = \frac{MS_{reg}}{MS_{res}} = \frac{SS_{reg}/1}{SS_{res}/(n-2)} = \frac{(SS_{reg}/\sigma^2)/1}{(SS_{res}/\sigma^2)/(n-2)}$

$$F \equiv \frac{\chi_1^2/1}{\chi_{n-2}^2/n-2} \sim F_{1, n-2} \text{ υπό } H_0: \beta_1 = 0.$$

χ_1^2 και χ_{n-2}^2 ανεξάρτητες
 αν SS_{reg} ανεξάρτητο SS_{res} (SS_{reg} εξαρτάται από $\hat{\beta}_1$ το οποίο είναι ανεξάρτητο SS_{res})
 Ισχύει αλλά η απόδειξη ξεφεύγει από το μάθημα.

Συγκεντρωτικά:

► Για τον έλεγχο της $H_0: \beta_1 = 0$ κ.σ.σ. του τεστ είναι $F = \frac{MS_{reg}}{MS_{res}}$ με κατανομή

$F_{1, n-2}$ υπό H_0 και κ.π. της μορφής: μεγάλες τιμές του F , δηλ. $F \geq c \rightarrow$ υπόψιν αίτιο

Υπολογισμός κ.σ. c :

Μορφή της κρίσιμης περιοχής: Μεγάλες τιμές του F , δηλαδή $F \geq c$.

$$\alpha = P(\text{Απορ. } H_0 | H_0 \text{ αληθ.}) = P(F \geq c | F \sim F_{1, n-2}) = P(F_{1, n-2} \geq c) \Rightarrow c = F_{1, n-2, \alpha}$$

ΠΑΡΑΤΗΡΗΣΗ: Το F τεστ ισοδύναμο με το t -τεστ γιατί $F = \frac{MS_{reg}}{MS_{res}} = \frac{SS_{reg}}{SS_{res}} \rightarrow$
 $(t = \frac{\hat{\beta}_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} \sim t_{n-2})$

$$= \frac{\hat{\beta}_1^2 \sum (X_i - \bar{X})^2}{MS_{res}} = \left(\frac{\hat{\beta}_1}{\sqrt{\frac{MS_{res}}{\sum (X_i - \bar{X})^2}}} \right)^2 = \left(\frac{\hat{\beta}_1 - 0}{\sqrt{\widehat{Var}(\hat{\beta}_1)}} \right)^2 = t^2$$

Συντελεστής Συσχέτισης Pearson

Πληθυσμιακή μορφή: Αν X, Y είναι τ.μ. $\rightarrow \rho(X, Y) = \frac{Cov(X, Y)}{\sqrt{Var X} \sqrt{Var Y}}$

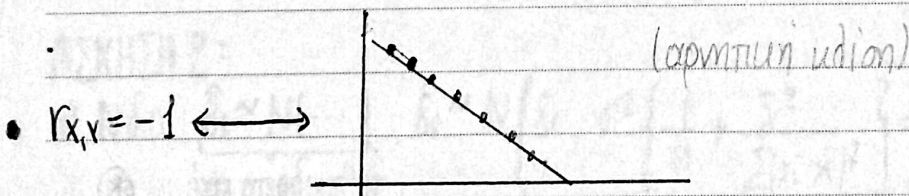
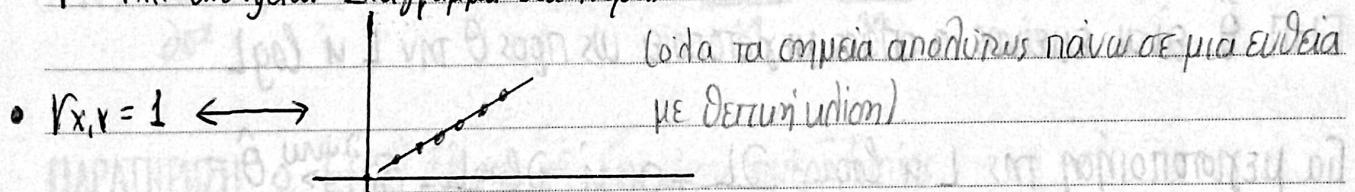
Δειγματική μορφή: Αν έχω τ.δ. X_1, \dots, X_n
 $Y_1, \dots, Y_n \rightarrow r_{X, Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \sqrt{\sum (Y_i - \bar{Y})^2}}$

Ο συντελεστής εμπίπτει του πληθυσμιακού συντελεστή συσχέτισης του Pearson και χρησιμοποιείται για ελέγχους της μορφής $H_0: \rho = 0$ ή $H_0: \rho = \rho_0$.

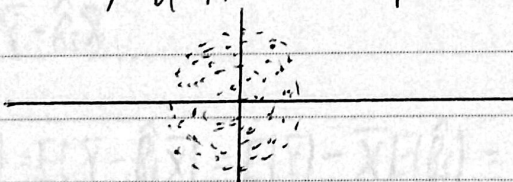
ΙΔΙΟΤΗΤΕΣ:

- 1) Είναι καθαρός αριθμός (απαλλαγμένος, δηλ. από μονάδες μέτρησης)
- 2) Συμμετρικός $r_{X, Y} = r_{Y, X}$
- 3) $-1 \leq r_{X, Y} \leq 1$ (από ανισότητα Cauchy-Schwarz)

Το $r_{X, Y}$ σχετίζεται: Διαγράμμα διασποράς.



• $r_{X, Y} = 0 \iff \nexists$ γραμμική σχέση μεταξύ των X και Y



Αν $r_{X, Y} = 0$ τότε \nexists γραμμική σχέση, μπορεί όμως να \exists άλλη είδους σχέση.