



ΓΕΩΠΟΝΙΚΟ ΠΑΝΕΠΙΣΤΗΜΙΟ ΑΘΗΝΩΝ
AGRICULTURAL UNIVERSITY OF ATHENS

Συσχέτιση - Συμμεταβολή

Κατσιλέρος Αναστάσιος

2017

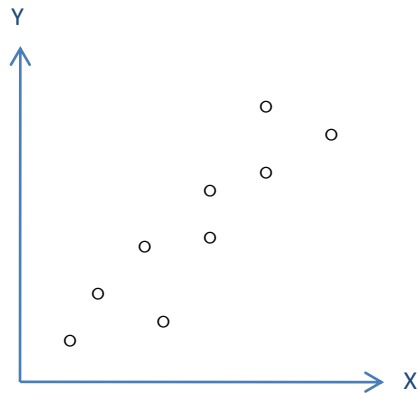
ΣΥΣΧΕΤΙΣΗ ΔΥΟ ΜΕΤΑΒΛΗΤΩΝ

Με την ανάλυση της συσχέτισης, μελετάται η ένταση (ή βαθμός) της σχέσης των δύο μεταβλητών ανεξάρτητα από την ύπαρξη σχέσης αιτίου-αποτελέσματος μεταξύ τους

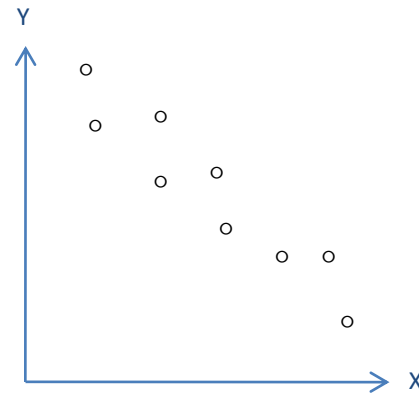
Για να εκτιμήσουμε την ένταση και την κατεύθυνση της σχέσης αυτής, υπολογίζουμε τον συντελεστή γραμμικής συσχέτισης Pearson's r .

Ο συντελεστής γραμμικής συσχέτισης του Pearson (r) κυμαίνεται από -1 έως 1 και είναι ανεξάρτητο από τις μονάδες μέτρησης. Ο βαθμός της σχέσης αυξάνεται όσο το r πλησιάζει την απόλυτη τιμή 1 ενώ τιμές κοντά στο μηδέν υποδεικνύει απουσία σχέσης μεταξύ των δύο μεταβλητών.

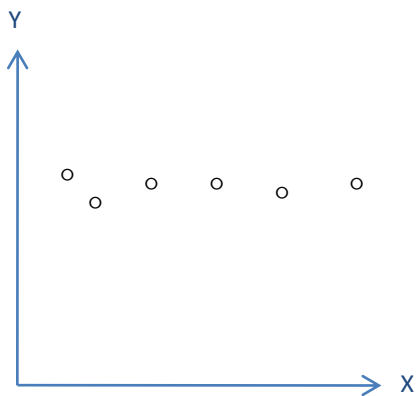
Πειραματικοί Σχεδιασμοί



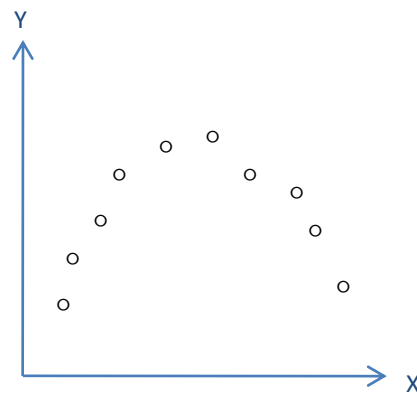
$r \approx 1$



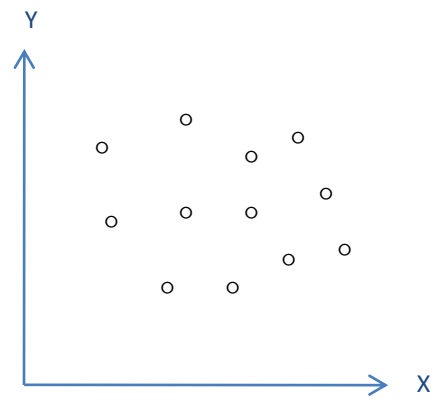
$r \approx -1$



$r = 0$



$r = 0$



$r = 0$

Προσδιορισμός της σχέσης δύο μεταβλητών:

$$r = \frac{S_{xy}}{S_x S_y}$$

Όπου:

$$s_{xy} = \text{Cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

οπότε

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Προσδιορισμός της σχέσης δύο μεταβλητών:

$$r = \frac{A\Gamma_{XY}}{\sqrt{AT_x AT_y}}$$

Όπου:

$$A\Gamma_{XY} = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \sum_{i=1}^n X_i Y_i - \frac{\sum_{i=1}^n (X_i) \sum_{i=1}^n (Y_i)}{n}$$

$$AT_x = \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n X_i^2 - \frac{(\sum_{i=1}^n X_i)^2}{n}$$

$$AT_y = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - \frac{(\sum_{i=1}^n Y_i)^2}{n}$$

Πειραματικοί Σχεδιασμοί

Παράδειγμα: Έλαια Χ και πρωτεΐνη Υ

Χ	Υ	ΧΥ
1	9	9
2	6	12
3	8	24
3	4	12
4	5	20
5	3	15
5	2	10
6	3	18
ΣΧ = 29	ΣΥ = 40	ΣΧΥ = 120
ΣΧ ² = 125	ΣΥ ² = 244	n = 8

$$A\Gamma_{XY} = \sum_{i=1}^n X_i Y_i - \frac{\sum_{i=1}^n (X_i) \sum_{i=1}^n (Y_i)}{n} = 120 - \frac{29 * 40}{8} = -25$$

$$A\Gamma_x = \sum_{i=1}^n X_i^2 - \frac{(\sum_{i=1}^n X_i)^2}{n} = 125 - \frac{29^2}{8} = 19,875$$

$$A\Gamma_y = \sum_{i=1}^n Y_i^2 - \frac{(\sum_{i=1}^n Y_i)^2}{n} = 244 - \frac{40^2}{8} = 44$$

$$r = \frac{A\Gamma_{XY}}{\sqrt{A\Gamma_x A\Gamma_y}} = \frac{-25}{\sqrt{19,875 * 44}} = -0,845$$

Έλεγχος σημαντικότητας του συντελεστή συσχέτισης r

Μηδενική υπόθεση $H_0: r = 0$ και εναλλακτική $H_1: r \neq 0$

A. Χρήση πινάκων ελέγχου σημαντικότητας του συντελεστή συσχέτισης r

Οι κρίσιμες τιμές των πινάκων εξαρτώνται σε μεγάλο βαθμό από τον αριθμό των παρατηρήσεων n , όσο το n αυξάνεται η κρίσιμη τιμή του πίνακα για το r μειώνεται και είναι πιθανό να απορριφθεί η μηδενική υπόθεση. Η κρίσιμη τιμή του σχετικού πίνακα για τον συντελεστή συσχέτισης r με $BE = n-2 = 8-2 = 6$ και για $\alpha = 0,05$ είναι 0,755.

Επειδή η τιμή που υπολογίσθηκε $- 0,845$ είναι μεγαλύτερη από την κρίσιμη τιμή του πίνακα, απορρίπτεται η μηδενική υπόθεση. Επομένως υπάρχει στατιστικά σημαντική αρνητική σχέση μεταξύ των μεταβλητών X και Y .

Β. Έλεγχος σημαντικότητας με την δοκιμασία του t

$$t_r = \frac{(r - 0)}{\sqrt{(1 - r^2)/(n - 2)}}$$

$$t_r = \frac{(0,845 - 0)}{\sqrt{(1 - 0,845^2)/(8 - 2)}} = 3,8769$$

Η κρίσιμη τιμή του πίνακα t με ΒΕ= n-2= 8-2= 6 και για α = 0,05 είναι 2,44. Επειδή η τιμή που υπολογίσθηκε είναι μεγαλύτερη από την κρίσιμη τιμή, απορρίπτεται η μηδενική υπόθεση.

Γ. Για μεγάλο αριθμό ζευγών παρατηρήσεων n ($n > 50$) χρησιμοποιείται η μετατροπή του z .

$$z = \frac{1}{2} \ln \left[\frac{1+r}{1-r} \right]$$

Η δοκιμασία είναι η εξής:

$$t_s = \frac{z - 0}{T\Sigma_z} = z\sqrt{n-3}$$

όπου το τυπικό σφάλμα:

$$T\Sigma_z = \frac{1}{\sqrt{n-3}}$$

και συγκρίνεται με την κρίσιμη τιμή των πινάκων του t για $BE = \infty$

Όρια εμπιστοσύνης συντελεστή συσχέτιση

Για να υπολογιστούν τα όρια εμπιστοσύνης μετατρέπουμε το r σε z , υπολογίζουμε τα όρια εμπιστοσύνης του z και μετά τα μετατρέπουμε στην κλίμακα του r .

$$\alpha = 0,05 \quad r = -0,845 \quad z = \frac{1}{2} \ln \left[\frac{1+r}{1-r} \right] = \frac{1}{2} \ln \left[\frac{1+(-0,845)}{1-(-0,845)} \right] = -1,238$$

$$KO = z - t_{(a,\infty)} T\Sigma z = -1,238 - 1,96(1/\sqrt{5}) = -2,16$$

$$AO = z + t_{(a,\infty)} T\Sigma z = -1,238 + 1,96(1/\sqrt{5}) = -0,36$$

Χρησιμοποιώντας τον πίνακα μετατροπής των τιμών z σε κλίμακα r :

$$KO \cong -0,97 \quad \text{και} \quad AO \cong -0,35$$

Πειραματικοί Σχεδιασμοί

```
> attach(Pearson)
> cor.test(X, Y, method = c("pearson"))
```

Pearson's product-moment correlation

data: X and Y

t = -3.8769, df = 6, p-value = 0.0082

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

-0.9713868 -0.3480893

sample estimates:

cor

-0.8453958

Σύγκριση δυο συντελεστών συσχέτισης

Για σύγκριση δυο συντελεστών συσχέτισης τον σκοπό αυτό χρησιμοποιούμε το κριτήριο του z .

$$z = \frac{z_1 - z_2}{\sigma_{(z_1 - z_2)}}$$

Όπου:

$$\sigma_{(z_1 - z_2)} = \sqrt{\frac{1}{n_1 - 3} + \frac{1}{n_2 - 3}}$$

Όταν δεν απορρίπτεσαι η μηδενική υπόθεση, υπολογίζουμε τον κοινό συντελεστή συσχέτισης ως εξής:

$$z_w = \frac{(n_1 - 3)z_1 + (n_2 - 3)z_2}{(n_1 - 3) + (n_2 - 3)}$$

και στην συνέχεια η τιμή z_w μετατρέπεται στην κλίμακα r .

ΑΠΛΗ ΓΡΑΜΜΙΚΗ ΣΥΜΜΕΤΑΒΟΛΗ

Η ανάλυση συμμεταβολής ή παλινδρόμησης είναι μια τεχνική προσδιορισμού της σχέσης μεταξύ δύο ή περισσότερων μεταβλητών, με σκοπό την πρόβλεψη μιας από αυτές, μέσω των άλλων.

Η σχέση αυτή μπορεί να οφείλεται εξ ολοκλήρου ή εν μέρει στην επίδραση της μιας μεταβλητής (ανεξάρτητη X) στην άλλη μεταβλητή (εξαρτημένη Y), υπάρχει δηλαδή σχέση αιτίου-αποτελέσματος.

Το γραμμικό μοντέλο για την απλή γραμμική συμμεταβολή είναι το εξής:

$$Y_i = a + b X_i + e_i \quad i = 1, 2, \dots, n$$

όπου

Y_i = η i παρατήρηση της εξαρτημένης μεταβλητής

X_i = η i παρατήρηση της ανεξάρτητης μεταβλητής

a, b = είναι δυο άγνωστες σταθερές

e_i = το υπόλοιπο, $N(0, \sigma^2)$

Πειραματικοί Σχεδιασμοί

Με την μέθοδο της απλής γραμμική συμμεταβολής ψάχνουμε να βρούμε μια ευθεία η οποία θα προσαρμόζει καλύτερα στα δεδομένα μας. Για τον σκοπό αυτό χρησιμοποιούμε την μέθοδο των ελαχίστων τετραγώνων, η οποία συνίσταται στον προσδιορισμό των παραμέτρων a και b , έτσι ώστε το άθροισμα των τετραγώνων των αποκλίσεων των σημείων από την ευθεία να είναι ελάχιστο.

$$\sum_{i=1}^n e_i = \sum_{i=1}^n (Y_i - a - bX_i)$$

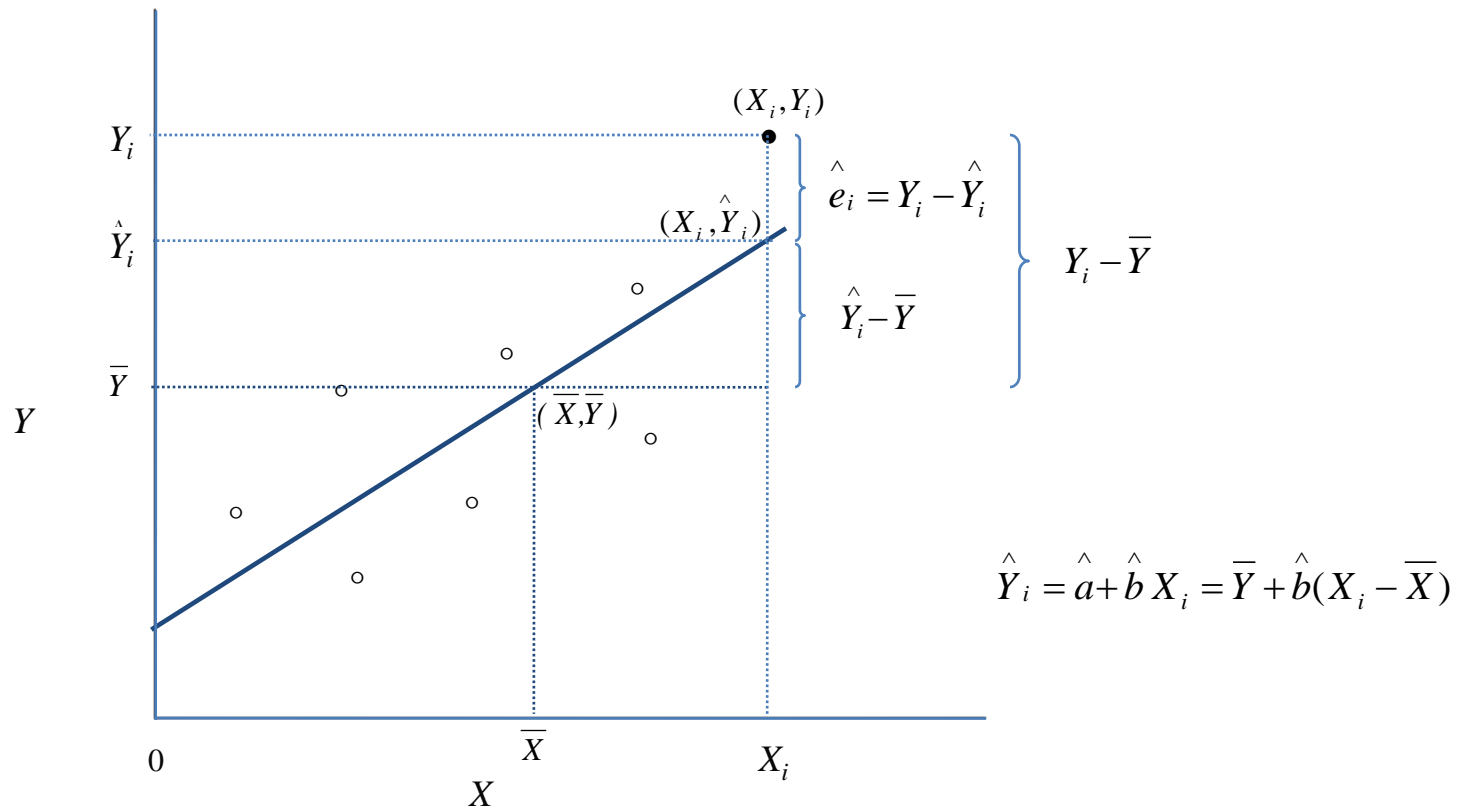
Οι τιμές των παραμέτρων a και b , που ελαχιστοποιούν την σχέση καλούνται εκτιμήτριες ελαχίστων τετραγώνων, συμβολίζονται με \hat{a} και \hat{b} και δίνονται από τις σχέσεις:

$$\hat{a} = \bar{Y} - \hat{b} \bar{X} \quad \text{και} \quad \hat{b} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{A\Gamma_{XY}}{A\Gamma_x}$$

Η ευθεία καλείται ευθεία ελαχίστων τετραγώνων ή ευθεία συμμεταβολής της Y πάνω στη X .

$$\hat{Y}_i = \hat{a} + \hat{b} X_i$$

Πειραματικοί Σχεδιασμοί



α = η τομή της γραμμής στον άξονα των Y

b = ο συντελεστής συμμεταβολής ή κλίση της ευθείας

Για μεταβολή του X κατά μία μονάδα έχουμε μεταβολή του Y κατά b μονάδες

Όρια εμπιστοσύνης των παραμέτρων της συμμεταβολής

Το s^2 παρέχει μία εκτίμηση της διασποράς σ^2 των τυχαίων αποκλίσεων των σημείων από την ευθεία γραμμή και υπολογίζεται από τη σχέση:

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$s_a^2 = s^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\text{ATX}} \right) \quad s_b^2 = \frac{s^2}{\text{ATX}}$$

Επομένως τα διαστήματα εμπιστοσύνης για τις παραμέτρους a , b ορίζονται από τις σχέσεις:

$$\hat{a} \pm s_a t_{(n-2, a/2)} \quad \text{και} \quad \hat{b} \pm s_b t_{(n-2, a/2)}$$

Δοκιμασίες σημαντικότητας

Για τον έλεγχο σημαντικότητας των δυο παραμέτρων χρησιμοποιείται η δοκιμασία του t για $BE = n - 2$.

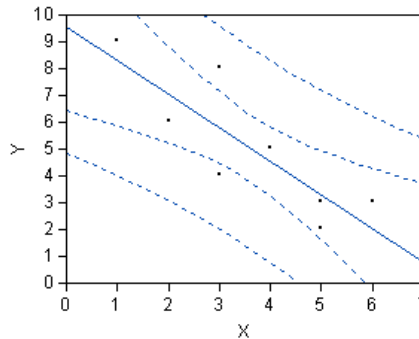
$$t = \frac{\hat{a}}{s_a} = \frac{\hat{a}}{s \sqrt{\left(\frac{1}{n} + \frac{\bar{X}^2}{ATx} \right)}} \quad \text{και} \quad t = \frac{\hat{b}}{s_b} = \frac{\hat{b}}{s \sqrt{\frac{1}{ATx}}}$$

Το διάστημα μέσης πρόβλεψης του Y για συγκεκριμένη τιμή του X , είναι:

$$\hat{a} + \hat{b} X \pm s \sqrt{\left(\frac{1}{n} + \frac{\bar{X}^2}{\text{ATX}} \right)} t_{(n-2, \alpha/2)}$$

Το διάστημα ατομικής πρόβλεψης του Y για συγκεκριμένη τιμή του X , είναι:

$$\hat{a} + \hat{b} X \pm s \sqrt{\left(1 + \frac{1}{n} + \frac{\bar{X}^2}{\text{ATX}} \right)} t_{(n-2, \alpha/2)}$$



Οι ζώνες εμπιστοσύνης για τους μέσους είναι στενότερες από ότι στις ατομικές παρατηρήσεις

Προϋποθέσεις της ανάλυσης συμμεταβολής

1. Υπάρχει γραμμική σχέση μεταξύ X και Y (γραμμικότητα).
2. Οι κατανομές της Y έχουν ίδια διασπορά για όλα τα επίπεδα της X (ομοσκεδαστικότητα).
3. Οι τιμές της Y που αντιστοιχούν στα διάφορα επίπεδα της X είναι ανεξάρτητες μεταξύ τους.
4. Η κατανομή της Y για όλα τα επίπεδα της X είναι κανονική.

$$Y_i - \bar{Y} = (\hat{Y}_i - \bar{Y}) + (Y_i - \hat{Y}_i)$$

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Το ΑΤολ. εκφράζει τη συνολική μεταβλητότητα των παρατηρήσεων Y_i .

$$ATολ. = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - (\sum_{i=1}^n Y_i)^2 / n$$

Το ΑΤσυμ. εκφράζει τη μεταβλητότητα των προσαρμοσμένων τιμών \hat{Y}_i .

$$ATσυμ. = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \left(\sum_{i=1}^n X_i Y_i - \sum_{i=1}^n (X_i) \sum_{i=1}^n (Y_i) / n \right)^2 / \left(\sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2 / n \right)$$

Το ΑΤαποκλ. εκφράζει τη μεταβλητότητα των Y_i σε σχέση με τις αντίστοιχες τιμές \hat{Y}_i .

$$ATαποκλ. = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Πίνακας Ανάλυσης Συμμεταβολής

Πηγή παρ/τητας	ΒΕ	ΑΤ	ΜΤ	F
Συμμεταβολή	1	$\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$	$\frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{1}$	$\frac{\text{ΜΤ}_{\text{συμ}}}{\text{ΜΤ}_{\text{απ}}}$
Αποκλίσεις	n - 2	$\sum_{i=1}^n (Y_i - \hat{Y})^2$	$\frac{\sum_{i=1}^n (Y_i - \hat{Y})^2}{n - 2}$	
Σύνολο	n - 1	$\sum_{i=1}^n (Y_i - \bar{Y})^2$		

Παράδειγμα:

X	Y	XY
1	9	9
2	6	12
3	8	24
3	4	12
4	5	20
5	3	15
5	2	10
6	3	18
ΣX = 29	ΣY = 40	ΣXY = 120
ΣX ² = 125	ΣY ² = 244	n = 8

$$b = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n X_i Y_i - \sum_{i=1}^n (X_i) \sum_{i=1}^n (Y_i) / n}{\sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2 / n} = \frac{120 - 29 * 40 / 8}{125 - 29^2 / 8} = \frac{-25}{19,875} = -1,257$$

$$\hat{a} = \bar{Y} - \hat{b} \bar{X} = 40/8 - (-1,257) * 29/8 = 9,55$$

$$\hat{Y} = \hat{a} + \hat{b} X = 9,55 - 1,257 X$$

Έλεγχος ύπαρξης γραμμικής σχέσης μεταξύ Y και X .

Δοκιμασία του F

$$ATολ. = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - \left(\sum_{i=1}^n Y_i\right)^2 / n = 244 - 40^2 / 8 = 44$$

$$ATσυμ. = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \frac{\left(\sum_{i=1}^n X_i Y_i - \sum_{i=1}^n (X_i) \sum_{i=1}^n (Y_i) / n\right)^2}{\sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i\right)^2 / n} = \frac{(120 - 29 * 40 / 8)^2}{125 - 29^2 / 8} = \frac{(25)^2}{19,875} = 31,447$$

$$ΑΤαποκλ. = ΑΤολ. - ΑΤσυμ. = 44 - 31,447 = 12,553$$

Πίνακας Ανάλυσης Συμμεταβολής

Πηγή παρ/τητας	BE	AT	MT	F	Fπιν
Συμμεταβολή	1	31,447	31,447	15,03	5,98
Αποκλίσεις	6	12,553	2,0922		
Σύνολο	7	44			

Επειδή το F για τη συμμεταβολή είναι μεγαλύτερο από την κρίσιμη τιμή, απορρίπτεται η μηδενική υπόθεση και επομένως υπάρχει γραμμική σχέση μεταξύ X και Y .

Έλεγχος ύπαρξης γραμμικής σχέσης μεταξύ Y και X .

Δοκιμασία του t

$$s^2 = \frac{AT_Y - \frac{A\Gamma_{XY}^2}{AT_X}}{n-2} = \frac{44 - 31,447}{6} = 2,092$$

$$s_b^2 = \frac{s^2}{AT_X} = \frac{2,092}{19,875} = 0,10$$

$$t = \frac{\hat{b}}{s_b} = \frac{-1,257}{\sqrt{0,105}} = \frac{-1,257}{0,3244} = -3,88$$

Πειραματικοί Σχεδιασμοί

X	Y	\hat{Y}	$e = Y - \hat{Y}$	$(Y - \hat{Y})^2$	$Y - \bar{Y}$	$(Y - \bar{Y})^2$
1	9	8,3	0,7	0,5	4	16
2	6	7,0	-1,0	1,1	1	1
3	8	5,8	2,2	4,9	3	9
3	4	5,8	-1,8	3,2	-1	1
4	5	4,5	0,5	0,2	0	0
5	3	3,3	-0,3	0,1	-2	4
5	2	3,3	-1,3	1,6	-3	9
6	3	2,0	1,0	1,0	-2	4
$\Sigma X = 29$	$\Sigma Y = 40$			12,55		44

Πειραματικοί Σχεδιασμοί

```
> attach(Regression)
> fit=lm(Y~X)
> library(car)
> dwt(fit)
```

lag	Autocorrelation	D-W Statistic	p-value
1	-0.7070084	3.297526	0.088

Alternative hypothesis: rho != 0

Πειραματικοί Σχεδιασμοί

> summary(fit)

Call:

lm(formula = Y ~ X)

Residuals:

Min	1Q	Median	3Q	Max
-1.7862	-1.1006	0.1006	0.7704	2.2138

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.5597	1.2825	7.454	0.000301 ***
X	-1.2579	0.3245	-3.877	0.008200 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.446 on 6 degrees of freedom

Multiple R-squared: 0.7147, Adjusted R-squared: 0.6671

F-statistic: 15.03 on 1 and 6 DF, p-value: 0.0082

> anova(fit)

Analysis of Variance Table

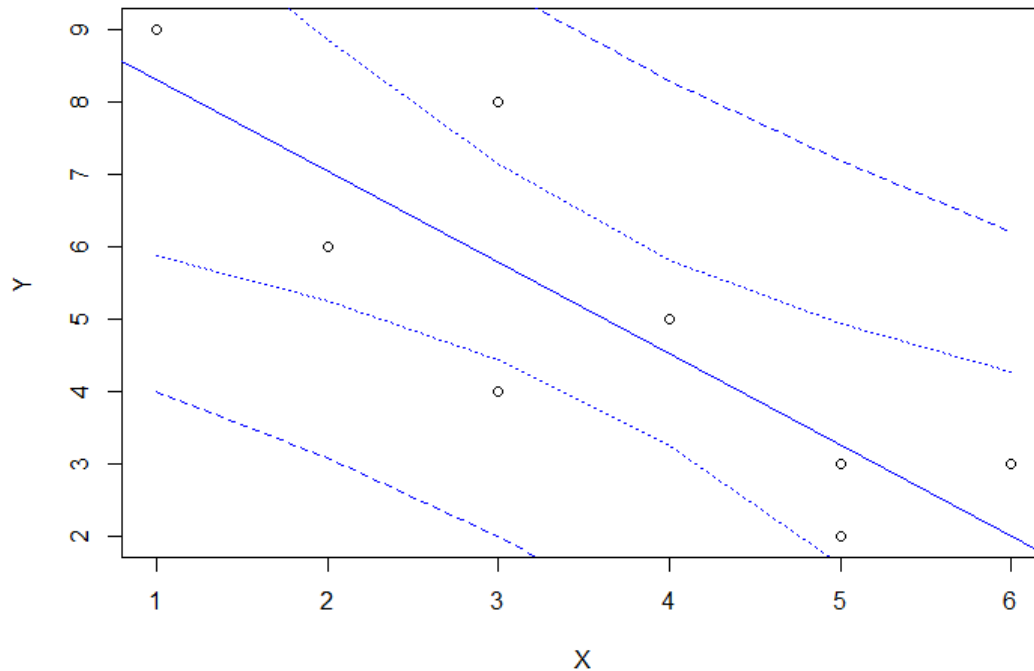
Response: Y

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
X	1	31.447	31.4465	15.03	0.0082 **
Residuals	6	12.553	2.0922		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Πειραματικοί Σχεδιασμοί

```
> plot(X, Y)
> abline(lm(Y~X), col="blue")
> conf_interval=predict(fit, interval="confidence",level = 0.95)
> lines(X, conf_interval[,2], col="blue", lty=2)
> lines(X, conf_interval[,3], col="blue", lty=2)
> pred_interval= predict(fit, interval="prediction", level = 0.95)
> lines(X, pred_interval[,2], col="blue", lty=2)
> lines(X, pred_interval[,3], col="blue", lty=2)
```



Πειραματικοί Σχεδιασμοί

> conf_interval

	fit	lwr	upr
1	8.301887	5.8710462	10.732727
2	7.044025	5.2467368	8.841314
3	5.786164	4.4400250	7.132302
4	5.786164	4.4400250	7.132302
5	4.528302	3.2420220	5.814582
6	3.270440	1.6098600	4.931020
7	3.270440	1.6098600	4.931020
8	2.012579	-0.2504108	4.275568

Σύγκριση των συντελεστών συμμεταβολής

$$t = \frac{b_1 - b_2}{S_{(b_1 - b_2)}}$$

$$S_{(b_1 - b_2)} = s_{yp}^2 \sqrt{\left[\frac{1}{AT_{1XX}} + \frac{1}{AT_{2XX}} \right]}$$

$$s_{yp}^2 = \frac{AT_{1YY} - b_1 A\Gamma_{1XY} + AT_{2YY} - b_1 A\Gamma_{2XY}}{n_1 + n_2 - 4}$$

Και συγκρίνεται με την τιμή του πίνακα t για BE= $n_1 + n_2 - 2$