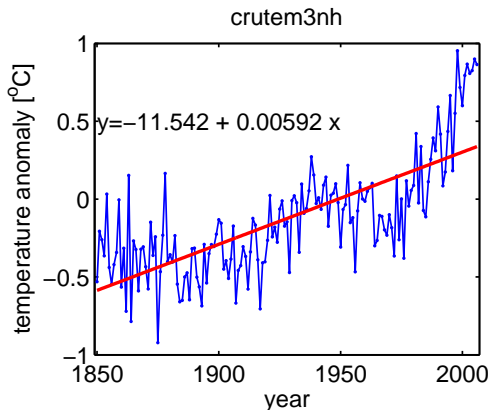


# Συσχέτιση και Παλινδρόμηση

Δημήτρης Κουγιουμτζής

10 Μαΐου 2011



**Συσχέτιση:** γραμμική / μη-γραμμική,  
μηδενική / ασθενής / ισχυρή

**Παλινδρόμηση:** γραμμική / μη-γραμμική,  
απλή / πολλαπλή

Δύο τ.μ.  $X$  και  $Y$  συσχετίζονται:

- ▶ Η μία επηρεάζει την άλλη
- ▶ Επηρεάζονται και οι δύο από κάποια άλλη

$X$ : χρόνος ως την αποτυχία ενός στοιχείου κάποιας μηχανής

$Y$ : ταχύτητα του κινητήρα της μηχανής

$X$ : χρόνος ως την αποτυχία ενός στοιχείου κάποιας μηχανής

$Y$ : θερμοκρασία του στοιχείου της μηχανής

$\sigma_X^2, \sigma_Y^2$ : διασπορά

συνδιασπορά των  $X$  και  $Y$

$$\sigma_{XY} = \text{Cov}(X, Y) = E(X, Y) - E(X)E(Y),$$

## συντελεστής συσχέτισης Pearson $\rho$

$$\rho \equiv \text{Corr}(X, Y) = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

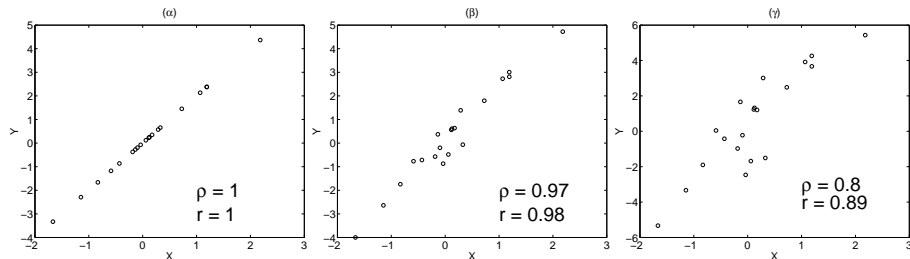
- ▶  $\rho \in [-1, 1]$
- ▶  $\rho = 1$ : τέλεια θετική συσχέτιση
- ▶  $\rho = 0$ : καμιά (γραμμική) συσχέτιση
- ▶  $\rho = -1$ : τέλεια αρνητική συσχέτιση
- ▶  $\rho$  'κοντά' στο  $-1$  ή  $1 \rightarrow$  ισχυρή συσχέτιση
- ▶  $\rho$  'κοντά' στο  $0 \rightarrow$  οι τ.μ. είναι πρακτικά ασυσχέτιστες
- ▶  $\rho$  δεν εξαρτάται από τη μονάδα μέτρησης των  $X$  και  $Y$
- ▶  $\rho$  είναι συμμετρικός ως προς τις  $X$  και  $Y$ .

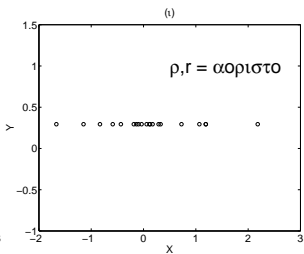
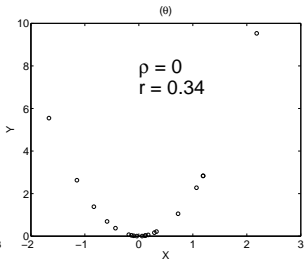
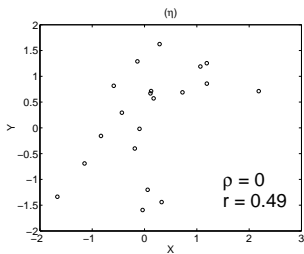
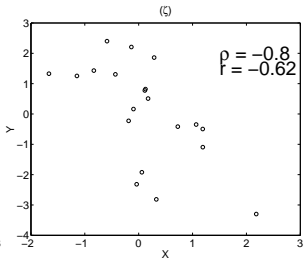
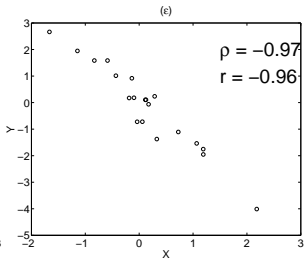
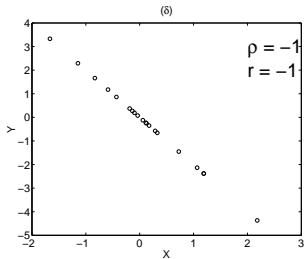
# Δειγματικός συντελεστής συσχέτισης

Παρατηρήσεις των δύο τ.μ.  $X$  και  $Y$ :

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

**διάγραμμα διασποράς**





# Σημειακή εκτίμηση του $\rho$

Εκτίμηση διασποράς

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

Εκτίμηση συνδιασποράς

$$s_{XY} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n-1} \left( \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \right)$$

**δειγματικός συντελεστής συσχέτισης (Pearson)**

$$\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} \rightarrow \hat{\rho} \equiv r = \frac{s_{XY}}{s_X s_Y}$$

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

## Συντελεστής προσδιορισμού $r^2$

(ή σε ποσοστά  $100r^2\%$ ):

Δηλώνει το ποσοστό μεταβλητότητας που μπορούμε να ερμηνεύσουμε για τη μια τ.μ. όταν γνωρίζουμε την άλλη.



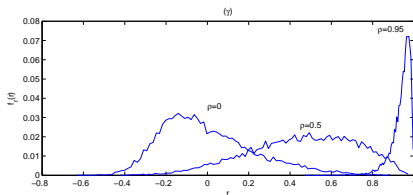
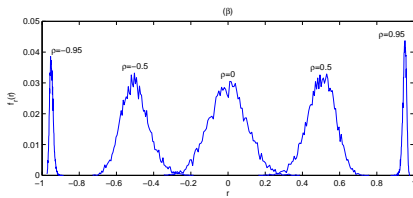
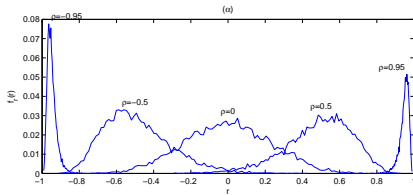
Η κατανομή του εκτιμητή  $r$  εξαρτάται από:

- ▶ Την τιμή του  $\rho$
- ▶ Το μέγεθος του δείγματος  $n$
- ▶ Την κατανομή των τ.μ.  $X$  και  $Y$ .

$(X, Y) \sim$  διμεταβλητή  
κανονική κατανομή  
 $n = 20$

$(X, Y) \sim$  διμεταβλητή  
κανονική κατανομή  
 $n = 100$

$X' = X^2$  και  $Y' = Y^2$   
 $\rho' = \text{Corr}(X', Y') = \rho^2$   
 $n = 20$



Fisher μετασχηματισμός

$$z = \tanh^{-1}(r) = 0.5 \ln \frac{1+r}{1-r}$$

Όταν το δείγμα είναι μεγάλο και από διμεταβλητή κανονική κατανομή  $\implies z \sim N(\mu_z, \sigma_z^2)$

$$\mu_z \equiv E(z) = \tanh^{-1}(\rho)$$

$$\sigma_z^2 \equiv \text{Var}(z) = 1/(n-3).$$

Μπορούμε λοιπόν να υπολογίσουμε διάστημα εμπιστοσύνης και να κάνουμε έλεγχο υπόθεσης χρησιμοποιώντας την κανονική κατανομή του  $z$ .

# Διάστημα εμπιστοσύνης για το συντελεστή συσχέτισης

$(1 - \alpha)\%$  διάστημα εμπιστοσύνης για το  $\rho$

1. Μετασχηματισμός του  $r$  στο  $z$  ( $\tanh^{-1}$ )
2.  $z \pm z_{1-\alpha/2} \sqrt{1/(n-3)}$  είναι το  $(1 - \alpha)\%$  διάστημα εμπιστοσύνης για  $\zeta$ ,  $[\zeta_l, \zeta_u]$
3. Αντίστροφος μετασχηματισμός για τα άκρα του διαστήματος  $\zeta_l$  και  $\zeta_u$

$$r_l = \tanh(\zeta_l) = \frac{\exp(2\zeta_l) - 1}{\exp(2\zeta_l) + 1}, \quad r_u = \frac{\exp(2\zeta_u) - 1}{\exp(2\zeta_u) + 1}$$

# Έλεγχος μηδενικής συσχέτισης

Έλεγχος από το διάστημα εμπιστοσύνης του  $\rho$

Αν  $[r_l, r_u]$  δεν περιέχει το 0  $\implies$  οι δύο τ.μ. συσχετίζονται.

Έλεγχος υπόθεσης  $H_0: \rho = 0$

κατανομή του  $r$  κάτω από την  $H_0$

$$t = r \sqrt{\frac{n-2}{1-r^2}} \sim t_{n-2},$$

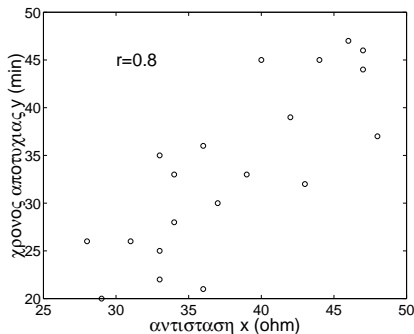
Απόφαση ελέγχου από το  $t$

$$\rho = 2 * (1 - F(t))$$

## Παράδειγμα: Συσχέτιση αντίστασης / χρόνου αποτυχίας

A/A ( $i$ )	Αντίσταση $x_i$ (ohm)	Χρόνος αποτυχίας $y_i$ (min)
1	28	26
2	29	20
3	31	26
4	33	22
5	33	25
6	33	35
7	34	28
8	34	33
9	36	21
10	36	36
11	37	30
12	39	33
13	40	45
14	42	39
15	43	32
16	44	45
17	46	47
18	47	44
19	47	46
20	48	37

# Παράδειγμα: Συσχέτιση αντίστασης / χρόνου αποτυχίας



$$\bar{x} = 38$$

$$\bar{y} = 33.5$$

$$\sum_{i=1}^{20} x_i^2 = 29634$$

$$\sum_{i=1}^{20} y_i^2 = 23910$$

$$\sum_{i=1}^{20} x_i y_i = 26305.$$

## Παράδειγμα: Συσχέτιση αντίστασης / χρόνου αποτυχίας

$$r = \frac{26305 - 20 \cdot 38 \cdot 33.5}{\sqrt{(29634 - 20 \cdot 38^2) \cdot (23910 - 20 \cdot 33.5^2)}} = 0.804.$$

Η μεταβλητότητα της μιας τ.μ. (αντίσταση ή χρόνος αποτυχίας) μπορεί να εξηγηθεί από τη συσχέτιση της με την άλλη κατά ποσοστό

$$r^2 \cdot 100 = 0.804^2 \cdot 100 = 64.64 \rightarrow \simeq 65\%.$$



$(1 - \alpha)\%$  διάστημα εμπιστοσύνης για το  $\rho$

1. Μετασχηματισμός του  $r$  στο  $z$ ,  $z = 1.110$
2. 95% διαστήματος εμπιστοσύνης του  $z$ :  $\zeta_l = 0.634$ ,  
 $\zeta_u = 1.585$
3. Αντίστροφος μετασχηματισμός

$$r_l = 0.561, \quad r_u = 0.919$$

- ▶  $r_l > 0$ ,  $\implies$  η αντίσταση και ο χρόνος αποτυχίας συσχετίζονται
- ▶ Έλεγχος υπόθεσης  $H_0: \rho = 0$   
στατιστικό  $t = 5.736 \implies p = 2 * (1 - F(t)) = 0.0000194$

# Συσχέτιση και γραμμικότητα

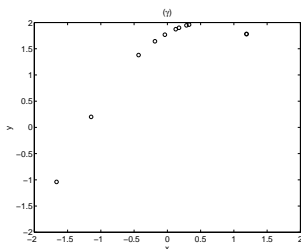
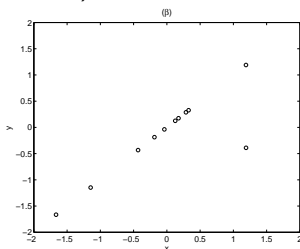
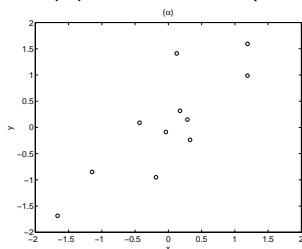
$\rho$  και  $r$  κατάλληλα για γραμμική συσχέτιση και κανονικότητα

Τρία δείγματα των  $(X, Y)$  με  $r = 0.84$ .

(α)  $(X, Y)$  από διμεταβλητή κανονική κατανομή

(β)  $Y = X$  για όλα εκτός από ένα ζευγάρι

(γ)  $Y = 2 - 0.6(X - 0.585)^2$



# Απλή Γραμμική Παλινδρόμηση

Συντελεστής συσχέτισης: γραμμική σχέση δύο τ.μ.  $X$  και  $Y$

**παλινδρόμηση**: εξάρτηση της τ.μ.  $Y$  από τη  $X$

$Y$ : **εξαρτημένη** μεταβλητή  $\Leftarrow$  είναι τ.μ.

$X$ : **ανεξάρτητη** μεταβλητή  $\Leftarrow$  δεν είναι τ.μ.

Μια μόνο ανεξάρτητη μεταβλητή: **απλή παλινδρόμηση**

Παράδειγμα: σε μια μονάδα παραγωγής ηλεκτρικής ενέργειας από λιγνίτη, για να προσδιορίσουμε το κόστος της παραγωγής ενέργειας, μελετάμε την εξάρτηση του από το κόστος του λιγνίτη.  $X$ ;  $Y$

Εξάρτηση είναι γραμμική: **απλή γραμμική παλινδρόμηση**

# Το πρόβλημα της απλής γραμμικής παλινδρόμησης

Γενικά:  $F_Y(y|X = x)$  για κάθε τιμή  $x$  της  $X$

Περιορίζουμε το πρόβλημα σε  $E(Y|X = x)$

Υπόθεση εργασίας:

$$E(Y|X = x) = \beta_0 + \beta_1 x$$

**απλή γραμμική παλινδρόμηση της  $Y$  στη  $X$**

$\beta_0 = ?$ ,  $\beta_1 = ?$

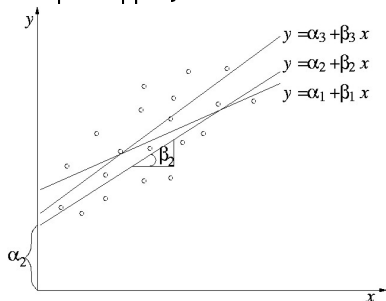
$\beta_0$ :  $y$  για  $x = 0$ , **διαφορά ύψους**

$\beta_1$ : συντελεστής του  $x$ , **κλίση** της ευθείας παλινδρόμησης ή **συντελεστής παλινδρόμησης**

# Το πρόβλημα της απλής γραμμικής παλινδρόμησης

Δείγμα:  $\{(x_1, y_1), \dots, (x_n, y_n)\}$

Πολλές ευθείες που προσαρμόζονται σε αυτό



Για  $x_i$  της  $X$  αντιστοιχούν διαφορετικές τιμές  $y_i$  της  $Y$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

**σφάλμα παλινδρόμησης:**  $\epsilon_i = y_i - E(Y|X = x_i)$

- ▶ Η μεταβλητή  $X$  είναι ελεγχόμενη.
- ▶ Η σχέση είναι πράγματι γραμμική.
- ▶  $E(\epsilon_i) = 0$  και  $\text{Var}(\epsilon_i) = \sigma_\epsilon^2$  για κάθε τιμή  $x_i$  της  $X$  ή ισοδύναμα

$$\text{Var}(Y|X = x) \equiv \sigma_{Y|X}^2 = \sigma_\epsilon^2,$$

Υποθέτουμε κανονική κατανομή

$$Y|X = x \sim N(\beta_0 + \beta_1 x, \sigma_\epsilon^2).$$

όχι απαραίτητο για σημειακή εκτίμηση των παραμέτρων.

# Σημειακή εκτίμηση παραμέτρων της απλής γραμμικής παλινδρόμησης

Το πρόβλημα: εκτίμηση των τριών παραμέτρων παλινδρόμησης:

1.  $\beta_0$ ,
2.  $\beta_1$ ,
3.  $\sigma_\epsilon^2$ .

Εκτίμηση των  $\beta_0$ ,  $\beta_1$  με τη μέθοδο των **ελαχίστων τετραγώνων**:

το άθροισμα των τετραγώνων των κατακόρυφων αποστάσεων των σημείων από την ευθεία να είναι το ελάχιστο.

# Εκτίμηση των παραμέτρων της ευθείας παλινδρόμησης

Ελαχιστοποίηση του αθροίσματος των τετραγώνων των σφαλμάτων:

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n \epsilon_i^2 \quad \text{ή} \quad \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.$$

Λύση:

$$\left. \begin{aligned} \frac{\partial \sum (y_i - \beta_0 - \beta_1 x_i)^2}{\partial \beta_0} = 0 \\ \frac{\partial \sum (y_i - \beta_0 - \beta_1 x_i)^2}{\partial \beta_1} = 0 \end{aligned} \right\} \begin{aligned} \sum_{i=1}^n y_i &= n\beta_0 + \beta_1 \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i y_i &= \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \end{aligned}$$

Εκτίμηση για την κλίση

$$\hat{\beta}_1 \equiv b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}.$$

$$s_{XY} = \frac{S_{xy}}{n-1} \quad \text{και} \quad s_X^2 = \frac{S_{xx}}{n-1}$$



Εκτίμηση του σταθερού όρου ως

$$\hat{\beta}_0 \equiv b_0 = \frac{\sum_{i=1}^n y_i - b_1 \sum_{i=1}^n x_i}{n}$$

$$b_1 = \frac{s_{XY}}{s_X^2}, \quad b_0 = \bar{y} - b_1 \bar{x}.$$

**ευθεία ελαχίστων τετραγώνων**

$$\hat{y} = b_0 + b_1 x$$

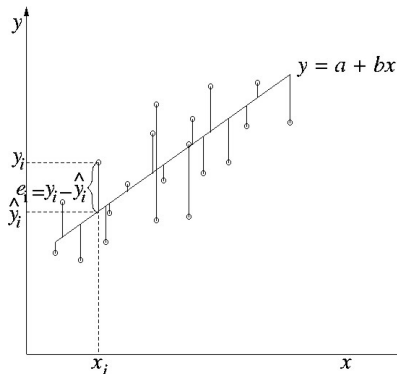
# Εκτίμηση της διασποράς των σφαλμάτων παλινδρόμησης

$y_i$ : παρατηρούμενη τιμή για  $x_i$

$\hat{y}_i$ : εκτιμώμενη τιμή από την ευθεία ελαχίστων τετραγώνων για  $x_i$ .

σφάλμα ελαχίστων τετραγώνων ή **υπόλοιπο**

$$e_i = y_i - \hat{y}_i = y_i - b_0 - b_1 x_i$$



# Εκτίμηση της διασποράς των σφαλμάτων παλινδρόμησης

$e_i$ : εκτίμηση του σφάλματος παλινδρόμησης  $e_i = y_i - \beta_0 - \beta_1 x_i$ .

Εκτίμηση της διασποράς του  $e_i$ :

$$s_\epsilon^2 \equiv \hat{\sigma}_\epsilon^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

$n - 2$ : βαθμοί ελευθερίας

$$s_\epsilon^2 = \frac{n-1}{n-2} \left( s_Y^2 - \frac{s_{XY}^2}{s_X^2} \right) = \frac{n-1}{n-2} (s_Y^2 - b_1^2 s_X^2)$$

1. Η ευθεία ελαχίστων τετραγώνων περνάει από το σημείο  $(\bar{x}, \bar{y})$

$$b_0 + b_1\bar{x} = \bar{y} - b_1\bar{x} + b_1\bar{x} = \bar{y}.$$

Η ευθεία ελαχίστων τετραγώνων μπορεί να οριστεί ως

$$y_i - \bar{y} = b_1(x_i - \bar{x}).$$

2. Η σημειακή εκτίμηση των  $\beta_0$  και  $\beta_1$  με τη μέθοδο των ελαχίστων τετραγώνων δεν προϋποθέτει σταθερή διασπορά και κανονική κατανομή της  $Y|X$ .
3. Για κάθε  $x$  της  $X$  η πρόβλεψη της  $y$  της  $Y$ :

$$\hat{y} = b_0 + b_1x$$

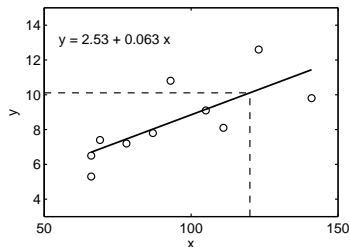
$x$  πρέπει να ανήκει στο εύρος τιμών της  $X$  από το δείγμα.

## Παράδειγμα: Απολαβή ρεύματος τρανζίστορ

Θέλουμε να μελετήσουμε σε ένα ολοκληρωμένο κύκλωμα την εξάρτηση της απολαβής ρεύματος κρυσταλλολυχνίας από την αντίσταση του στρώματος της κρυσταλλολυχνίας.

A/A ( $i$ )	Αντίσταση στρώματος $x_i$ (ohm/cm)	Απολαβή ρεύματος $y_i$
1	66	5.3
2	66	6.5
3	69	7.4
4	78	7.2
5	87	7.8
6	93	10.8
7	105	9.1
8	111	8.1
9	123	12.6
10	141	9.8

# Παράδειγμα: Απολαβή ρεύματος τρανζίστορ



$$\bar{x} = 93.9$$

$$\bar{y} = 8.46$$

$$\sum_{i=1}^{10} x_i^2 = 94131$$

$$\sum_{i=1}^{10} y_i^2 = 757.64$$

$$\sum_{i=1}^{10} x_i y_i = 8320.2$$

$$s_{XY} = 41.81 \quad s_X^2 = 662.1 \quad s_Y^2 = 4.66.$$

## Παράδειγμα: Απολαβή ρεύματος τρανζίστορ

$$b_1 = \frac{41.81}{662.1} = 0.063$$

$$b_0 = 8.46 - 0.063 \cdot 93.9 = 2.53.$$

$$s_\epsilon^2 = \frac{9}{8}(4.66 - 0.063^2 \cdot 41.81) = 2.271.$$

1.  $b_1 = 0.063$ : απολαβή ρεύματος για αύξηση της αντίστασης στρώματος κατά 1 ohm/cm
2.  $b_0 = 2.53$ : απολαβή ρεύματος όταν δεν υπάρχει αντίσταση στρώματος ( $x = 0$ )
3.  $s_\epsilon^2 = 2.271 \implies s_\epsilon = 1.507$ : το τυπικό σφάλμα της εκτίμησης της παλινδρόμησης .

## Παράδειγμα: Απολαβή ρεύματος τρανζίστορ

Πρόβλεψη απολαβή ρεύματος μπορεί να γίνει για κάθε αντίσταση στρώματος κρυσταλλολυχνίας στο διάστημα  $[66, 141]$  ohm/cm.

Για αντίσταση στρώματος  $x = 120$  ohm/cm

$$\hat{y} = 2.53 + 0.063 \cdot 120 = 10.11.$$



# Σχέση του συντελεστή συσχέτισης και παλινδρόμησης

Παλινδρόμηση:  $X$  ελεγχόμενη και  $Y$  τυχαία

Συσχέτιση:  $X$  και  $Y$  τυχαίες, αλλά υπολογίζουμε το  $r$  και για  $X$  ελεγχόμενη.

$$r = \frac{s_{XY}}{s_X s_Y} \text{ και } b_1 = \frac{s_{XY}}{s_X^2} \implies$$

$$r = b_1 \frac{s_X}{s_Y} \quad \text{ή} \quad b_1 = r \frac{s_Y}{s_X}.$$

$r$  και  $b_1$  εκφράζουν ποσοτικά τη γραμμική συσχέτιση των  $X$  και  $Y$

$b_1$  εξαρτάται από τη μονάδα μέτρησης των  $X$  και  $Y$   
Σχέση των  $r$  και  $b_1$

- ▶  $r > 0 \Leftrightarrow b_1 > 0$
- ▶  $r < 0 \Leftrightarrow b_1 < 0$
- ▶  $r = 0 \Leftrightarrow b_1 = 0$

Σχέση  $r^2$  και  $s_\epsilon^2$

$$s_\epsilon^2 = \frac{n-1}{n-2} s_Y^2 (1-r^2) \quad \text{ή} \quad r^2 = 1 - \frac{n-2}{n-1} \frac{s_\epsilon^2}{s_Y^2}.$$

Όσο μεγαλύτερο είναι το  $r^2$  (ή το  $|r|$ ) τόσο μειώνεται το  $s_\epsilon^2$

## Παράδειγμα: Απολαβή ρεύματος τρανζίστορ (συνέχεια)

Συντελεστής συσχέτισης της απολαβής ρεύματος και της αντίστασης στρώματος

$$r = \frac{s_{XY}}{s_X s_Y} = \frac{41.81}{\sqrt{662.1 \cdot 4.66}} = 0.753$$

Το  $r$  δηλώνει την σχετικά ασθενή θετική συσχέτιση

$b_1 = 0.063$  εξηγεί το βαθμό εξάρτησης;

$s_\epsilon^2 = 2.249$  εξηγεί το βαθμό εξάρτησης;

# Διάστημα εμπιστοσύνης των παραμέτρων της απλής γραμμικής παλινδρόμησης

$b_1$  και  $b_0$  είναι εκτιμητές των  $\beta_1$  και  $\beta_0$

κατανομή των  $b_1$  και  $b_0$ ;

$b_1$  γραμμικός συνδυασμός των τ.μ.  $y_1, \dots, y_n$

$$\mu_{b_1} \equiv E(b_1) = \beta_1$$

$$\sigma_{b_1}^2 \equiv \text{Var}(b_1) = \frac{\sigma_\epsilon^2}{S_{xx}} \implies \sigma_{b_1} = \sigma_\epsilon / \sqrt{S_{xx}}$$

Εκτίμηση:

$$s_{b_1} = \frac{s_\epsilon}{\sqrt{S_{xx}}}.$$

# Διάστημα εμπιστοσύνης των παραμέτρων της απλής γραμμικής παλινδρόμησης

$Y$  ακολουθεί κανονική κατανομή  $\implies b_1$  ακολουθεί κανονική κατανομή

$(1 - \alpha)\%$  διάστημα εμπιστοσύνης του  $\beta_1$ :

$$b_1 \pm t_{n-2, 1-\alpha/2} s_{b_1} \quad \text{ή} \quad b_1 \pm t_{n-2, 1-\alpha/2} \frac{s_\epsilon}{\sqrt{S_{xx}}}$$

Αντίστοιχα για  $b_0$

$$\sigma_{b_0} = s_\epsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

$(1 - \alpha)\%$  διάστημα εμπιστοσύνης του  $\beta_0$ :

$$b_0 \pm t_{n-2, 1-\alpha/2} s_\epsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

# Έλεγχος υπόθεσης για τις παραμέτρους της απλής γραμμικής παλινδρόμησης

$$H_0: \beta_1 = \beta_1^0$$

Στατιστικό παραμετρικού ελέγχου

$$t = \frac{b_1 - \beta_1^0}{s_b} = \frac{(b_1 - \beta_1^0)\sqrt{S_{xx}}}{s_\epsilon}, \quad t \sim t_{n-2}$$

Ιδιαίτερο ενδιαφέρον έχει  $H_0: \beta_1 = 0$  ή η  $Y$  δεν εξαρτάται από την  $X$ .

$$H_0: \beta_0 = \beta_0^0$$

Στατιστικό παραμετρικού ελέγχου

$$t = \frac{b_0 - \beta_0^0}{s_\epsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}}, \quad t \sim t_{n-2}$$

$\hat{y} = b_0 + b_1x$ : εκτιμητής της  $E(Y|X = x) = \beta_0 + \beta_1x$  για κάποιο  $x$

$\hat{y}$ : γραμμικός συνδυασμός των τ.μ.  $y_1, \dots, y_n$ , των  $x_1, \dots, x_n$  και  $x$ .

$$\mu_{\hat{y}} \equiv E(\hat{y}) = E(Y|X = x) = \beta_0 + \beta_1x.$$

$$\sigma_{\hat{y}}^2 \equiv \text{Var}(\hat{y}) = \sigma_{\epsilon}^2 \left( \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}} \right)$$

Εκτίμηση της τυπικής απόκλισης του  $\hat{y}$

$$s_{\hat{y}} = s_{\epsilon} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}.$$

## Διαστήματα πρόβλεψης

$\hat{y}$  για κάποιο  $x$  ακολουθεί κανονική κατανομή

**$(1 - \alpha)\%$  διάστημα εμπιστοσύνης της μέσης τιμής του  $Y$  για κάποιο  $x$**

$$\hat{y} \pm t_{n-2, 1-\alpha/2} s_{\hat{y}} \quad \text{ή} \quad (b_0 + b_1 x) \pm t_{n-2, 1-\alpha/2} s_{\epsilon} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$$

Λέγεται και **διάστημα της μέσης πρόβλεψης**: τα όρια της πρόβλεψης για τη μέση (αναμενόμενη) τιμή της  $Y$  για κάποιο  $x$

**$(1 - \alpha)\%$  διάστημα πρόβλεψης για μια παρατήρηση  $y$  της  $Y$  για κάποιο  $x$**

$$\hat{y} \pm t_{n-2, 1-\alpha/2} \sqrt{s_{\epsilon}^2 + s_{\hat{y}}^2} \quad \text{ή} \quad (b_0 + b_1 x) \pm t_{n-2, 1-\alpha/2} s_{\epsilon} \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$$



## Παράδειγμα: Απολαβή ρεύματος τρανζίστορ (συνέχεια)

Υπολογίσαμε  $b_1 = 0.063$ ,  $b_0 = 2.53$ ,  $s_\epsilon = 1.507$

Ακρίβεια / σημαντικότητα εκτιμήσεων και πρόβλεψης:

$\beta_1$ :  $s_{b_1} = 0.0195$

95% διάστημα εμπιστοσύνης για  $\beta_1$

$$0.063 \pm 2.306 \cdot 0.0195 \Rightarrow [0.018, 0.108].$$

$\beta_1$  σημαντικά διάφορο του 0

Στατιστικό ελέγχου για  $H_0: \beta_1 = 0$

$$t = \frac{0.063}{0.0195} = 3.235$$

$t > t_{0.975,8} = 2.306 \implies H_0$  απορρίπτεται.

Η τιμή  $t = 3.235$  αντιστοιχεί σε  $p = 0.012$

## Παράδειγμα: Απολαβή ρεύματος τρανζίστορ (συνέχεια)

$$\beta_0: s_{b_0} = 1.894$$

95% διάστημα εμπιστοσύνης για  $\beta_0$

$$2.53 \pm 2.306 \cdot 1.894 \Rightarrow [-1.837, 6.898]$$

$\beta_0$  μπορεί να είναι 0

Στατιστικό ελέγχου για  $H_0: \beta_1 = 0$

$$t = \frac{2.53}{1.894} = 1.336$$

$t < t_{0.975,8} = 2.306 \implies H_0$  δεν απορρίπτεται.

Η τιμή  $t = 1.336$  αντιστοιχεί σε  $p = 0.218$

## Παράδειγμα: Απολαβή ρεύματος τρανζίστορ (συνέχεια)

διάστημα πρόβλεψης για αντίσταση στρώματος  $x = 120 \text{ ohm/cm}$

$\hat{y}$ :

$$s_{\hat{y}} = 1.507 \sqrt{\frac{1}{10} + \frac{(120 - 93.9)^2}{9 \cdot 662.1}} = 0.698$$

95% διάστημα πρόβλεψης της  $\hat{y}$  για  $x = 120$

$$10.108 \pm 2.306 \cdot 0.698 \Rightarrow [8.499, 11.717]$$

παρατήρηση  $y$ :

95% διάστημα πρόβλεψης για μια (μελλοντική) παρατήρηση  $y$   
για  $x = 120$

$$10.108 \pm 2.306 \cdot 1.507 \sqrt{1 + \frac{1}{10} + \frac{(120 - 93.9)^2}{9 \cdot 662.1}} \Rightarrow [6.279, 13.937]$$

# Παράδειγμα: Απολαβή ρεύματος τρανζίστορ (συνέχεια)

