

23-2-2023

Μοντέλα Παλινδρόμησης

Y : εξαρτημένη μεταβλητή / απόκριση
dependent variable / response variable

X_1, X_2, \dots, X_p ανεξάρτητες μεταβλητές / παρόγοι
independent variables / predictors / factors

Y : ποσοτική (scale)

X_1, X_2, \dots, X_p : (ορισμένα) αριθμητικά
(αλλιώς γενικεύουμε)

Στατιστικό Μοντέλο

$$E(Y | X_1 = x_1, X_2 = x_2, \dots, X_p = x_p) = f(x_1, x_2, \dots, x_p)$$

συνάρτηση παλινδρόμησης

regression function

(αγνώστη ...)

Παράδειγμα

Y = score

X_1 = age

X_2 = χρόνος ασκησης

$$f(52, 120) = E(Y | X_1 = 52, X_2 = 120) =$$

= μέσο όρο ατόμων ηλικίας 52 είναι
που αφορούν 120 min/εβδ.

Παραδείγματα

1) Αν $f(x_1, \dots, x_p) = C = \text{σταθερό}$

Η συνάρτηση της Y με αυτή ως ανεξ. μεταβλητές

2) f : άγνωστη συνάρτηση, πρέπει να "εξεπαιδευτεί" από δεδομένα

π.χ. 1) $f(x_1, x_2) = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2$

2) $f(x_1, x_2) = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_1^2 + \alpha_4 x_2^2 + \alpha_5 x_1 x_2$

$\alpha_0, \alpha_1, \dots$

αγνωστές παράμετροι
που πρέπει να εκτιμηθούν

3) $f(x_1, x_2, \dots, x_p) = C$

$E(Y) = C$ (άγνωστη)

↳ εκτίμηση των $E(Y)$
για το θ των πιθανοτήτων

εκτίμηση
μ
μέσης
τιμής

4) $f(x) = a_0 + a_1 \sin(2\pi f x)$

a_0, a_1, f : άγνωστες παράμετροι
ή πιο a_0, a_1 : άγνωστα ϵ η f γνωστά.

$$6) f(x) = \alpha_0 e^{\alpha_1 x_1 + \alpha_2 x_2}$$

Γραμμικά Μοντέλα

Γενική μορφή:

$$E(Y) = b_0 + b_1 x_1 + \dots + b_p x_p$$

γραμμική συνάρτηση των b_0, b_1, \dots, b_p !!

x_1, x_2, \dots, x_p πρέπει να είναι οποιοδήποτε συναρτήσει αλληλ ανεξαρτητών γραμμικά n περι.

Παραδείγματα

$$E(Y) = b_0 + b_1 x + b_2 x^2 \quad (x = \text{age})$$

αγνωστων γραμμικά τετράγα

$$x_1 = x, \quad x_2 = x^2$$

$$E(Y) = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_1^2 + b_4 x_1 x_2$$

$x_1 \quad x_2 \quad x_3 \quad x_4$

γραμμικά

$$E(Y) = b_0 + b_1 e^{x_1} + b_2 \sin(20f x_2) \quad f: \text{γνωστή}$$

$x_1 \quad x_2$

$$E(Y) = b_0 + e^{b_1 x_1} + e^{b_2 x_2}$$

b_0, b_1, b_2 : άγνωστα

μη γραμμικό

$$E(Y) = b_0 + b_1 x_1 + b_1^2 x_2$$

μη γραμμικό
μοτέτζο

Τύποι παλινδρόμησης

$$E(Y) = b_0 + b_1 x$$

Γραμμική παλινδρόμηση

(Γραμμική σχέση ως

$E(Y)$ ως προς x

$$E(Y) = b_0 + b_1 x + b_2 x^2 + b_3 x^3$$

Πολυωνομική παλινδρόμηση

Πολυωνομική σχέση

EY ως προς x

Γραμμικά μοτέτζο

Απλό Μονοπαραγοντικό Μοντέλο

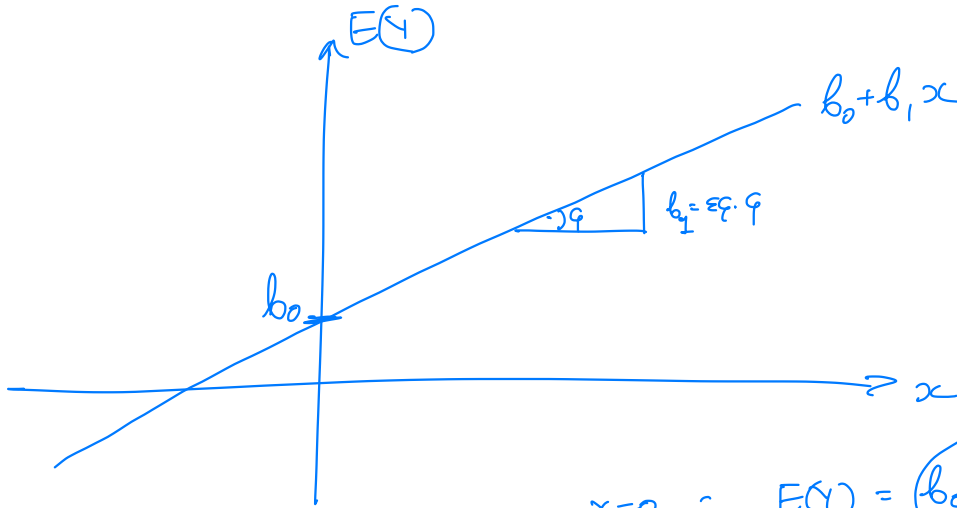
X : ανεξ. μεταβλητή

Y : εξαρτ. "

παρμελει δτιον μεταξυ x κ' εγ υλοθεση

$$E(Y|X=x) = b_0 + b_1 x$$

b_0
 b_1 } αγνωστες παραμετροι



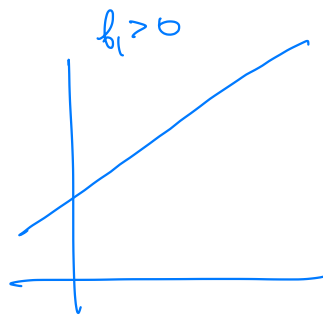
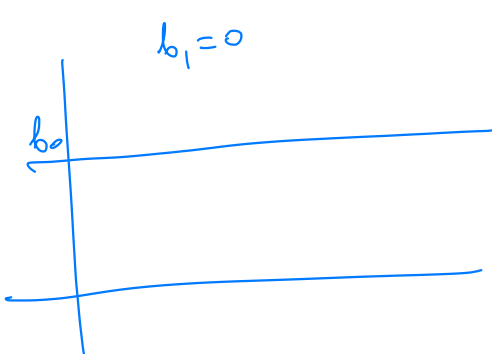
$x=0 : E(Y) = b_0 = \text{σταθ. όρος (intercept)}$

b_1 : κλιση (ρυθμός μεταβολής)

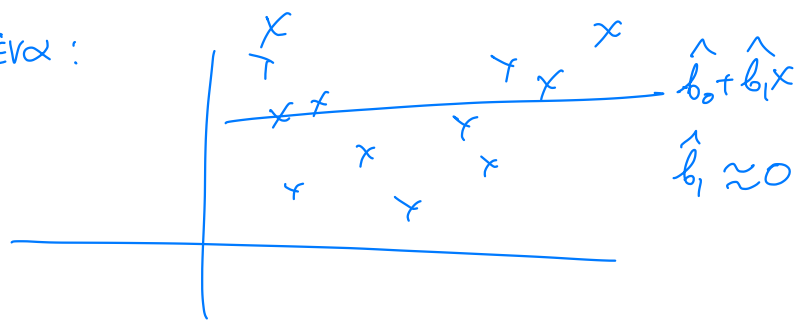
Αν $x = x_0 : E(Y) = b_0 + b_1 x_0$

$x = x_0 + 1 : E(Y) = b_0 + b_1 (x_0 + 1) = b_0 + b_1 x_0 + b_1$

b_1 : μεταβολή της $E(Y)$ για κάθε μονάδα μεταβολής της X .



Αν δεδομένα:



Εναλλακτική μορφή μοντέλου παλινδρόμησης

$$E(Y|X=x) = b_0 + b_1 x$$

αίθρων ϵ :
παλινδρόμηση

ισοδύναμα: $Y = b_0 + b_1 x + \epsilon$

ϵ : τυχαία διατάραξη

$$E(\epsilon) = 0$$

Σε ένα δείγμα n παρατηρήσεων

$$\sigma^2(\epsilon) ?$$

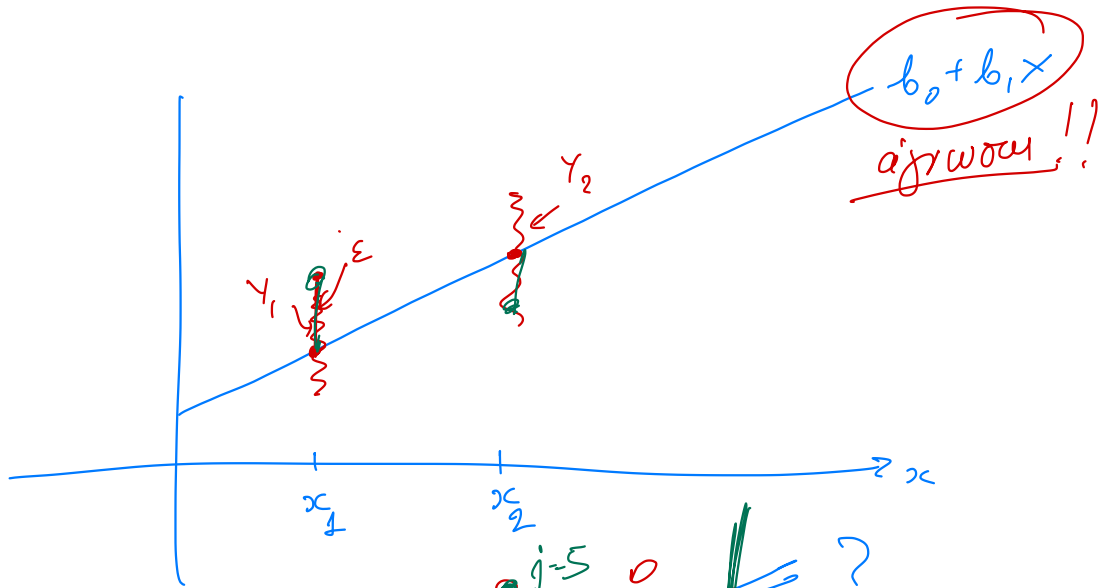
X	Y
x_1	y_1
\vdots	\vdots
x_n	y_n

$$y_j = b_0 + b_1 x_j + \epsilon_j$$

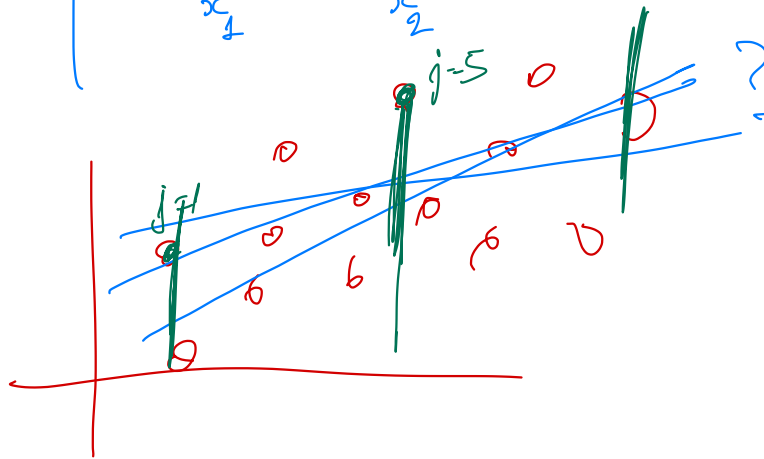
η j παρατήρηση έχει
 $x = x_j$ κ' η αντίκριση
 ως είναι τυχαία μεταβλητή
 y_j με μέση τιμή $b_0 + b_1 x_j$
 κ' μια τυχαία απόκλιση ϵ_j

Ο ληψυφύλι των απόψεων με $X = x_j$ έχει ως

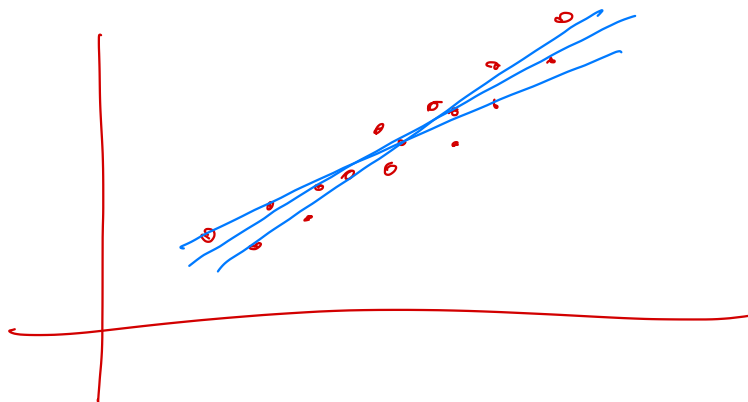
απόκριση y_j τυχαία με μέση τιμή $b_0 + b_1 x_j$



$\sigma^2(\epsilon) \gg$

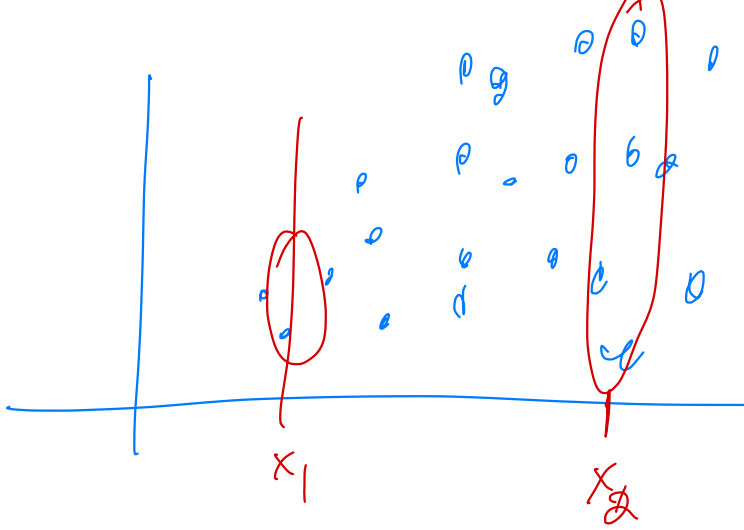


$\sigma^2(\epsilon) \ll$



Υποθέσεις Παλινδρόμησης

- ① $Y_j = b_0 + b_1 x_j + \epsilon_j \quad j=1, 2, \dots, n \quad \leftarrow$
- ② $E(\epsilon_j) = 0 \quad (E(Y_j) = b_0 + b_1 x_j) \quad \leftarrow \checkmark$
- ③ $\text{Var}(\epsilon_j) = \sigma^2 \quad \forall j \quad (\text{σταθερή ως προς } j) \quad \leftarrow \text{ομοσκεδαστικότητα}$
- ④ $\epsilon_1, \epsilon_2, \dots$ αυτοχόρτα (οι παρατηρήσεις στο δείγμα αυτοχόρτες!!) \leftarrow



Εξερευνητικότητα

Παραβιάζει τις υποθέσεις

4) Παραβιάζεται π.χ. όταν

1) x_1, x_2, \dots διαδοχικές χρονικές μετρήσεις. χρονο-
στίρες
 π.χ. x_1, x_2, \dots ως μια σειρά (αυτο-
συσχέτιση)
 y_1, y_2, \dots ανεξαρτητές.

Μοτίβα χρονοσειρών

2) π.χ.
 (x_1, x_2, x_3)
 (y_1, y_2, y_3) ; (x_4, x_5, x_6)
 (y_4, y_5, y_6) ; - - -
- - -

ομάδα 1 ομάδα 2

Μοτίβα Επαναλαμβανόμενων Μετρήσεων
 (Longitudinal / Repeated measurements)

Με βάση τις υποθέσεις 1-4 \Rightarrow αποδεικνύονται
αμερομηψία
Εκτιμώση

ΔΕ, ΕΛΕΥΧΟΥΣ ??

5) $\varepsilon_j \sim \mathcal{N}(0, \sigma^2)$ 4: ε_j ανεξάρτητα

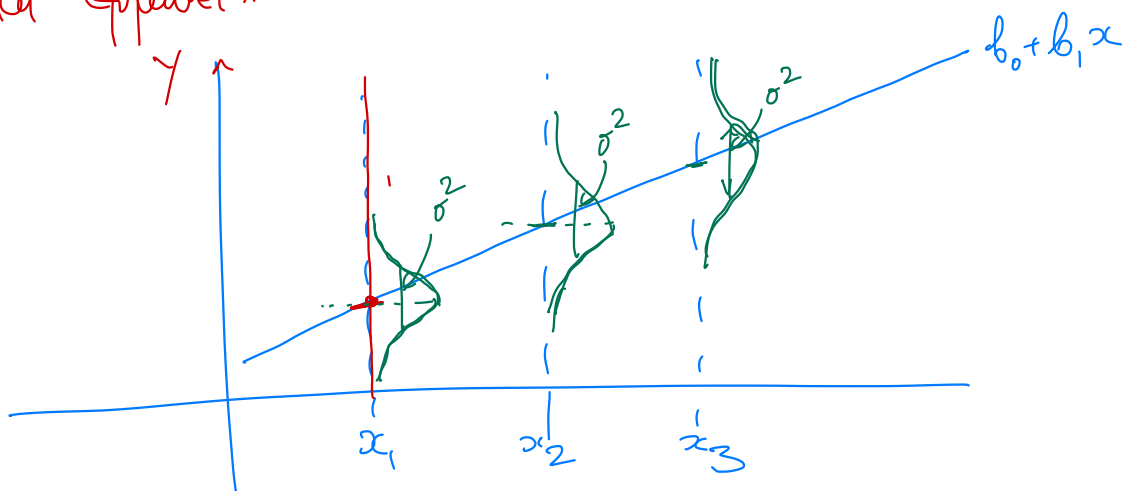
Γενικευμένο γραμμικό μοντέλο: άλλες υποθέσεις
κατανομής

Οι υποθέσεις 1-5 γίνονται για
2ο γενικό γραμμικό μοντέλο

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_p X_p + \underline{\underline{\varepsilon}}$$

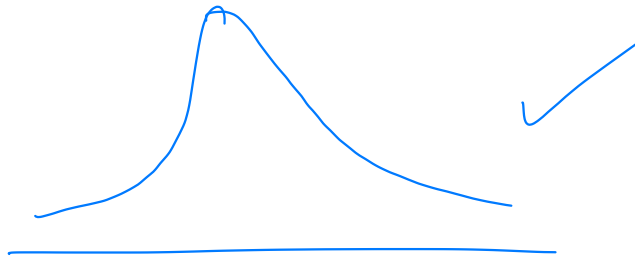
Άγνωστες παράμετροι $\left\{ b_0, b_1, \dots, b_p, \sigma^2 \right\}$

Τεωρητική ερμεία

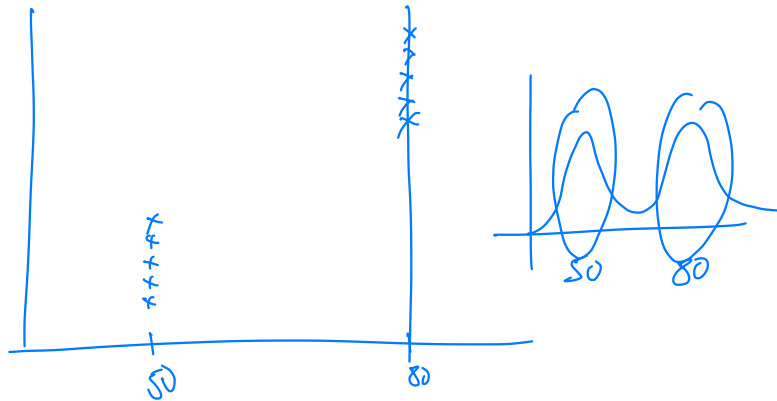
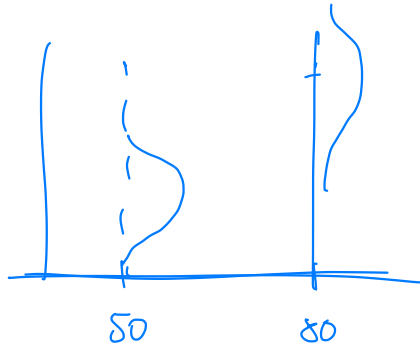


$$Y_j \sim N(b_0 + b_1 x_j, \sigma^2)$$

Y_j από την ίδια κανονική κατανομή



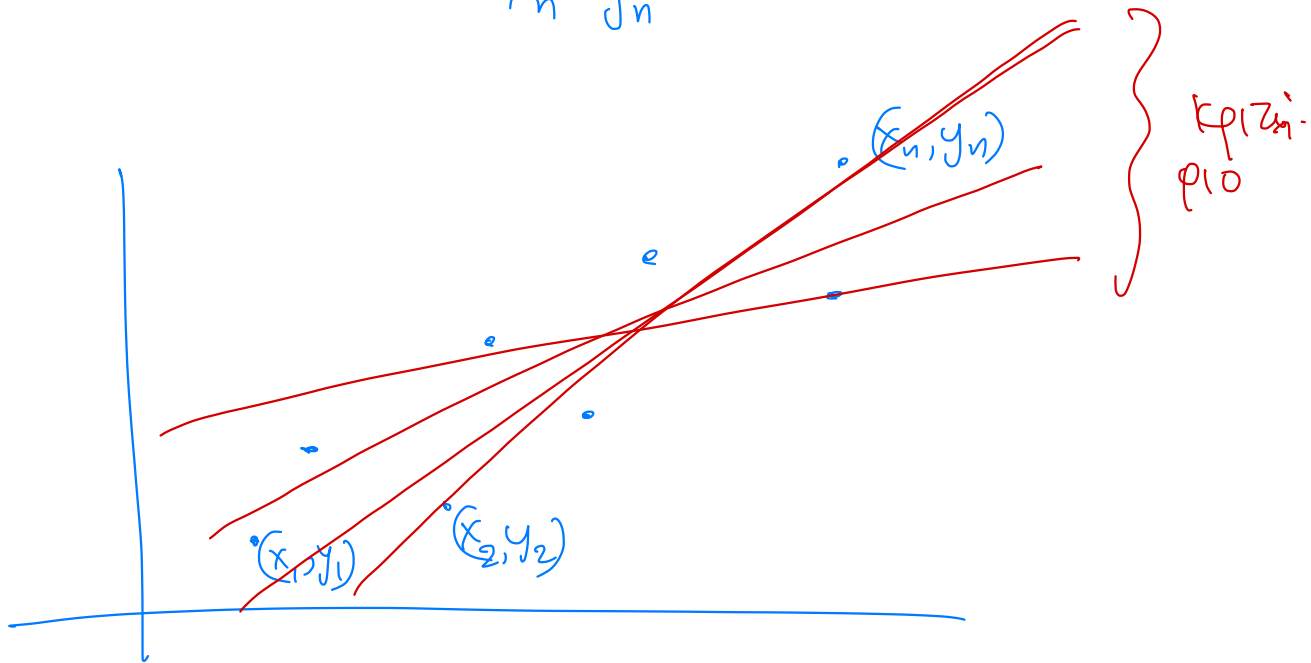
Y_j από διαφορετικές κανονικές π.κ. με $x_1 = 50$
 $x_2 = 80$



Εξέλιξη Παραμέτρων

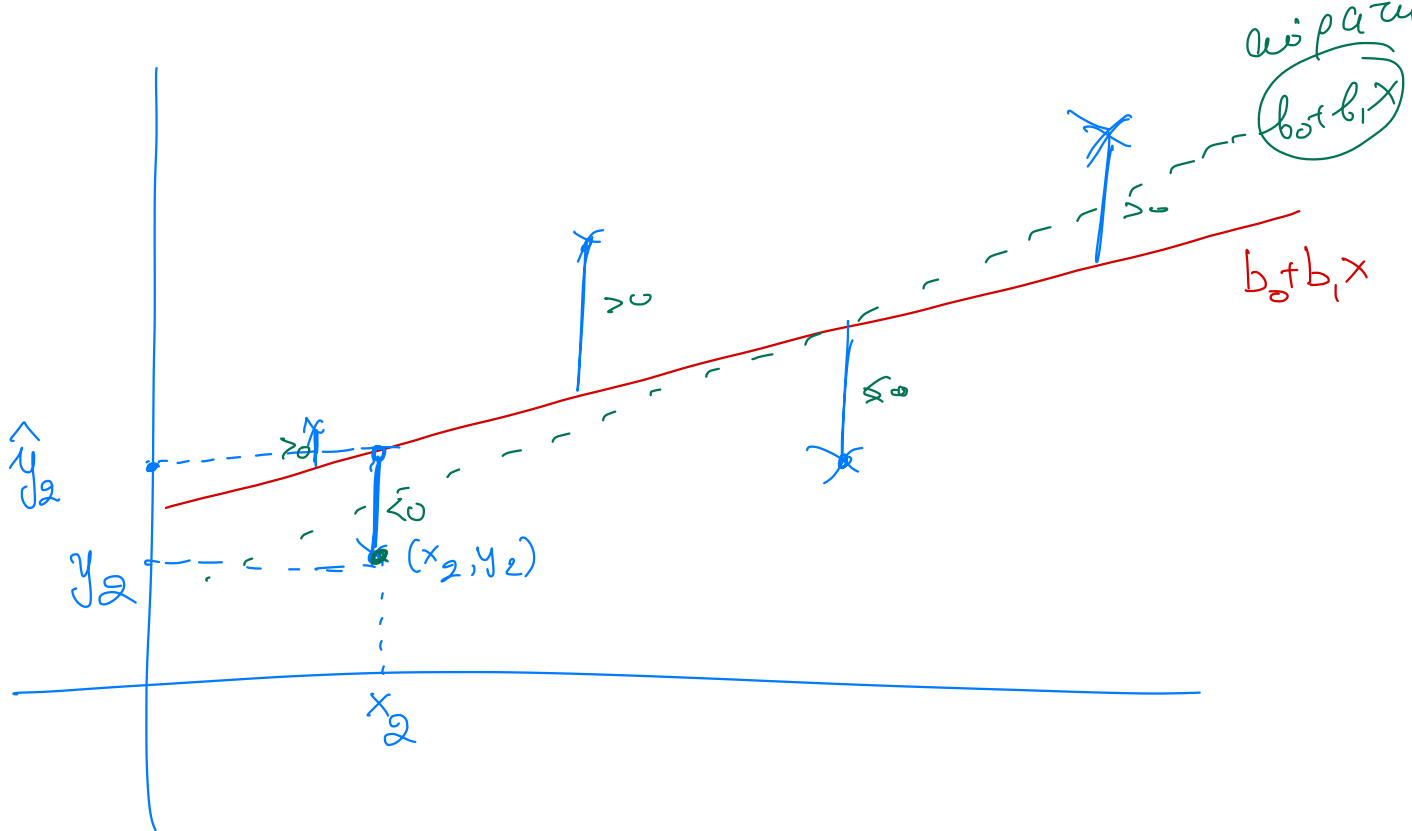
Δείγμα

x_1 y_1
 x_2 y_2
 \vdots \vdots
 x_n y_n



Κριτήριο Ελαχίστων Τετραγώνων

(\approx μέγιστη πιθανοφάνεια λόγω $\mathcal{N}(0, \sigma^2)$)



y_2 : πραγματική τιμή

$\hat{y}_2 = b_0 + b_1 \cdot x_2$: προβλεπόμενη τιμή της Y από το πρότυπο $y = b_0 + b_1 x$ για $x = x_2$

Απόκλιση της j : $(y_j - \hat{y}_j) = e_j$: residual (κατάλοιπο)

Κριτήριο τετραγωνικής απόκλισης

$$SSE(b_0, b_1) = \sum_{j=1}^n (y_j - \hat{y}_j)^2 = \sum_{j=1}^n (y_j - (b_0 + b_1 x_j))^2$$

(sum of square errors)

Θέλουμε είναι τα b_0, b_1 :

$SSE(b_0, b_1)$: ελάχιστο

$$\frac{\partial SSE}{\partial b_1} = 0, \quad \frac{\partial SSE}{\partial b_0} = 0 \Rightarrow \dots$$

⇒ ... ⇒

$\hat{b}_0, \hat{b}_1, \dots$ Εκτιμήσεις
Επιχ. παραμέτρων

$$SSE(\hat{b}_0, \hat{b}_1) \leq SSE(b_0, b_1) \quad \forall b_0, b_1$$

\hat{b}_0, \hat{b}_1 : τιμές

$$\hat{b}_1 = \frac{\sum x_j y_j - \frac{1}{n} (\sum x_j) (\sum y_j)}{\sum x_j^2 - \frac{1}{n} (\sum x_j)^2}$$

Παρατηρήσεις

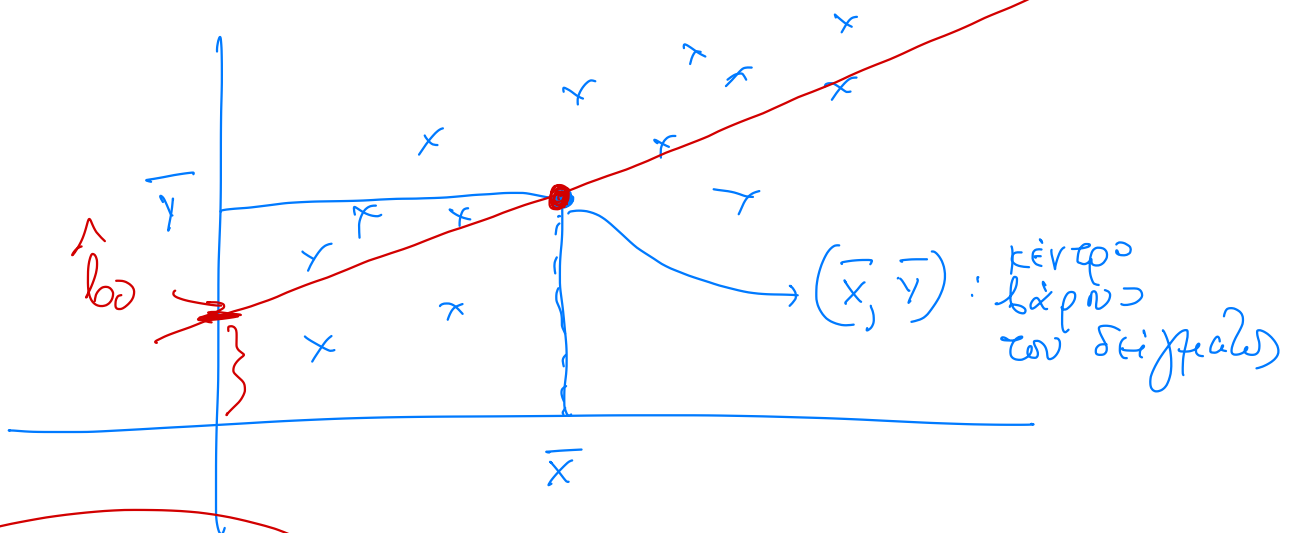
①

$$\hat{b}_0 + \hat{b}_1 \cdot \bar{X} = \bar{Y}$$

$$\bar{X} = \frac{x_1 + \dots + x_n}{n}$$

$$\bar{Y} = \frac{y_1 + \dots + y_n}{n}$$

$$\hat{b}_0 + \hat{b}_1 \cdot \bar{X}$$

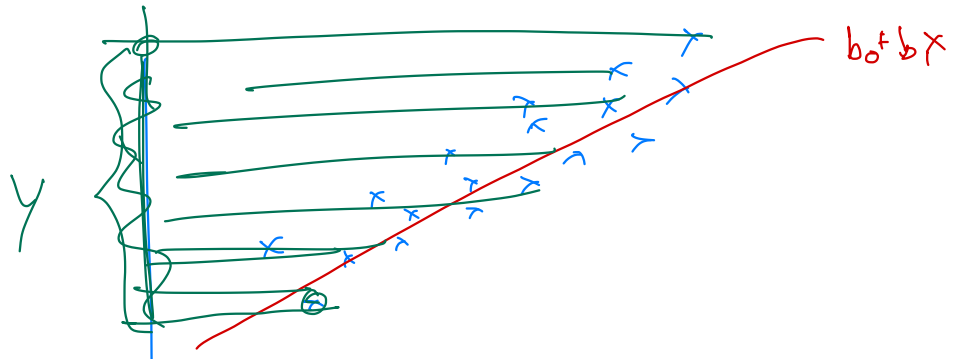


$$\hat{Y} = \hat{b}_0 + \hat{b}_1 \cdot X$$

: Εξίσωση παραμ. ευθείας με βάση δεδομένα

Παράδειγμα

① Δείγμα 1



② Δείγμα 2

