

2023-12-06

## Μοντέλο Παλινδρόμησης

$Y$  = εξαρτημένη μεταβλητή / απόκριση  
dependent variable / response variable

$X_1, X_2, X_3, \dots, X_p$  : ανεξάρτητες μεταβλητές / παράγοντες  
Independent variables / predictors / factors

$Y$  : ποσοτική μεταβλητή (scale)

$X_1, \dots, X_p$  : ποσοτικές αρχικά (θα το γενικεύσουμε)

---

## Στατιστικό Μοντέλο Παλινδρόμησης

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

(αγνωστή  
συνάρτηση)

συνάρτηση  
παλινδρόμησης  
(regression  
function)

π.χ.  $Y = \text{score}$

$X_1 = \text{age}$

$X_2 = \text{exercise time}$

$$E(Y | X_1 = x_1, X_2 = x_2) = f(x_1, x_2)$$

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

= μέσο όρος αρίθμων ηφελίας 52 ετών  
που ασκούνται 120 min ανά εβδομάδα

## Παρατηρήσεις

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

Δε υπάρχει συσχέτιση <sup>με  $Y$</sup>  με κάποια/ες από αυτές  
ως μεταβλητές

②  $f$ : άγνωστη συνάρτηση που πρέπει να εκτελεστεί  
από τα δεδομένα.

[στη ορολογία του machine learning  
θέλουμε να μάθουμε την  $f$  από δεδομένα]

π.χ.

①  $f(x_1, x_2) = a_0 + a_1 x_1 + a_2 x_2$

$a_0, a_1, a_2$  : άγνωστες παράμετροι

γραμμική παλινδρόμηση

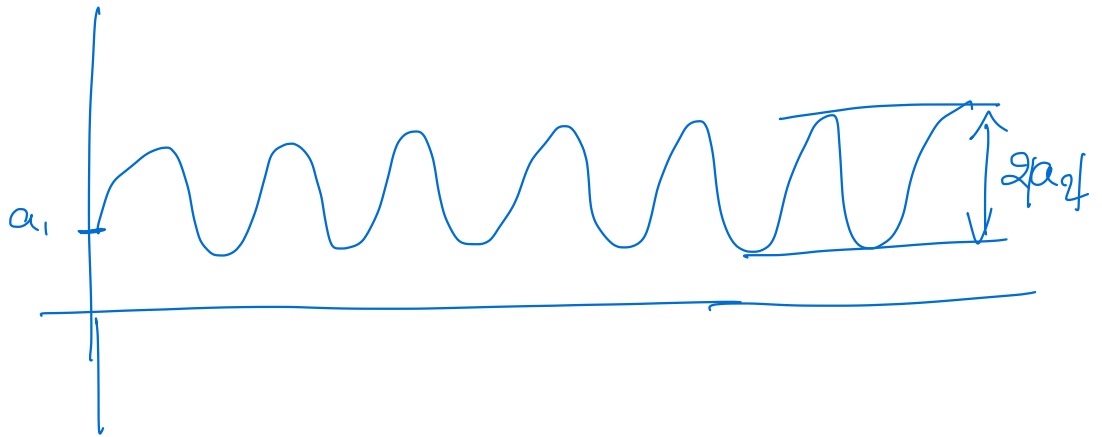
$f$  : γραμμική ως προς  $x_1, x_2$

$$\textcircled{2} \quad f(x_1, \dots, x_p) = C = \text{σταθερά}$$

$C$ : άγνωστο.

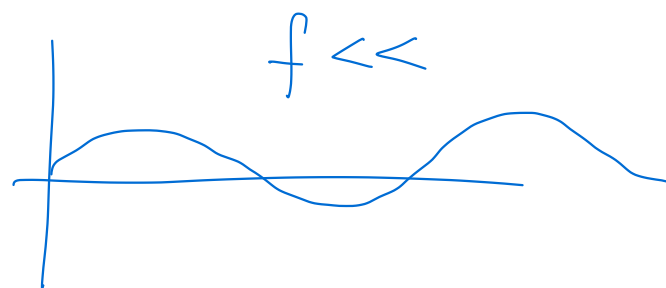
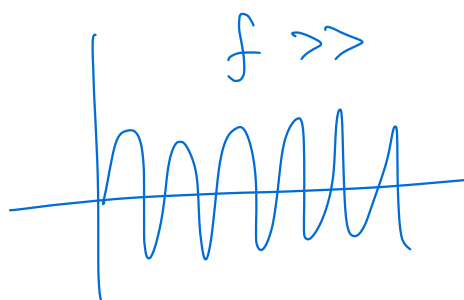
$$\textcircled{3} \quad C = E(Y) \Rightarrow \left. \begin{array}{l} \text{Εξίσωση του } C \\ \updownarrow \\ \text{Εξίσωση της μέσης} \\ \text{τιμής του } Y \text{ σε} \\ \text{όλο τον πληθυσμό} \end{array} \right\}$$

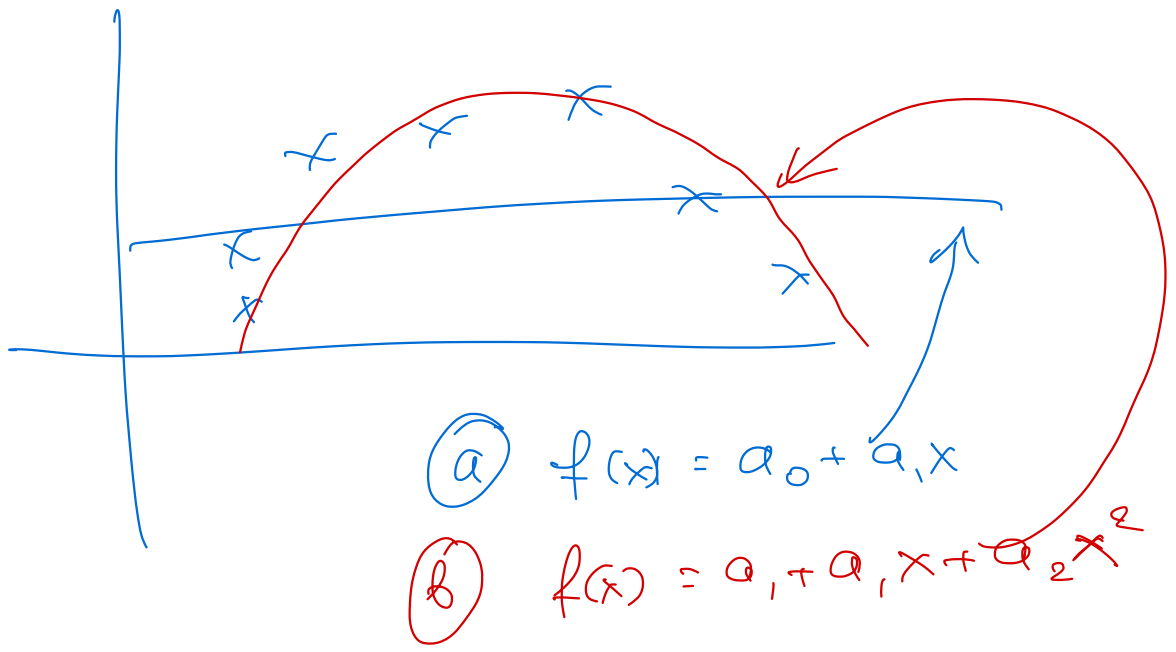
$$\textcircled{4} \quad f(x) = a_1 + a_2 \sin(2x)$$



$$\textcircled{5} \quad f(x) = a_1 + a_2 \cdot \sin(f \cdot x)$$

$a_1, a_2, f$ : άγνωστα.

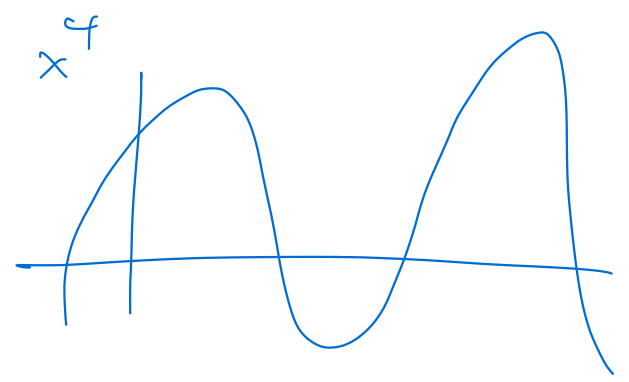
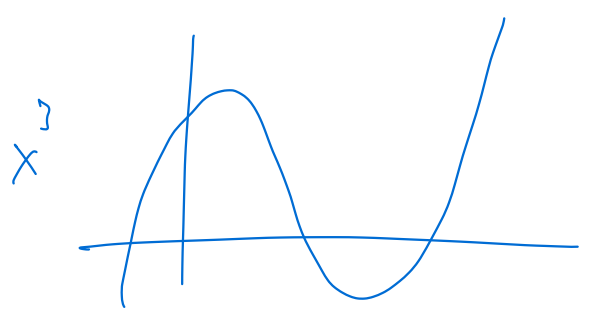




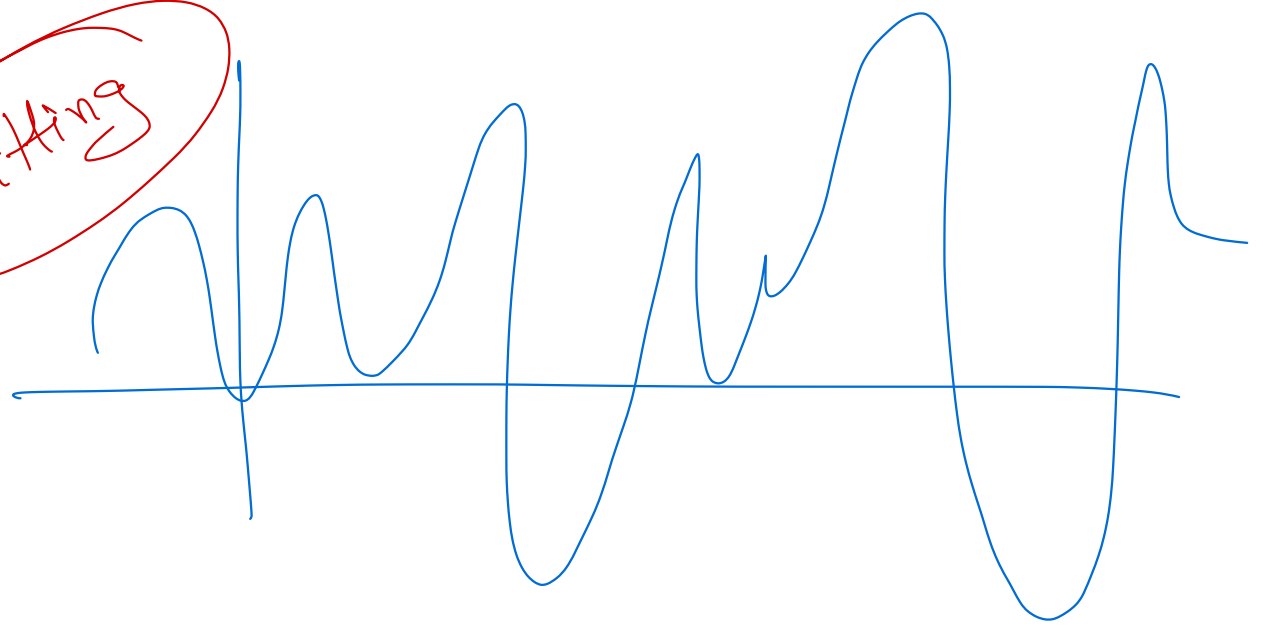
για να μη θεωρώ ες ακρίβ

$$f(x) = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots + a_{100}x^{100}$$

Πολυωνομική Λαβενδρόμια



overfitting



$$f(x_1, x_2, x_3) = a_0 + a_1 x_1 + a_2 x_1 x_2 + a_3 x_1^2 x_2 + a_4 x_1^3 x_2^2 \dots$$

Πολυωνομική παλινδρόμηση

$$f(x) = a_0 + a_1 \sin(2x) \quad \text{ημιτονοειδής}$$

$$f(x) = a_1 + a_2 e^{bx_1} \quad \text{εκθετική} \dots$$

## Γραμμικό Μοντέλο

$$E(y | x_1, \dots, x_p) = f(x_1, \dots, x_p) =$$

$$= b_0 + b_1 g_1(x_1, \dots, x_p) + b_2 g_2(x_1, \dots, x_p) + \dots + b_k g_k(x_1, \dots, x_p)$$

$b_0, b_1, \dots, b_k$  : άγνωστος παράμετροι

$g_1, g_2, \dots, g_k$  : γνωστές συναρτήσεις

Η  $f$  είναι γραμμική ως προς τις  
 άγνωστος παράμετρος

Γραμμικό μοντέλο

π.χ.  $f(x) = a_0 + a_1 \sin 2x$   
 $f_1(x)$   
 αγωρα

γραμμικό  
 μοτίβο  
 όχι γραμμική  
 μετασχηματισμός

$f(x) = a_0 + a_1 \sin(a_2 x)$   
 $a_0, a_1, a_2$ : αγωρα

όχι γραμμικό  
 μοτίβο

$f(x) = b_0 + b_1 e^{\delta_1 x} + b_2 e^{\delta_2 x}$

μη  
 γραμμικό

$\delta_1, \delta_2, b_0, b_1, b_2$   
 αγωρα

$f(x_1, x_2) = a_1 + a_2 x_1 + a_3 x_2 =$

γραμ. μοτίβο  
 γραμ. μετασχηματισμός

$f(x_1, x_2) = a_0 + a_1 x_1 + a_1^2 x_2$

όχι γραμ. μοτ.  
 γραμμική μετασχηματισμός

$f(x_1, x_2) = a_0 + a_1 x_1 + a_2 x_2$

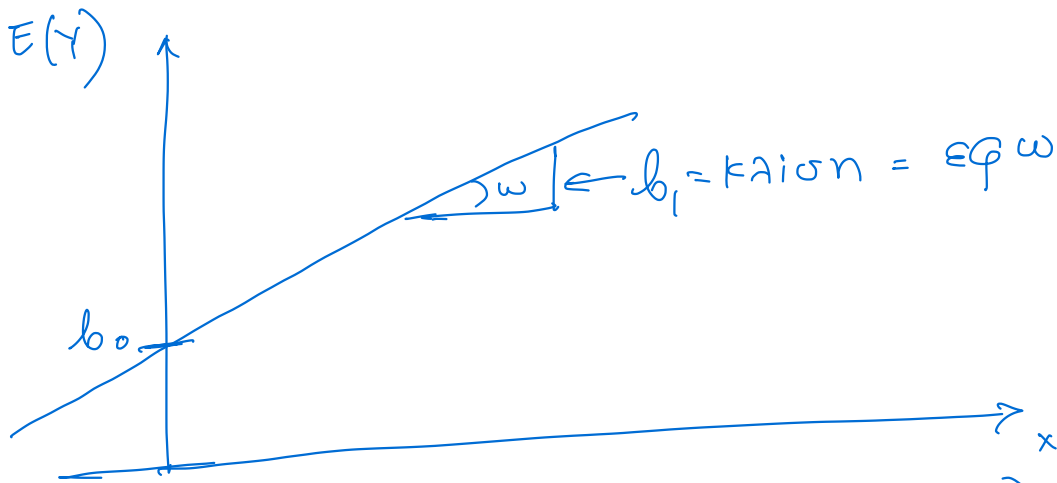
μη αβιοριστός  $a_2 = a_1^2$

# Άντιστοιχία Μονοπαραγοντικού Μοντέλου

$X$  : ανεξάρτητη μεταβλητή

$Y$  : εξαρτημένη μεταβλητή

Μοντέλο  $E(Y|X=x) = b_0 + b_1 x$   $b_0$  } άγνωστος  
 $b_1$  } παράμετρος



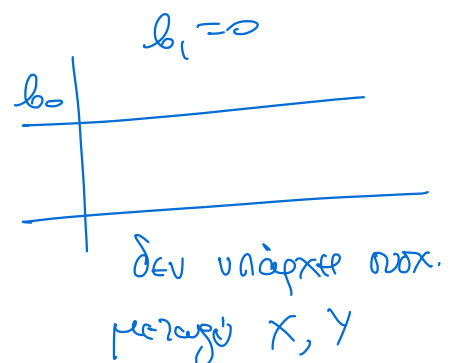
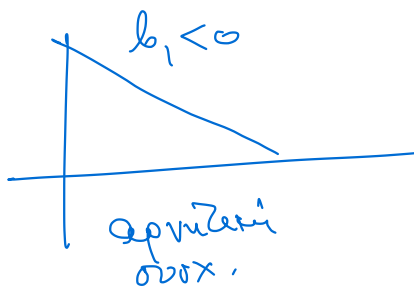
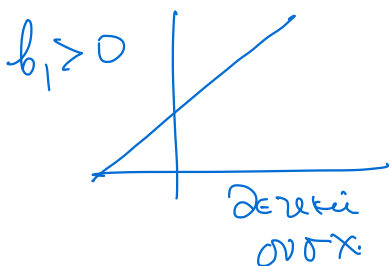
$b_0$  : σταθερός όρος (intercept)  
 $b_1$  : κλίση (slope).

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

$$E(Y|X=x+1) = b_0 + b_1 x + b_1$$

$b_1$  = μεταβολή της μέσης τιμής του  $Y$   
ανά μονάδα αλλαγής του  $X$ .

Ενίσχυση του  $X$

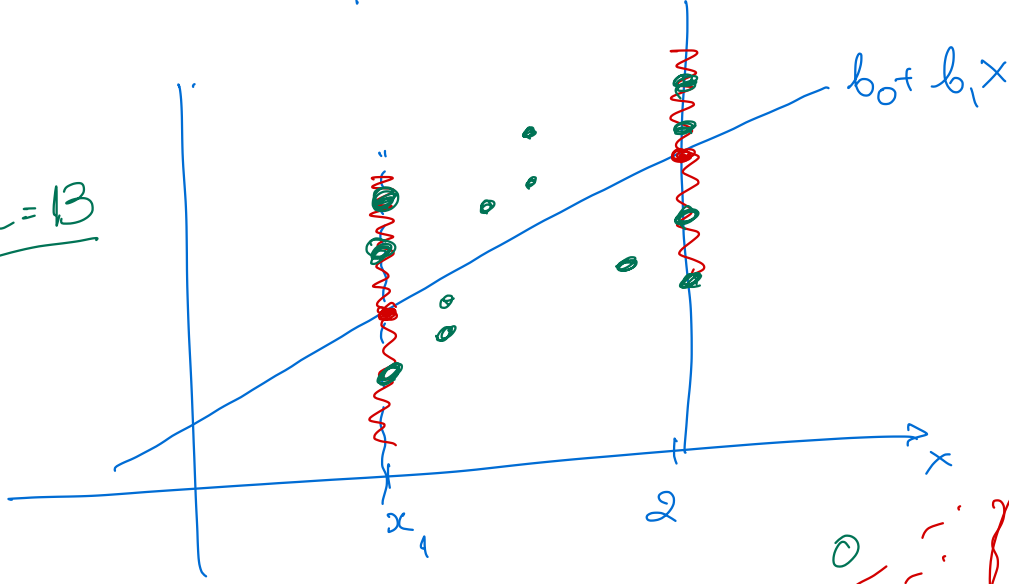


Μοντέλο  $E(Y|X=x) = b_0 + b_1 x$

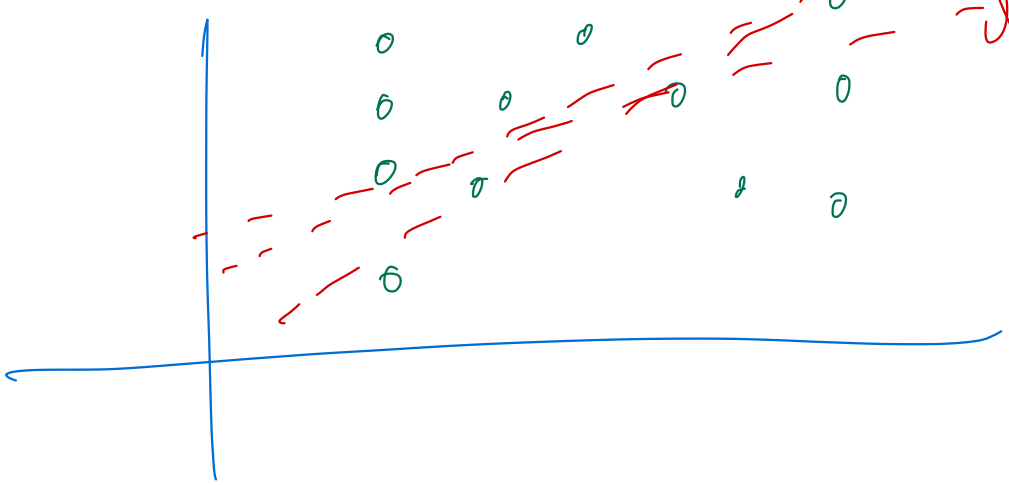
Ισοδύναμα  $Y = b_0 + b_1 x + \varepsilon$  , ε τυχαία μεταβλητή  
 $E(\varepsilon) = 0$ .

ε: τυχαία διακρίμανση / ανόρθωση

Δείγμα  
μεγέθους  $n=13$



ποια είναι  
n ενδείξεις  
 $b_0 + b_1 x$ ?



Δείγμα

X	Y
$x_1$	$y_1$
$x_2$	$y_2$
$\vdots$	
$x_n$	$y_n$

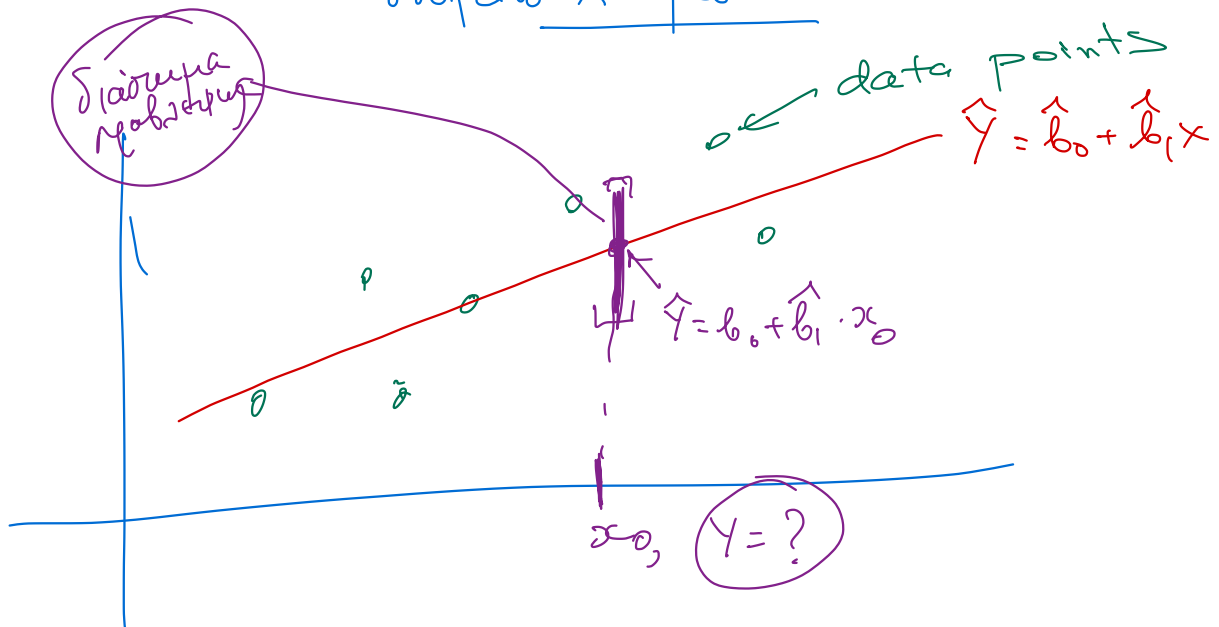
$j = 1, \dots, n$  παρατηρήσεις



① Από το δείγμα υπολογίζουμε  $\hat{b}_0, \hat{b}_1$  εκτιμήσεις  
 ως  $b_0, b_1$

$\hat{Y} = \hat{b}_0 + \hat{b}_1 x$  : εκτίμηση ως  
 συνάρτησης παραμόρφωσης

② Μπορούμε να χρησιμοποιήσουμε  $\hat{Y} = \hat{b}_0 + \hat{b}_1 x$   
 για να προβλέψουμε τιμές των  $Y$   
 για νέες παρατηρήσεις με  
 δεδομένα  $x$  μόνο



Στο πλαίσιο της μηχανικής μάθησης

μεθοδολογία

supervised learning (Επιβλεπόμενη μάθηση)

↓  
 στο δείγμα έχουμε και  $y$

# Υποθέσεις Παλινδρόμησης

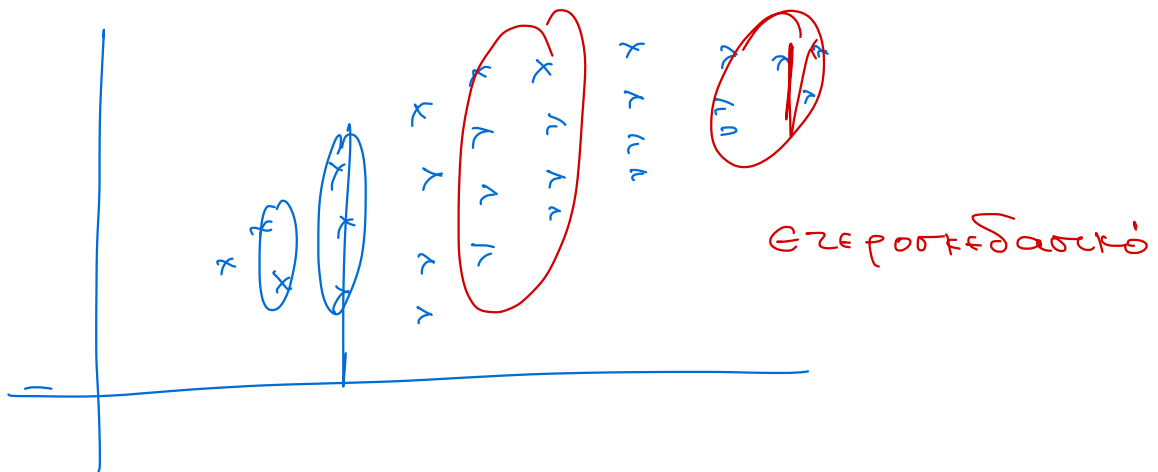
①  $Y_j = b_0 + b_1 x_j + \varepsilon_j \quad j=1, \dots, n$

②  $E(\varepsilon_j) = 0 \quad j=1, \dots, n$

③  $Var(\varepsilon_j) = \sigma^2 \quad \forall j$  (σταθερή / άγνωστη)

④  $\varepsilon_1, \varepsilon_2, \dots$  αυτοκίετα } οι παρατηρήσεις είναι αυτοκίετες μεταξύ τους

ομοσκεδαστικότητα



Η ④ παραβιάζεται π.χ. όταν οι παρατηρήσεις προκύπτουν από διαδοχικές χρονικές στιγμές (χρονοσειρά) ή από επαναλαμβανόμενες μετρήσεις στο ίδιο άτομο

↓  
longitudinal models

Με βάση τις υποθέσεις 1-4

μπορούμε να υπολογίσουμε ανεξάρτητα στατιστικές  $b_0, b_1$

ΔΕ? Ελάχιστος υποθέτω?

π.χ.  $H_0: b_1 = 0$   $H_1: b_1 \neq 0$

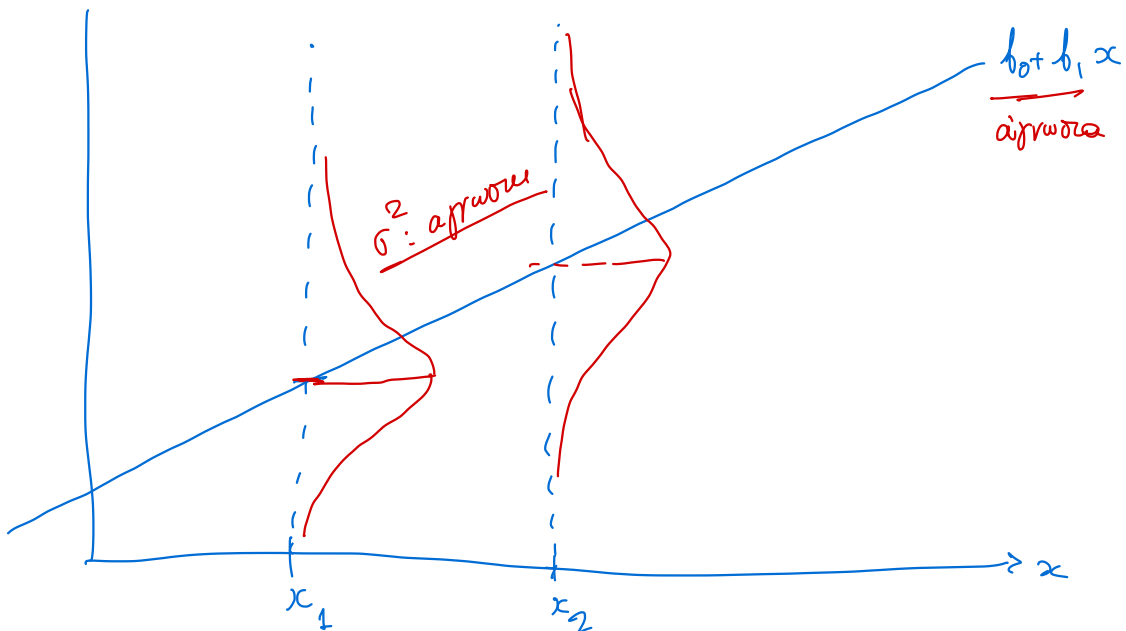
Χρησιμοποιούμε επιπλέον υπόθεση για την κατανομή του  $\varepsilon$

5)  $\varepsilon \sim N(0, \sigma^2)$  υπόθεση κανονικό τυχαίο

↑ όλα πραγματικά ποσά

Για άλλες κατανομές του  $\varepsilon \rightarrow$  Γενικευμένα Τετραγωνικά Ποσά

Για  $X=x \Rightarrow Y \sim N(b_0 + b_1 x, \sigma^2)$   
( $Y = b_0 + b_1 x + \varepsilon$ )

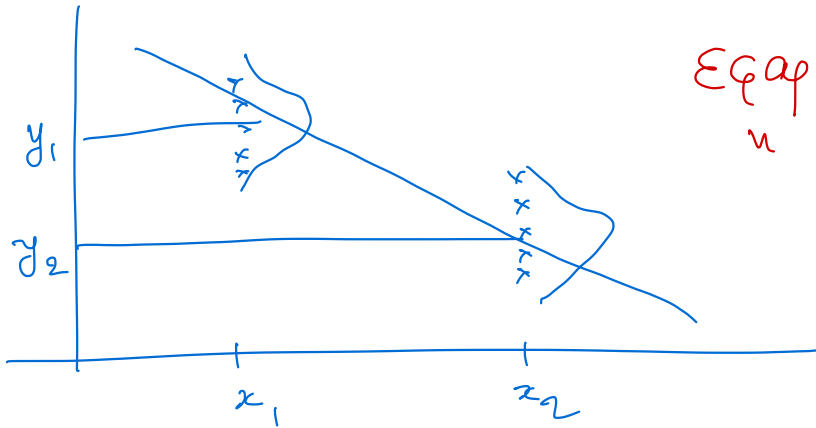


Όταν  $X=x$  = σταθερό  $Y|X=x$  ίδια κατανομή

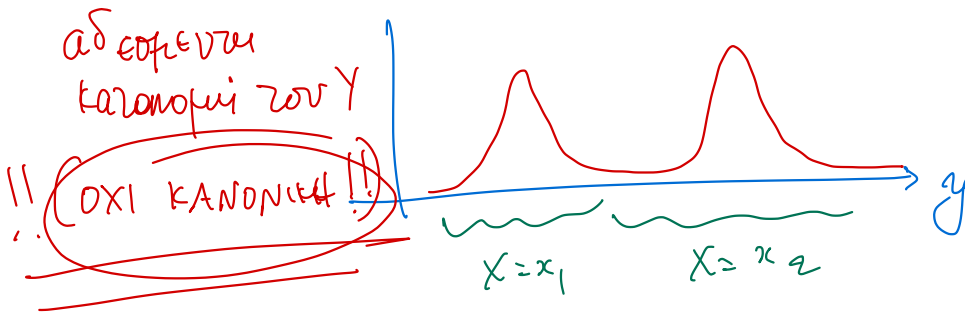
Για  $X_1 = x_1$ ,  $X_2 = x_2$

↓  
αίτην κατανομή      ↓  
αίτην.

Έστω  $X = x_1$  ή  $x_2$  σε στο 2ο δείγμα



εφαρμόζεται  
η αναμεταστροφή.



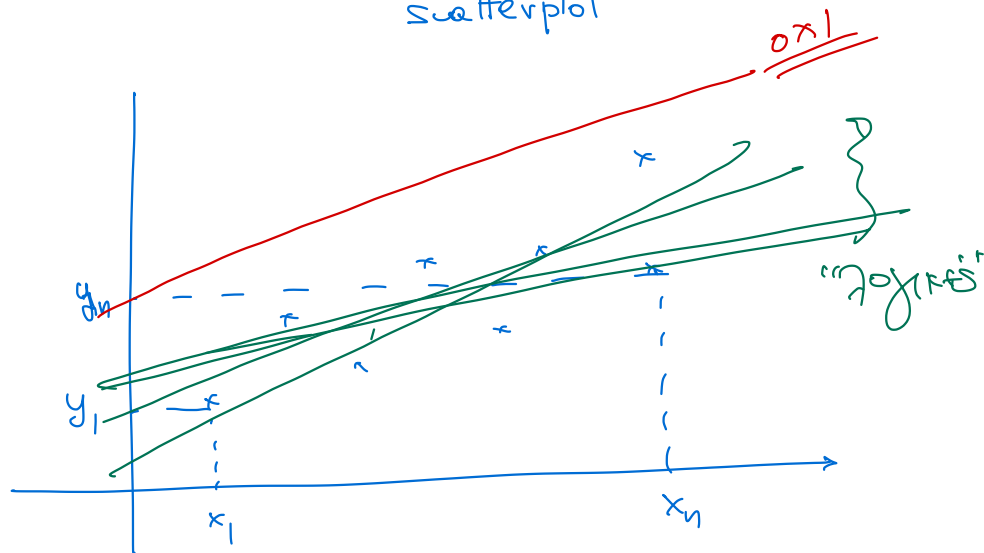
γενικό ισόδημα

Εκτίμηση Παραμέτρων στο Μονομεταβλητικό Μοντέλο

Δείγμα

x	y
$x_1$	$y_1$
...	...
$x_n$	$y_n$

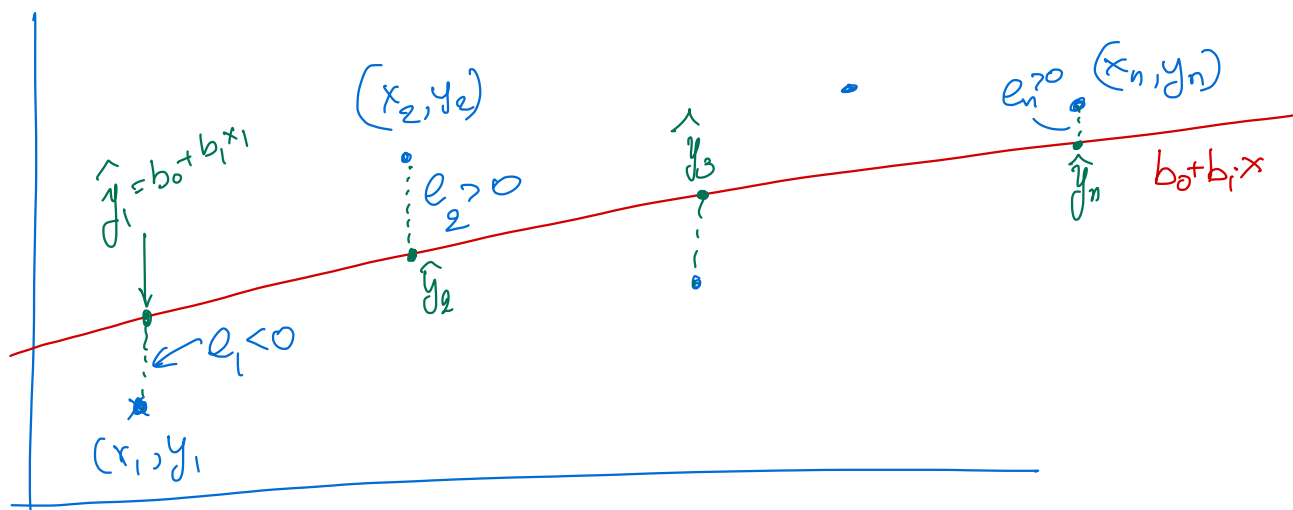
scatterplot



# Κριτήριο (και η Προσαρμογή)

## Κριτήριο Ελαχίστων τετραγώνων

[ισοδύναμο με  
μείνους τετραγώνων  
λόγω  $N(0, \sigma^2)$ ]



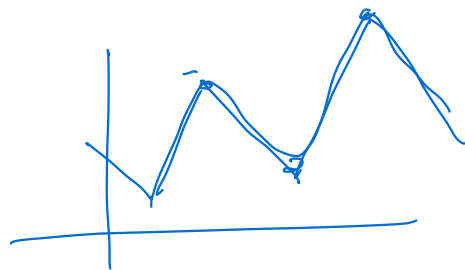
$\hat{y}_j$ : προβλεπόμενες τιμές από μόντελο  $b_0 + b_1 x$

$y_j$ : πραγματικές παρατηρήσεις

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

Κριτήριο 1 :  $\sum_{j=1}^n |e_j|$

δύσκολο μαθηματικό  
πρόβλημα



absolute regression

[μικ. μάλλον lasso]

Κριτήριο

$$\min \sum_{j=1}^n e_j^2$$

Μεθόδος  
ελαχίστων  
τετραγώνων

(least squares  
method)

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

↓  
Sum of  
Square errors

Θέλουμε ελαχίστα  $b_0, b_1$  έτσι ώστε

$$SSE(b_0, b_1) : \underline{\text{ελάχιστο}}$$

Αρκεί να αντιστοιχίσει

$$\frac{\partial SSE}{\partial b_0} = 0$$

$$\frac{\partial SSE}{\partial b_1} = 0$$

κανονικές  
εξισώσεις

$\Rightarrow \dots \Rightarrow$

$\left. \begin{array}{l} \hat{b}_0 \\ \hat{b}_1 \end{array} \right\}$  εκτιμήσεις ελαχίστων  
τετραγώνων των  $b_0, b_1$

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

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

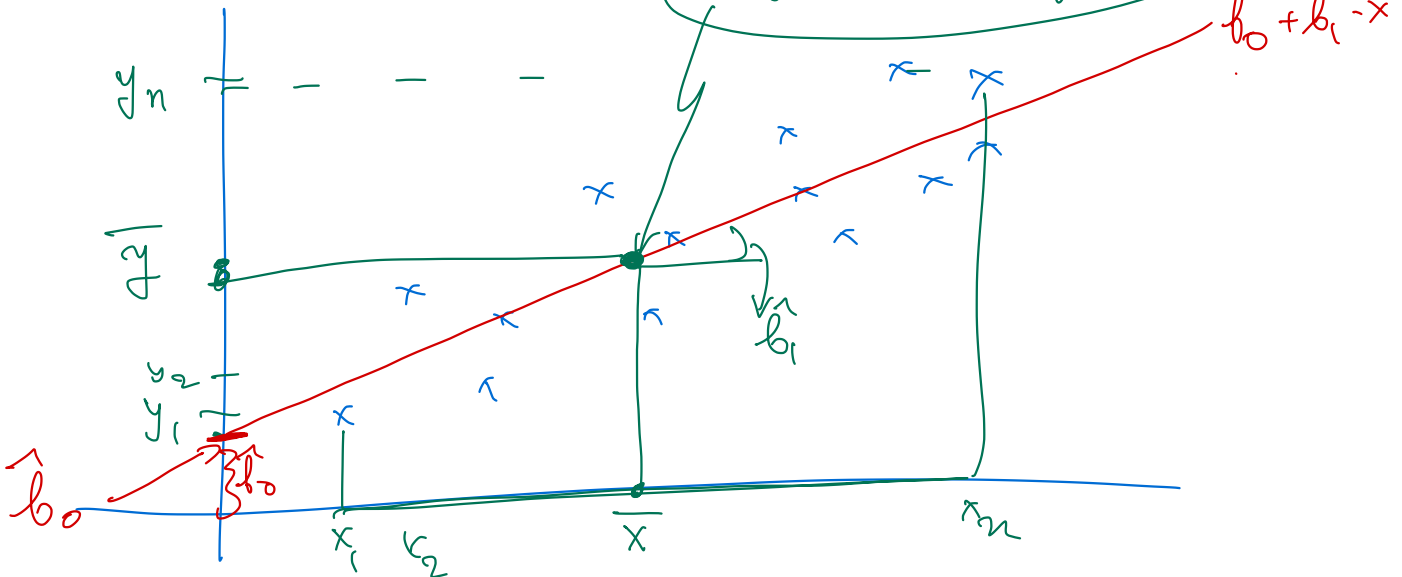
$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

$$\bar{y} = \frac{1}{n} \sum y_j$$

$$\bar{x} = \frac{1}{n} \sum x_j$$

$$\Rightarrow \hat{y} = \hat{b}_0 + \hat{b}_1 \cdot \bar{x}$$

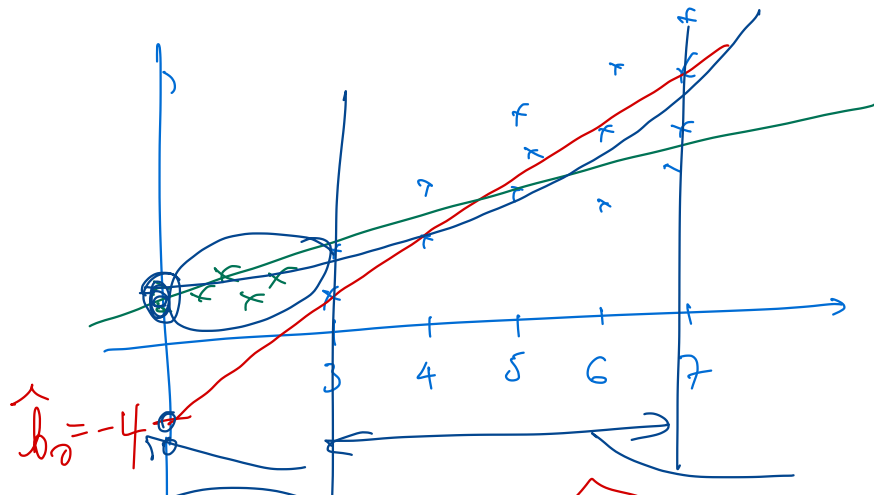
$(\bar{x}, \bar{y})$ : κέντρο βάρους των σημείων



### Ερμηνεία των $\hat{b}_0, \hat{b}_1$

$\hat{b}_1$ : κλίση : Για κάθε αύξηση του  $x$  κατά 1 μονάδα η  $E(y)$  αυξάνεται κατά  $\hat{b}_1$  κατά μέσο όρο (επίσης)

$$\hat{b}_0 = E(y | x=0)$$



$X = \text{age παιδιού}$   
 $Y = \text{βάρος}$   
 $n = 50$   
 $x \quad 3-7$

extrapolation  
 προβλεψη  
 για τιμές των  
 $x$  εκτός  
 των σημείων

$\hat{b}_0 = E(Y|X=0)$

