

28-2-2023

n. x. dataset exercise

$$Y = \text{score}$$

$$X = \text{age}$$

$$Y = b_0 + b_1 X + \varepsilon$$

R function lm : (linear model)

lm(model, . . .)

model : formula object ($y \sim x$)

$y \sim x$ \Leftrightarrow $y = \underline{b_0} + \underline{b_1}x$

analysis $\text{lm}(\text{score} \sim \text{age})$

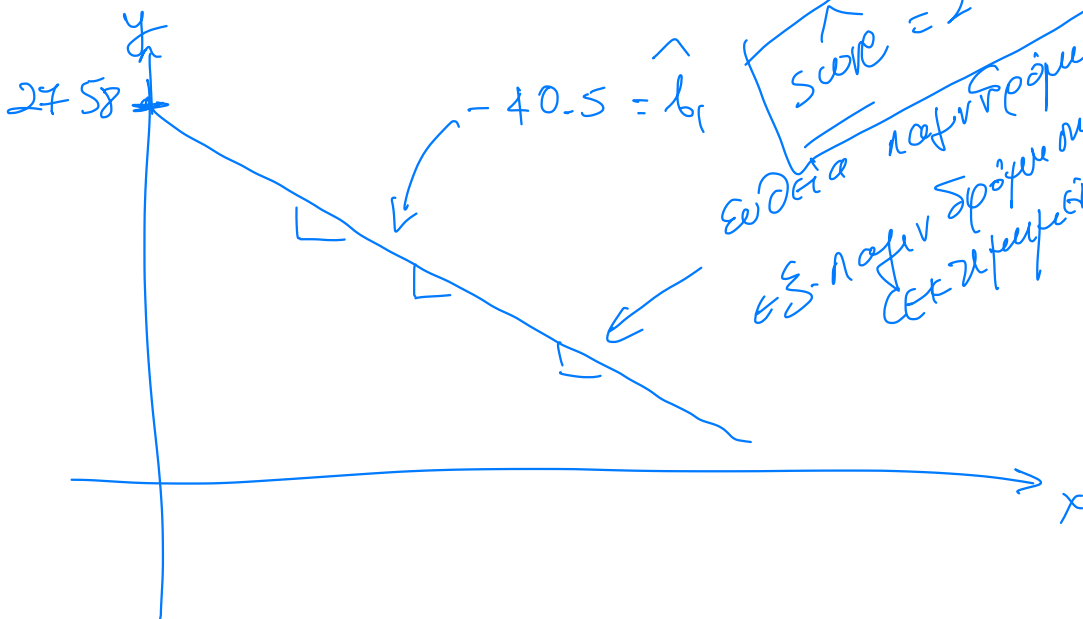
Εφαρμογή

model1 = score ~ age
lm(model1)

Παρατηρήσεις

$$\hat{b}_0 = 2758.22$$

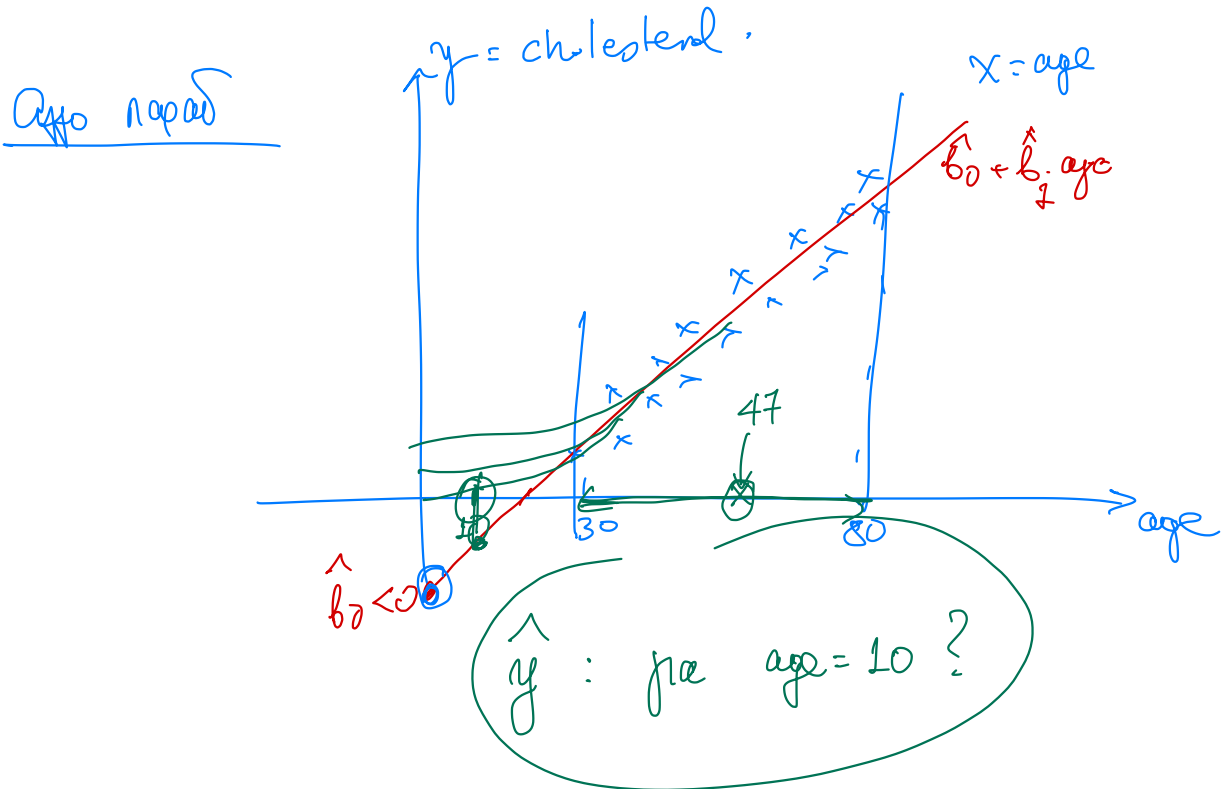
$$\hat{b}_1 = -40.51$$



① Ερμηνεία $\hat{b}_1 = -40.51$

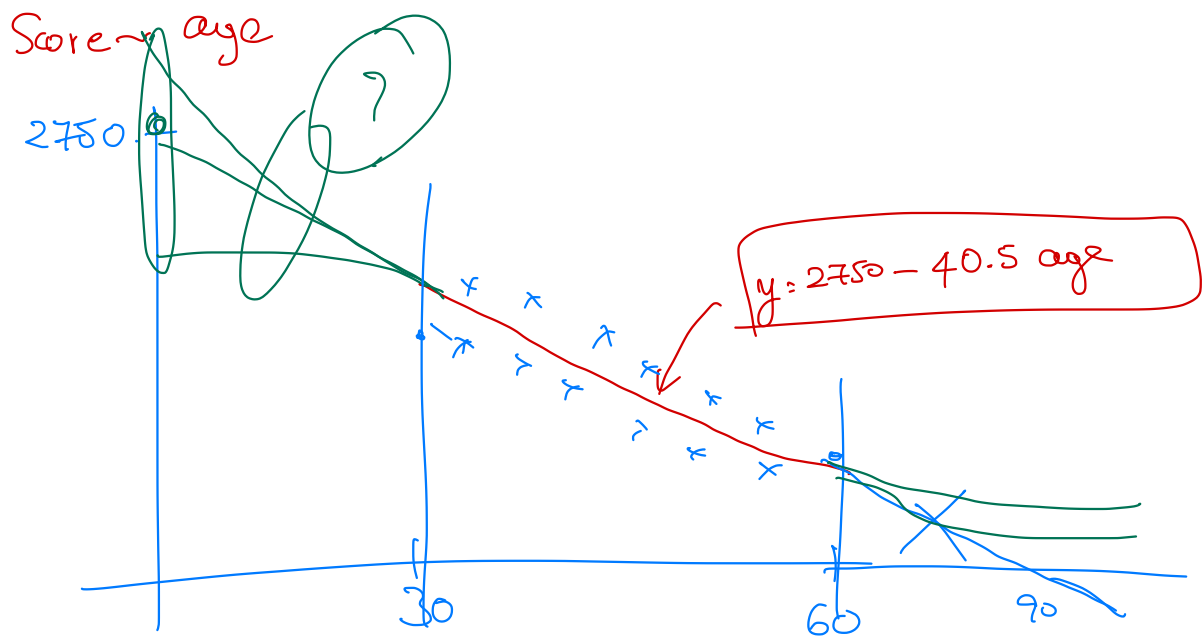
Για αύξηση της ηλικίας κατά 1 έτος
 η μέση τιμή των score μειώνεται κατά 40.51

Ερμηνεία $\hat{b}_0 = 2758$ (για $\text{age} = 0 \Rightarrow$
 $\hat{\text{Score}} = 2758$)



Πρόβλεψη του $E(y)$ για x εκτός των άμεσα διαθέσιμων σημείων του δείγματος \Rightarrow **extrapolation** X X !!

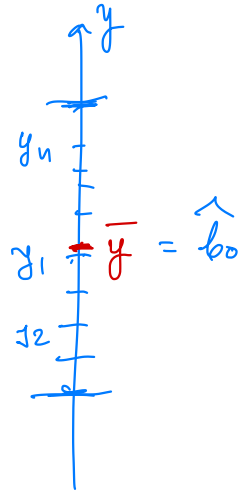
\hat{b}_0 : Έχει θετική ερμηνεία μόνο όταν υπάρχουν σημεία $x=0$ (ή και από 0) στο δείγμα.



Επιμέτρηση απροσδιοριστίας
(Ανάγκη διασποράς & Εξαρτητικότητας)

① Ανάλυση διασποράς (διαχωρισμός) (ANOVA)

$$Y = b_0 + \varepsilon \quad \left. \begin{array}{l} \varepsilon \sim \mathcal{N}(0, \sigma^2) \end{array} \right\} \Leftrightarrow$$



$$\Leftrightarrow Y \sim \mathcal{N}(b_0, \sigma^2)$$

b_0, σ^2 άγνωστα

$$\hat{b}_0 = \bar{y}$$

(MLE - LSE)

$$\hat{\sigma}^2 = \frac{SST}{n-1} \sim \chi^2_{n-1}$$

Y εξαρτ. μεταβλητή
δείγμα

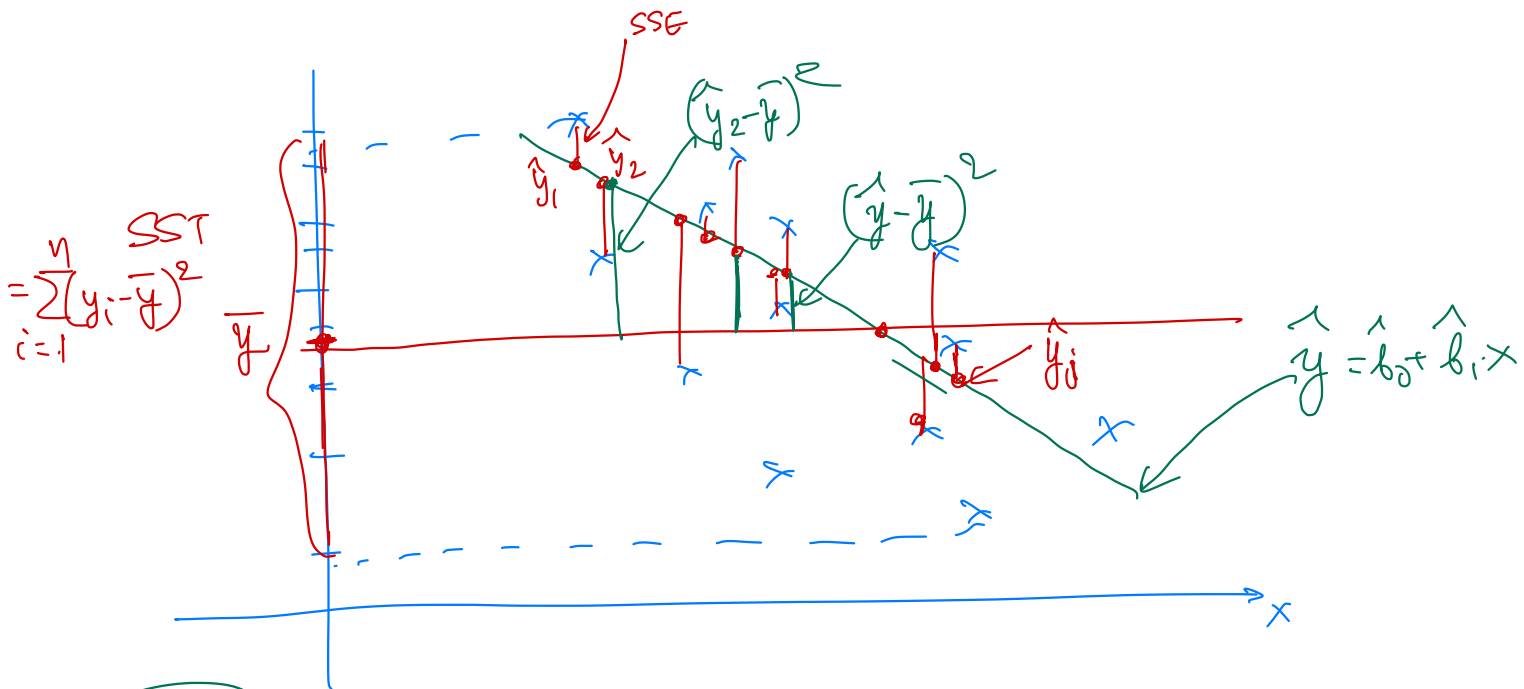
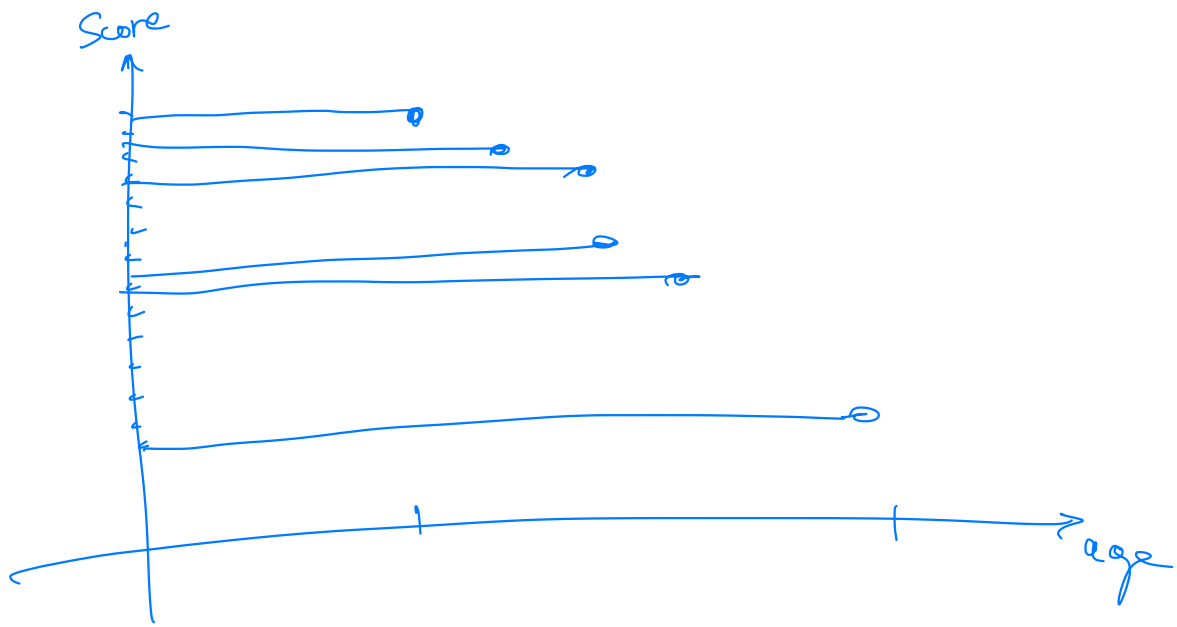
$$(y_1, y_2, \dots, y_n)$$

Διασπορά

$$\frac{1}{n-1} \sum (y_i - \bar{y})^2$$

συνολική διακύμανση εν Y στο δείγμα

Sum of Squares total (SST)



$$SSR = \sum (\hat{y}_i - \bar{y})^2$$

Or $SSR \approx 0 \Rightarrow$ a good model. $\sigma^2 \rightarrow 0$ or $\sigma^2 \rightarrow \infty$.

(εσ η μέρη μεταβλητά του \hat{y} και το X) = Sum of squares regression

$$SSR = SS_{\text{model}}$$

$$SSE = \sum_{j=1}^n (y_j - \hat{y}_j)^2 = \text{Sum of Squares error} = SS_{er}$$

(αν $SSE \approx 0 \Rightarrow$

μεταβλ. των Y στο δείγμα που παρατηρείται
αποτελείται από το μοντέλο $Y = b_0 + b_1 X$

$$SST = \sum_{j=1}^n (y_j - \bar{y})^2 \quad \text{συνολική}$$

$$SSR = \sum_{j=1}^n (\hat{y}_j - \bar{y})^2 \quad \text{εξηγημένη}$$

$$SSE = \sum_{j=1}^n (y_j - \hat{y}_j)^2 \quad \text{ανεξήγητη.}$$

ΘΕΩΡΗΜΑ Όταν τα \hat{b}_0, \hat{b}_1 εκτιμούνται με βάση LSE

τότε ισχύει

$$\boxed{SST = SSR + SSE}$$

ή ως αναγνώριση διακρίσεων
Analysis of Variance formula.

Παρένθεση

Ενώ από το κριτήριο Ελαχ. Τετρ.

θα μπορούσαμε π.χ. να ελαχιστοποιήσουμε

$$= \sum_{j=1}^n |y_j - (b_0 + b_1 x_j)| \quad \begin{array}{l} \text{συνολική} \\ \text{απόλυτη απόσταση.} \end{array}$$

b_0^*, b_1^* absolute distance estimates

SSR, SSE for choice b_0^*, b_1^*

$$SST = SSR + SSE$$

LSE :

$$SST = \underbrace{SSR} + \underbrace{SSE}$$

αντ. προσδιορίζεται $R^2 = \frac{SSR}{SST}$ = % μεταβολών Y σε σχέση που εξηγείται από το μοντέλο $Y = b_0 + b_1 X$

1) $0 \leq R^2 \leq 1$

2) $R^2 = r^2$, όπου $r = \text{Corr}(x, y)$
(στο μονοπαραγοντικό μοντέλο $Y = b_0 + b_1 X$)

$$SSR = 10255910$$

$$SSE = 1485592$$

$$SST = 11741503$$

$$R^2 = \frac{SSR}{SST} = 0.873$$

mean squares

n = sample size

Πινακας ANOVA		Degrees of freedom df	MS
	SS		
X (Model)	SSR	1	$MSR = \frac{SSR}{df_{reg}}$
Error	SSE	n-2	$MSE = \frac{SSE}{df_{error}}$
Total	SST	n-1 (df _{total})	

$$Y = b_0 + b_1 X + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

$$df_{er} = n - 2$$

$$\frac{SSE}{df_{er}} = \frac{SSE}{n-2} = \hat{\sigma}^2$$

απεριόριστη

$$\sum_{j=1}^n (y_j - b_0 - b_1 x_j)^2$$

απώττα

$$\sum_{j=1}^n (y_j - \hat{b}_0 - \hat{b}_1 x_j)^2 = SSE$$

2 απαιτήσεις
αυξη. και ελαττώσεως
προς axis το βλγτα

$$df_{er} = n - 2$$

$$MSE = \frac{SSE}{n-2} = \hat{\sigma}^2$$

απερ. εστ. του σ^2

$$\sim \chi^2_{n-2}$$

Βαθμοί Ελευθερίας

Στο γενικό γραμμικό μοντέλο

$$Y = b_0 + b_1 X_1 + \dots + b_p X_p + \varepsilon$$

απ. αγνώστων παραμέτρων = $p+1$

$$df_{er} = n - (\# \text{ β στο μοντέλο}) = n - (p+1)$$

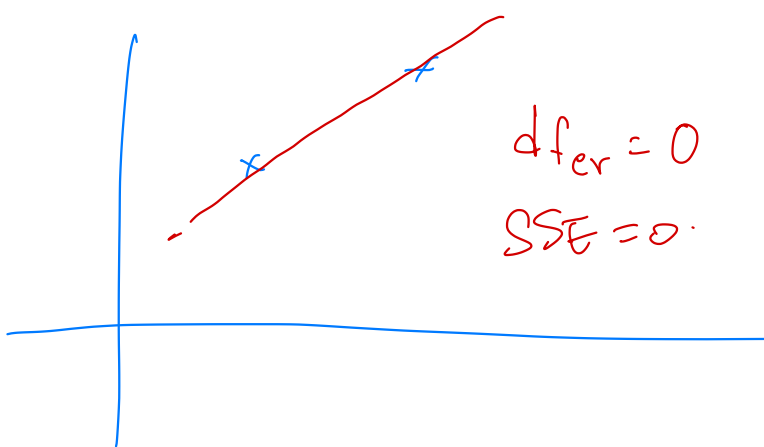
$$df_{reg} = (\# \text{ β στο μοντέλο}) - 1 = p$$

$$df_{er} + df_{reg} = n - 1 = df_{total}$$

$$SSE + SSR = SST$$

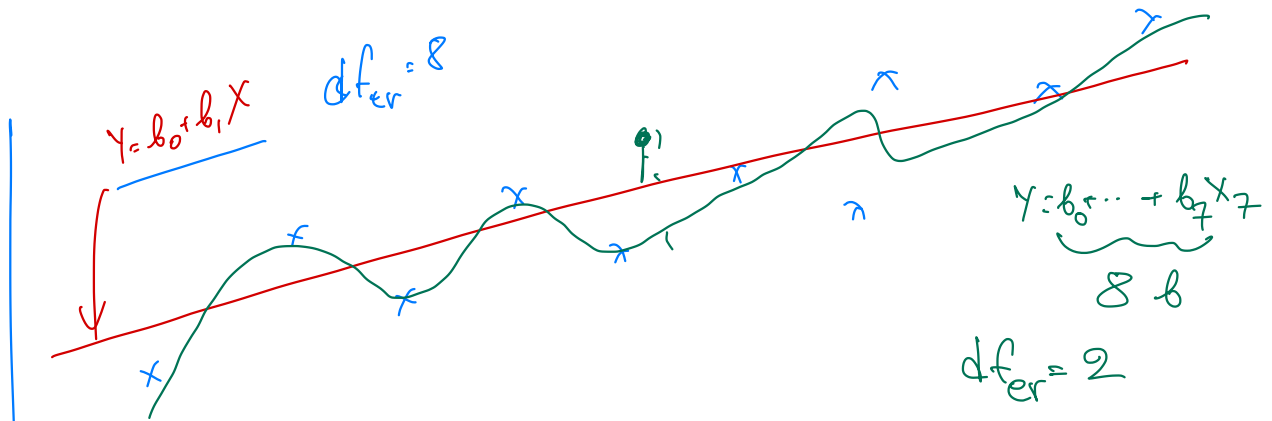
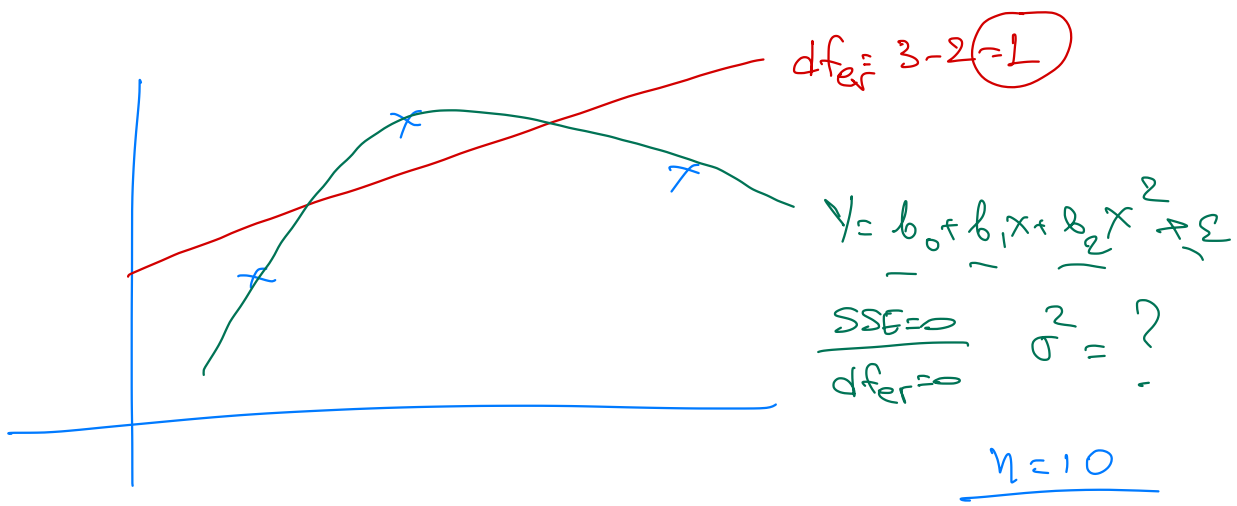
$df_{error} = n - (\# \text{ β})$ = "υποκατάσταση" ("ισοδυναμία") του μεγέθους διαγράμματος

Γενικά $\hat{\sigma}^2 = MSE = \frac{SSE}{df_{er}}$

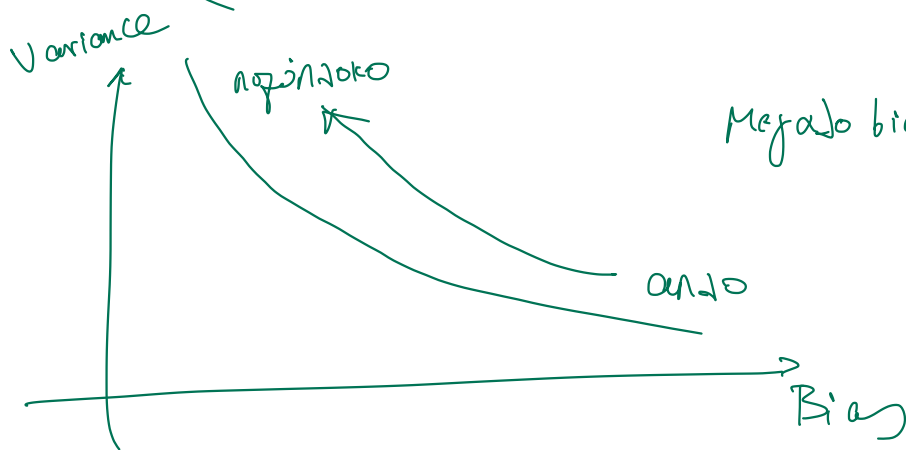


$$df_{er} = 0$$
$$SSE = 0$$

$\frac{n=2}{\hat{\sigma}^2 = \frac{0}{0}}$ (?)



ML: (Bias vs Variance tradeoff)



Mejor de bias: "kano" fit

Εκτιμώμεται για us παραμέτρους

\hat{b}_1 : εκτίμηση των $b_1 \sim \mathcal{N}(b_1, \sigma_{\hat{b}_1}^2)$

\hat{b}_0 : " των $b_0 \sim \mathcal{N}(b_0, \sigma_{\hat{b}_0}^2)$

① \hat{b}_0, \hat{b}_1 : ανεξάρτητες εκτιμήσεις των b_0, b_1

$$E(\hat{b}_0) = b_0$$

$$E(\hat{b}_1) = b_1$$

Διότι έχουμε



Δείγμα 2
 $\hat{b}_{1,2}$

Δείγμα N
 $b_{1,N}$

με το ίδιο n
ε' 2α ίδια x_1, \dots, x_n

$$\frac{1}{N} (\hat{b}_{1,1} + \dots + \hat{b}_{1,N}) \rightarrow b_1$$

Τις ΔΕ για b_1 :

$$\hat{b}_1 : E(\hat{b}_1) = b_1$$

$S_{\hat{b}_1}^2$: δεγ. διασπορά \hat{b}_1

$$\left(\text{απόσ.} \right) \sigma^2(\bar{y}) = \frac{s^2}{n}$$

$$S_{\hat{b}_1} = \sqrt{S_{\hat{b}_1}^2} = \text{ζων. σφάλμα των } \hat{b}_1$$

$$\frac{\hat{b}_1 - b_1}{S_{\hat{b}_1}} \sim t_{df, \epsilon}$$

} \Rightarrow ΔΕ

$$\frac{\hat{b}_0 - b_0}{S_{\hat{b}_0}} \sim t_{df, \text{er.}}$$

}

$$\Delta E b_1 : \left(\hat{b}_1 - t_{\alpha/2, df_{er}} \cdot S_{\hat{b}_1} \leq b_1 \leq \hat{b}_1 + t_{\alpha/2, df_{er}} \cdot S_{\hat{b}_1} \right)$$

$$\Delta E b_0 : \hat{b}_0 \dots \leq b_0 \leq \hat{b}_0 \dots$$

$$H_0 : b_1 = 0 \quad H_a : b_1 \neq 0$$

t-test, $t = \frac{\hat{b}_1 - 0}{S_{\hat{b}_1}}$

Ελεγχος
Συμμετατόνωσης

accept H_0 : αν $|t| \leq t_{\alpha/2, df_{er}} \quad p \geq \alpha$
 reject H_0 : α $|t| > t_{\alpha/2, df_{er}} \quad p < \alpha$

$p = p\text{value}$