

2023-12-15

Μοντέλο Πολλαπλής Παλινδρόμησης

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_k X_k + \varepsilon$$

$\varepsilon \sim \mathcal{N}(0, \sigma^2)$, ανεξάρτητα
μεταξύ παρατηρήσεων

$$Y | X_1, X_2, \dots, X_k \sim \mathcal{N}(b_0 + b_1 X_1 + \dots + b_k X_k, \sigma^2)$$

Δείγμα

j	X_1	X_2	...	X_k	Y
1	X_{11}	X_{21}	...	X_{k1}	Y_1
2	X_{12}	X_{22}	...	X_{k2}	Y_2
⋮	⋮				
n	X_{1n}	X_{2n}	...	X_{kn}	Y_n

X_1, \dots, X_k → διαφορετικές μεταβλητές
→ συσχέτισης αλληλ μεταβλητών

n.x. ①
$$y = b_0 + b_1 \cdot \underset{x_1}{\text{age}} + b_2 \cdot \underset{x_2}{\text{height}} + b_3 \cdot \underset{x_3}{\text{weight}}$$

②
$$y = b_0 + b_1 \cdot \underset{x_1}{\text{age}} + b_2 \cdot \underset{x_2 = x_1^2}{\text{age}^2}$$

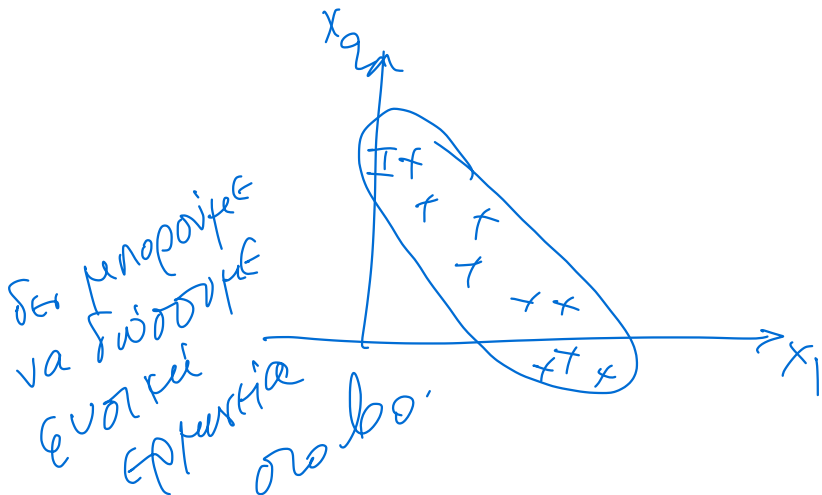
③
$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 \underset{x_3 = x_1 x_2}{x_1 x_2}$$

Ερμηνεία

$$y_0 = b_0 + b_1 x_1 + b_2 x_2 + \varepsilon$$

① $b_0 = E(y | x_1=0, x_2=0)$

εχει ερμηνεία σαν
 στο δείγμα υπάρχουν
 παρατηρήσεις με $x_1 = x_2 = 0$



$$\textcircled{2} \quad y_0 = b_0 + b_1 x_1 + b_2 x_2$$

b_1 = Μέση μεταβολή της Y αν η x_1 αυξηθεί κατά 1 μονάδα, κ' οι υπόλοιπες μετα-
βλητές παραμείνουν σταθερές

π.χ. $x_1 = \text{age}^{(\text{έτη})}$, $x_2 = \text{height}$

b_1 : επίρροη στο Y της αύξησης αυτής κατά 1 έτος, κρατώντας ύψος = σταθ.

$$y = b_0 + b_1 x_1 + b_2 x_2$$

$$b_1 = \frac{\partial y}{\partial x_1} \leftarrow$$

$$\textcircled{3} \quad y = b_0 + b_1 x + b_2 x^2$$

\uparrow
 x_1

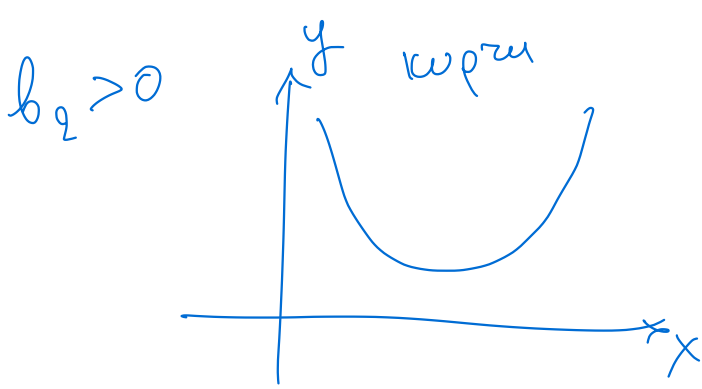
\uparrow
 x_2

δεν είναι δυνατό να αυξηθεί το x_1 κατά 1 κ' το x_2 να μείνει σταθερό

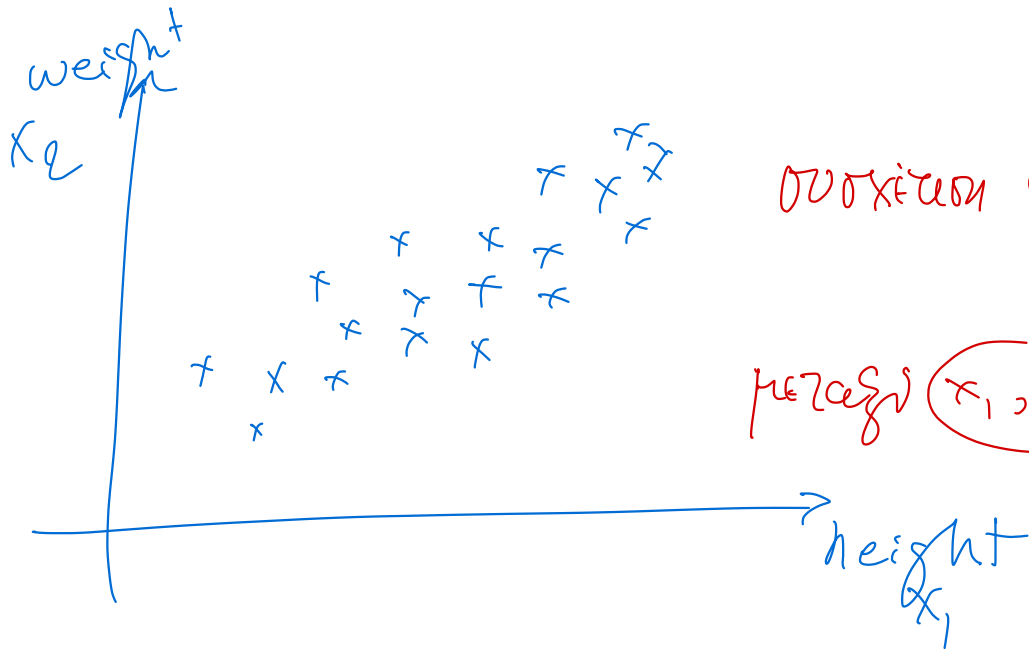
$$\frac{\partial y}{\partial x} = b_1 + 2b_2 x$$

b_1 δεν έχει φυσική ερμηνεία

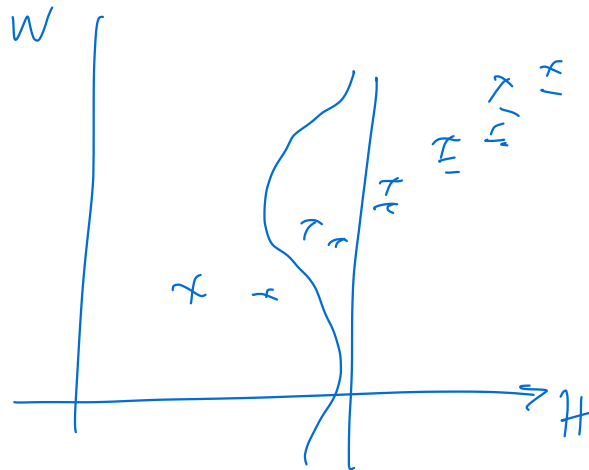
b_2 επίρροη

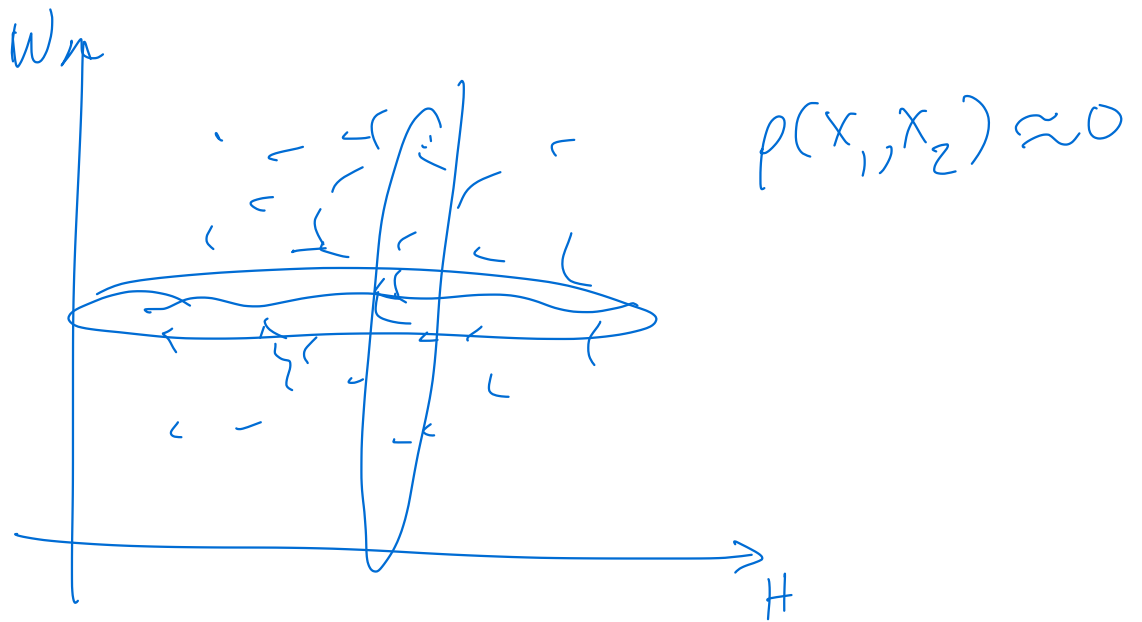


Ⓐ
$$Y = b_0 + b_1 \cdot \text{height} + b_2 \cdot \text{weight} + \varepsilon$$



α.ν $\rho(x_1, x_2) \approx 1$





Μοντέλο $Y_0 = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k + \Sigma$

b_0, b_1, \dots, b_k : άγνωστες παραμέτρους ($k+1$)

Εκτίμηση ελαχ. τετραγώνων

$$\hat{y} = \hat{b}_0 + \hat{b}_1 + \dots + \hat{b}_k x_k \leftarrow \text{προβλεπόμενες τιμές}$$

$$SSE(b_0, b_1, \dots, b_k) = \sum_{j=1}^k \left(y_j - \underbrace{b_0 + b_1 x_{1j} + b_2 x_{2j} + \dots + b_k x_{kj}} \right)^2$$

$\min_{b_0, b_1, \dots, b_k} SSE \Rightarrow \dots \hat{b}_0, \hat{b}_1, \dots, \hat{b}_k$: Εκτιμήσεις ελαχ. τετραγώνων

Ανάλυση Διασποράς

$$SST = \sum_j (y_j - \bar{y})^2 \quad (\text{συνεξ. ποσότητας})$$

$$SSR = \sum_j (\hat{y}_j - \bar{y})^2 : \text{μεταβλητότητα του } Y \text{ που εξηγείται από το μοντέλο}$$

$$SSE = \sum_j (y_j - \hat{y}_j)^2 : \text{μεταβλητότητα που παραμένει ανεξήγητη}$$

$$SST = SSR + SSE$$

$$R^2 = \frac{SSR}{SST} = \% \text{ μεταβ. του } Y \text{ που εξηγείται από το μοντέλο}$$

Εκτιμώμενη για τα a b

Πίνακας ANOVA

	SS	df	MS
Model	SSR	df. = k mod	MSR = $\frac{SSR}{df_{mod}}$
Error	SSE	n - (k+1)	MSE = $\frac{SSE}{df_{er}}$
Total	SST	n - 1	

$$SSR \sim \chi^2_k$$

$$SSE \sim \chi^2_{n-k-1}$$

$$SST \sim \chi^2_{n-1}$$

$$F = \frac{MSR}{MSE}$$

$$df_{er} =$$

n - # παραμέτρων
b στο μοντέλο

$$= n - (k+1)$$

$\hat{\sigma}^2 = \text{MSE}$ απερόκλητη εκτιμήτρια του σ^2

$$E(\text{MSE}) = \sigma^2$$

Ελεγχος F (Ftest) για το συνολικό μοντέλο

$$H_0: b_2 = 0$$

$$H_1: b_2 \neq 0$$

X

$$H_0: b_1 = b_2 = \dots = b_k = 0$$

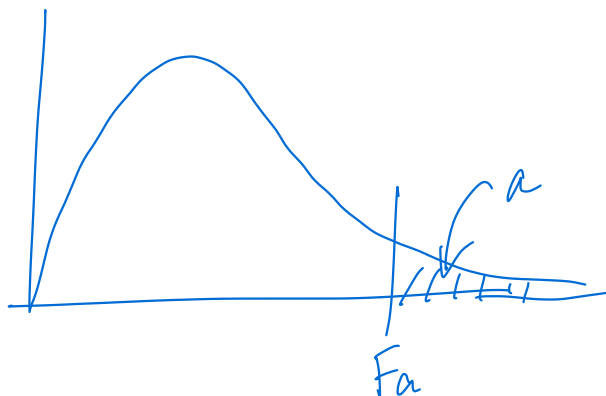
H_1 : τουλάχιστον ένα
από $b_1, \dots, b_k \neq 0$

Στατιστικό ελέγχου

$$F = \frac{\text{MSR}}{\text{MSE}}$$

Accept H_0 if $F \leq F_\alpha$ } κρίση $p \geq \alpha$
Reject H_0 if $F > F_\alpha$ } κρίση $p < \alpha$

$$\text{Av } H_0: F = \frac{\text{MSR}}{\text{MSE}}, \quad F_{df_{\text{mod}}, df_{\text{er}}}$$



Μονοπαράγοντικός Μοντέλο

F-test

$$H_0: b_1 = 0$$

$$H_1: b_1 \neq 0$$

t-test

$$H_0: b_1 = 0$$

$$H_1: b_1 \neq 0$$

} μπορεί
να δώσουν
διαφορετικό
αποτέλεσμα;

ΟΧΙ

Στο μονοπαράγοντικό

ισχύει

$$F = \frac{MSR}{MSE} = t^2$$

για το b_1

$$= \frac{b_1^2}{\frac{2}{S_{b_1}}}$$

$$\left(t = \frac{\hat{b}_1}{S_{b_1}} \right)$$

$$F_{\alpha, 1, n-2} = \left(t_{\alpha/2, n-2} \right)^2$$

$$F > F_{\alpha} \Leftrightarrow t^2 > t_{\alpha/2, n-2}^2 \Leftrightarrow |t| > t_{\alpha/2, n-2}$$

Reject H_0
με το F-test



reject H_0
με το t-test

Παράδειγμα

y, x_1, x_2, x_3

Model 1 : $y = b_0 + b_1 x_1 + \varepsilon$

$H_0: b_1 = 0, H_1: b_1 \neq 0$
 $p \approx 0$

SSR	3142
SSE	5235
SST	8377

$R^2 = 0.3751$

Model 2 : $y = b_0 + b_1 x_3 + \varepsilon$

SSR	2377
SSE	6000
SST	8377

$R^2 \approx 0.28$
 $H_0: b_1 = 0, b_1 \neq 0 (x_3)$
 $p \approx 0$

Model 3 : $y = b_0 + b_1 x_1 + b_3 x_3 + \varepsilon$

xx

"πρῆμερω" $SSR \approx 3142 + 2377$
 $R^2 \approx 0.38 + 0.28 \approx 66\%$

Πραγματοποίηση

SSR	=	3158
SSE	=	5219
		8378

$R^2 = 0,3770$

$H_0: b_1 = 0, b_1 \neq 0, p \approx 0.$
 $\rightarrow H_0: b_3 = 0, b_3 \neq 0, p = 0.588 (x_3)$

$$SSR(x_1, x_3) = 3158$$

$$SSR(x_1) = 3142$$

$$SSR(x_1, x_3) - SSR(x_1) \approx 16$$

$$SSR(x_3|x_1)$$

Επιπλέον μεταβι-
των Y που εξηγεί
από την x_3
όταν προστίθεται
στο x_1

$$SSR(x_3) = 2377$$

$$SSR(x_3|x_1) = 16$$

Η στατιστική σημαντικότητα μιας μεταβλητής
εξαρτάται από το μοντέλο στο οποίο περιλαμβάνεται

Κατά κανόνα αυτές οι διαφορές

οφείλονται σε συσχετίσεις μεταξύ ανεξάρτητων

μεταβλητών (πολυγραμμικότητα
multicollinearity)

Στο μοντέλο $y = b_0 + b_1 x_1 + \dots + b_k x_k + \varepsilon$

t-test $H_0 : b_1 = 0 \quad H_1 : b_1 \neq 0$

Ελέγχει τη στατ. σημαντικότητα του b_1
(συν. της x_1) δεδομένου ότι προσιδεζοται

Σεβασταία στο μοντέλο που ληφίεται
ομάς ως υπόθεση !!!

Οριστικά

π.χ.

Y $x_1 = \text{age}$
 $x_2 = \text{αρζ. νίση (BP)}$

$Y_1 = b_0 + b_1 \text{age} + b_2 \text{BP}$

$H_0 : b_2 = 0 \quad H_1 : b_2 \neq 0.$

Αν η BP είναι στατιστικά σημαντική

Αδρασηία της ηλικίας στο μοντέλο

Σε αυτό το μοντέλο εξετάζουμε τη επίρωση

της BP στο Y

ΕΧΟΝΤΑΣ ΕΓΓΕΓΓΡΑ ως προς
την ηλικία

Συνάρτηση

Y	BP
1	1
⋮	⋮

$$Y = b_0 + b_1 BP$$
$$P \approx 0.001$$

Πως ελέγχουμε αν το effect είναι συν-πραγματικότητα ή όχι;

①

Y	age	BP
⋮	⋮	⋮

$$Y_0 = b_0 + b_1 age + b_2 BP$$

$$H_0: b_2 = 0 \quad H_1: b_2 \neq 0$$

$$P \approx 0.2$$

② $Y_0 = b_0 + b_1 age + b_2 BP$; $H_0: b_2 = 0$, $H_1: b_2 \neq 0$

$$P = 0.005$$

BP είναι στατ. σμφ.
έχοντας ελεγχθεί ως
ΑΡΟ με ηλικία

Μεταβιβάσεις : age, υψος, βάρος : "φυσιολογικές" // μεταβιβάσεις //

Επιλογή στο μονοδιάστατο παράδειγμα

$$Y = b_0 + b_1 X_1 \quad : \quad R^2 = 0.3751 \quad P_{X_1} \approx 0$$

$$Y = b_0 + b_2 X_2 \quad : \quad R^2 = 0.2918 \quad P_{X_2} \approx 0$$

$$Y = b_0 + b_1 X_1 + b_2 X_2 \quad : \quad R^2 = 0.7255 \quad P_{X_1} \approx 0$$

$$P_{X_2} \approx 0$$

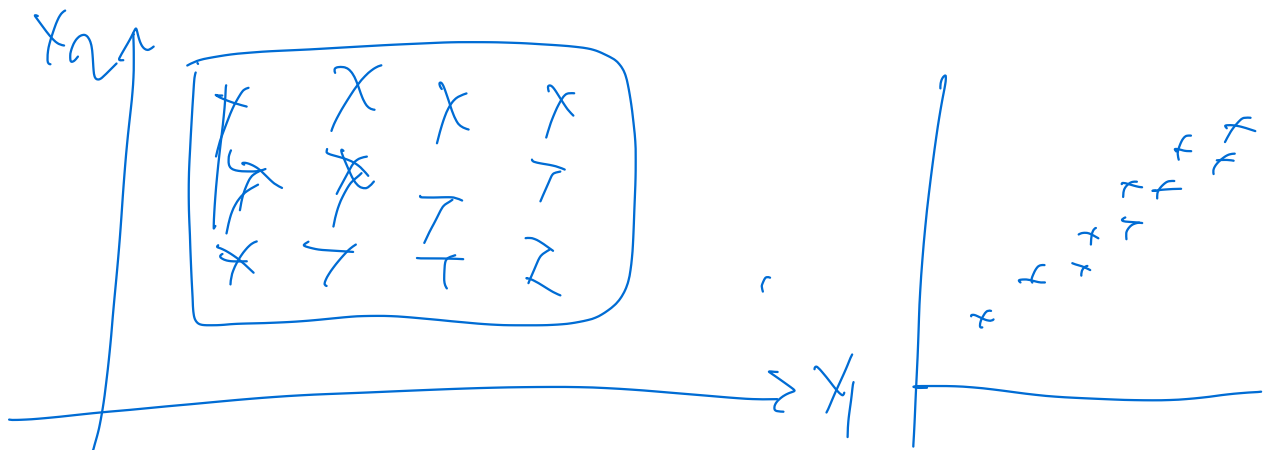
$$\text{SSR}(X_2 | X_1) = \text{SSR}(X_1, X_2) - \text{SSR}(X_1)$$

$$= 6077 - 3142 \approx 2900$$

$$\text{SSR}(X_2) = 2444$$

$$\rho(X_1, X_2) \approx 0$$

~~X_1, X_2
ορθογώνιες~~



$$Y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3$$

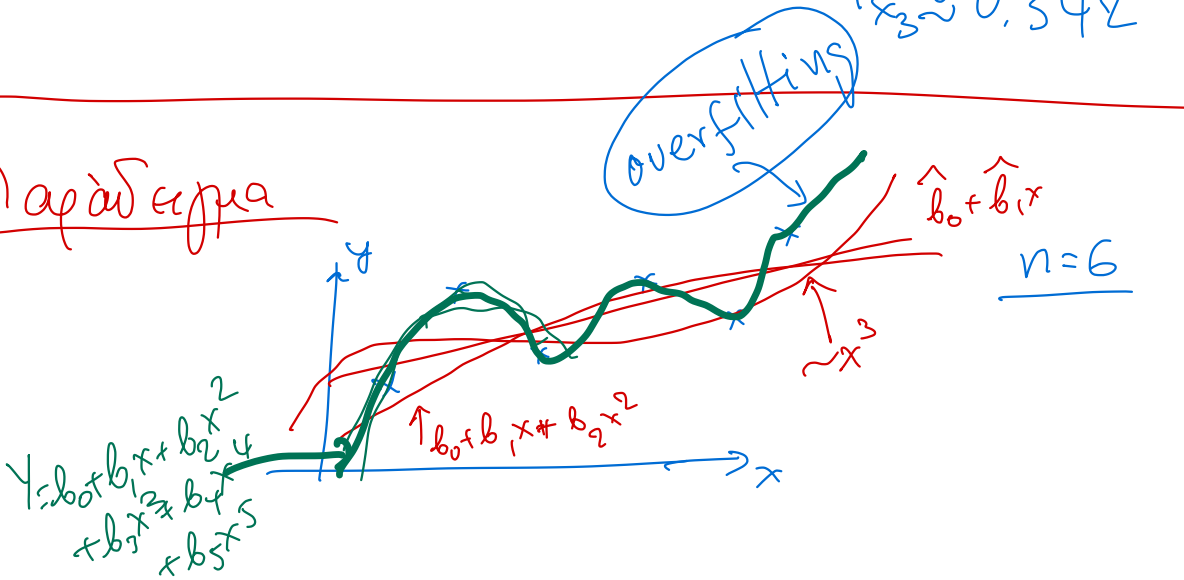
$$\underline{R^2 = 0.7280}$$

$$P_{x_1} \approx 0$$

$$P_{x_2} \approx 0$$

$$P_{x_3} \approx 0.342$$

Παράδειγμα



$$\textcircled{1} Y = b_0 + b_1 x + \epsilon$$

$$\textcircled{2} Y = b_0 + b_1 x + b_2 x^2 + \epsilon$$

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

$$+ b_4 x^4$$

$$Y = b_0 + b_1 x + b_2 x^2 + \dots + b_5 x^5$$

$$\textcircled{R^2 = 1}$$

$$SSE = 0$$

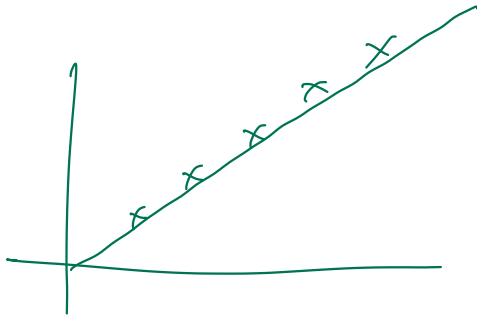
$$SSR = SST$$

$$P_{b_1} \approx 1$$

$$P_{b_2}$$

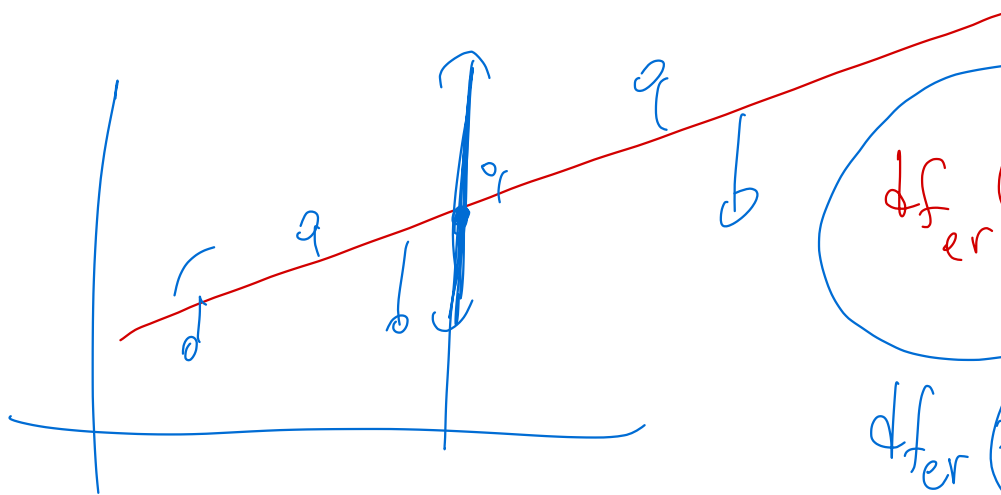
$$df_{er} = n - \#b = 6 - 6 = 0$$

$$\hat{\sigma}^2 = MSE = \frac{SSE}{df_{er}} = \frac{0}{0}$$



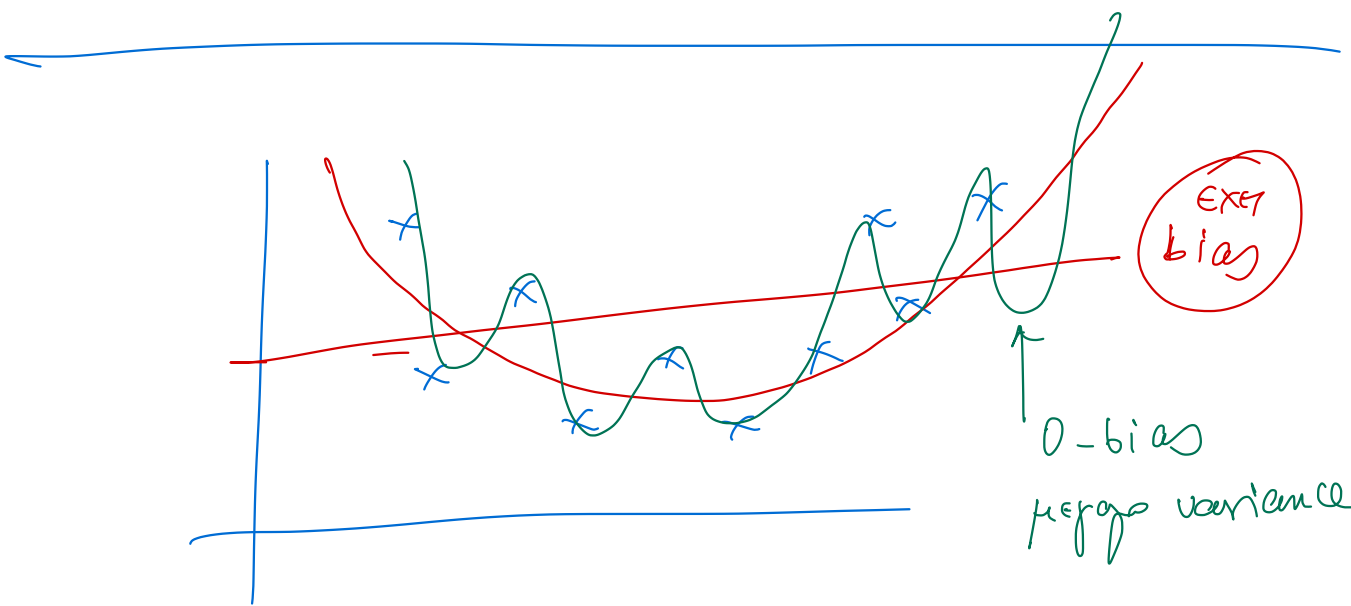
$$SSE = 0$$

$$df_{er} = n - 1$$



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

df_{er} (peraga)



EXTRA bias

0-bias

peraga variance