

14-1-2025

## Μοντέλο Πολλαπλής Παμδρόμησης

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

$\left. \begin{matrix} X_1 \\ \vdots \\ X_k \end{matrix} \right\}$  ανεξάρτητες μεταβλητές / παράγοντες / predictors

$Y$  : εξαρτημένη μεταβλητή / response

$\varepsilon \sim N(0, \sigma^2)$  ανεξάρτητες μεταξύ παρατηρήσεων

Δεδομένα

|          | $X_1$    | $X_2$    | $\dots$ | $X_k$    | $Y$   |
|----------|----------|----------|---------|----------|-------|
| 1        | $x_{11}$ | $x_{12}$ | $\dots$ | $x_{1k}$ | $y_1$ |
| 2        | $x_{21}$ | $x_{22}$ |         | $x_{2k}$ | $y_2$ |
| $\vdots$ |          |          |         |          |       |
| $j$      | $x_{j1}$ | $x_{j2}$ |         | $x_{jk}$ | $y_j$ |
| $\vdots$ |          |          |         |          |       |
| $n$      | $x_{n1}$ |          |         | $x_{nk}$ | $y_n$ |

$$\hat{y}_j = b_0 + b_1 x_{j1} + b_2 x_{j2} + \dots + b_k x_{jk}$$

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

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

$X_1, \dots, X_k = \left\{ \begin{array}{l} \text{διαφορετικά φυσικά μεγέθη} \\ \text{ανάπτυξης} \end{array} \right.$  από  $X_j$

πχ.  $X_2 = X_1^2$

①  $Y = b_0 + b_1 X_1 + b_2 X_1^2$

②  $Y = b_0 + b_1 X_1 + b_2 X_2 + b_3 \underbrace{X_1 X_2}_{X_3}$

③  $Y = b_0 + b_1 X_1 + b_2 \underbrace{\ln(X_1)}_{X_2}$

# Μέθοδος Ελαχ. Τετραγώνων $\#b = k+1$

$$\min_{b_0, b_1, \dots, b_k} SSE \rightarrow \underbrace{\hat{b}_0, \hat{b}_1, \dots, \hat{b}_k}_{\text{Εκτιμήσεις Ελαχ. Τετραγώνων}}$$

$$E(\hat{b}_0) = b_0 \quad \hat{b}_0 \sim N(b_0, \sigma_{\hat{b}_0}^2)$$

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

$$E(\hat{b}_k) = b_k \quad \hat{b}_k \sim N(b_k, \sigma_{\hat{b}_k}^2)$$

## Πίνακας ANOVA

|       | SS  | df                          | MS                                  |
|-------|-----|-----------------------------|-------------------------------------|
| Model | SSR | $\#b - 1 = df_{\text{mod}}$ | $MSR = \frac{SSR}{df_{\text{mod}}}$ |
| Error | SSE | $n - \#b = df_{\text{er}}$  | $MSE = \frac{SSE}{df_{\text{er}}}$  |
|       | SST | $n - 1$                     |                                     |

$$SSR = \sum_{j=1}^n (\hat{y}_j - \bar{y})^2 \quad : \text{μεταβ. που εξηγείται από το μοντέλο}$$

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

$$SST = SSR + SSE$$

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

$$E(MSE) = \sigma^2$$

# Παρατηρήσεις

①  $Y = b_0 + b_1 X + \varepsilon$  ,  $E(Y|X=x) = b_0 + b_1 x$ .

$b_1$ : κλίση (ρυθμός μεταβολής)

$b_0 = E(Y|X=0)$  [όταν ο αριθμητής στις τιμές του  $X$  στο δείγμα]

②  $E(Y|X_1=x_1, \dots, X_k=x_k) = b_0 + b_1 x_1 + \dots + b_k x_k$

$b_0 = E(Y|X_1=X_2=\dots=X_k=0)$

[όταν στο δείγμα υπάρχουν παρατηρήσεις με όλα τα  $X_j = 0$ ]

③  $b_1 =$  ρυθμός μεταβολής της  $EY$  ως προς  $X_1$  όταν τα  $X_2, \dots, X_k$  μένουν σταθερά

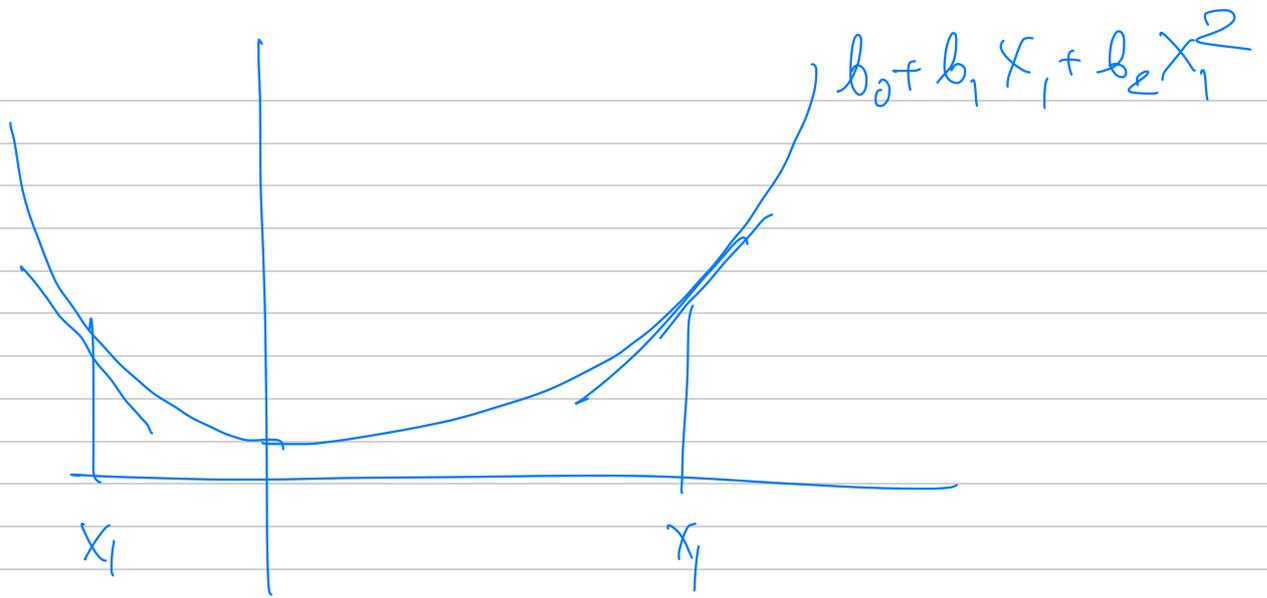
$b_1 = \frac{\partial EY}{\partial X_1}$  ,  $b_2 = \frac{\partial EY}{\partial X_2}$  ...

(?)

④ Παράδ:  $EY = b_0 + b_1 X_1 + b_2 X_1^2$

?  $b_1 = \frac{\partial EY}{\partial X_1}$  ~~X~~

$\frac{\partial EY}{\partial X_1} = b_1 + 2b_2 X_1$  !!



5)  $EY = b_0 + b_1 \cdot X_1 + b_2 X_2$

$X_1 = \text{age}$

$X_2 = \text{height (cm)}$

$X_1 \in [0, 10]$

$X_2 \in [30, 120]$

$b_1 = \frac{\partial EY}{\partial X_1}$  ( ? )

Множество  $\text{Cor}(X_1, X_2) > 0$



$$\text{Aw } \text{Cor}(X_1, X_2) \approx 1$$

$$\Rightarrow X_2 = \gamma_0 + \gamma_1 X_1 \Rightarrow X_1 = \frac{X_2 - \gamma_0}{\gamma_1}$$

Τότε όμως :  $Y = b_0 + b_1 X_1 + b_2 X_2 =$

$$= b_0 + b_1 X_1 + b_2 (\gamma_0 + \gamma_1 X_1)$$

$$= b_0 + b_2 \gamma_0 + (b_1 + b_2 \gamma_1) X_1$$

$$= \underline{\delta_0 + \delta_1 X_1}$$

$\dots$

$$= \theta_0 + \theta_1 X_2$$

Πολυσημασια (Multicollinearity)

Συσχεση μεταξύ ανεξαρτητων μεταβλητων.

# Παράδειγμα

# Exercise. Rdata.

$$Y = \text{score.}$$

$$X_1 = \text{age}$$

$$X_2 = \text{extime.}$$

①

$$Y = b_0 + b_1 \cdot \text{age}$$

$$\text{SSR} = 10\ 255\ 910$$

$$\text{SSE} = 1\ 485\ 592$$

$$\hline 11\ 741\ 503$$

$$R^2 = 0.8735$$

$$Y = b_0 + b_1 \cdot \text{extime}$$

$$\text{SSR} = 10\ 701\ 573$$

$$\text{SSE} = 1\ 039\ 929$$

$$\hline \text{SST} = 11\ 741\ 503$$

$$R^2 = 0.9114$$

$$Y = b_0 + b_1 \cdot \text{age} + b_2 \cdot \text{extime}$$

$$\text{SSR} = 10\ 808\ 248$$

$$\text{SSE} = 933\ 254$$

$$\hline \text{SST} = 11\ 741\ 503$$

$$R^2 = 0.9205$$

$$SSR(\text{age}) = 10255910$$

$$SSR(\text{age}, \text{extime}) = 10808248$$

$$SSR(\text{age}, \text{extime}) - SSR(\text{age}) = \underline{552338}$$
$$= SSR(\text{extime} | \text{age})$$

Γενικά  $SSR(\text{extime} | \text{age}) \neq SSR(\text{extime})$

γιατί;

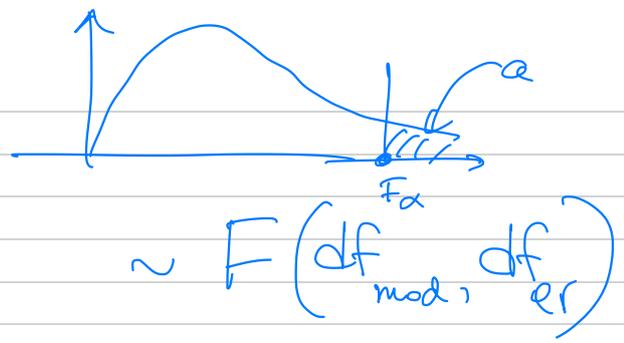
παρουσιάζουμε κολυζα !!

---

Η μεταβλητότητα του  $Y$  που εξηγείται από μια ανεξ. μεταβλητή  $X$ , ή η στατιστική σημαντικότητα της  $X$  ως predictor εξαρτάται από τις άλλες μεταβλητές που ήδη υπάρχουν στο μοντέλο παλινδρόμησης

1

$$F = \frac{MSR}{MSE} = \frac{\frac{SSR}{df_{mod}}}{\frac{SSE}{df_{er}}}$$



F statistic συγκριών ποσοτήτων

Ελεγχος  
F για  
συγκριτικό  
ποσοτό

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

$$H_1: \text{τουλάχιστον ένα } \neq 0$$

$$H_0: EY = b_0$$

(όχι το ποσοτό με σταθερούς ανεξάρτητες)

$$F > F_\alpha \Leftrightarrow \text{reject } H_0 \quad p < \alpha$$

$$F \leq F_\alpha \Leftrightarrow \text{accept } H_0 \quad p \geq \alpha$$

Είδος Απεικρίσεων (k=1)

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

$$F = \frac{MSR}{MSE} = \frac{\frac{SSR}{1}}{\frac{SSE}{df_{er}}} \sim F(1, df_{er})$$

$$H_0: b_1 = 0 \quad H_1: b_1 \neq 0$$

Οπως έχουμε δει ου

$$\left\{ \begin{array}{l} H_0: b_1 = 0 \quad H_1: b_1 \neq 0 \\ \text{Ολο μονοπαραγοντικό} \end{array} \right\} \text{t-test 1 για } b_1$$

$$t = \frac{\hat{b}_1}{\text{SE}_{\hat{b}_1}}$$

Στο μοντέλο  $Y = b_0 + b_1 X$

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

Από ν.θ.  $F(1, k) \stackrel{d}{=} t_k^2$

$$F_\alpha = t_{\alpha/2}^2$$

$$F \leq F_\alpha \Leftrightarrow t^2 \leq t_{\alpha/2}^2 \Leftrightarrow |t| \leq t_{\alpha/2}$$

accept  $H_0$

accept  $H_0$

2

## Εστειχος F για μέρος του μοντέλου

Model 1 (full model)

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_p X_p$$

$b_{p+1} X_{p+1} + \dots + b_k X_k$   
(p < k)

Model 2 (partial model)

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_p X_p$$

Model 2 nested στο Model 1 (ζωγράφισμα)

Ξέρουμε  $SSR(X_1, X_2, \dots, X_p, \dots, X_k) \geq SSR(X_1, \dots, X_p)$

$$\frac{SSE(X_1, X_2, \dots, X_p, \dots, X_k)}{SST} \leq \frac{SSE(X_1, \dots, X_p)}{SST}$$

Επομένως  $SSR(X_1, \dots, X_k) - SSR(X_1, \dots, X_p) =$

$$= SSR(X_{p+1}, \dots, X_k | X_1, \dots, X_p) \geq 0$$

$$\Rightarrow R_{full}^2 \geq R_{partial}^2$$

k-p ανεξαρτησίες

$$H_0: \beta_{p+1} = \beta_{p+2} = \dots = \beta_k = 0, \quad H_1: \text{τουλάχιστον ένα } \neq 0$$



εστειχος για μέρος του model 1

$$SSR_{full} + SSE_{full} = SSR_{par} + SSE_{par} \Rightarrow$$

$$\Rightarrow SSR_{full} - SSR_{par} = SSE_{par} - SSE_{full}$$

$$df_{er, full} = n - (k+1) = n - k - 1$$

$$df_{er, par} = n - (p+1) = n - p - 1 > n - k - 1$$

$$df_{er, par} - df_{er, full} = k - p$$

$$H_0 : b_{p+1} = b_{p+2} = \dots = b_k = 0, \quad H_1 : \text{zωταξισωv} \\ \text{vna} \neq 0$$

$$F = \frac{\frac{SSE_{par} - SSE_{full}}{df_{er, par} - df_{er, full} = k - p}}{MSE_{full}} \sim F(k - p, n - k - 1)$$

$$\text{accept } H_0 \text{ αν } F \leq F_{\alpha}(k - p, n - k - 1)$$

$$\text{reject } H_0 \text{ αν } F > F_{\alpha}(k - p, n - k - 1)$$

~~εισ. ανειζωμ~~  
1

ειλεχος F για ογο το μονελο ειδη  
ανειζωμ

$$\text{full } Y = b_0 + b_1 X_1 + \dots + b_k X_k \quad \text{full}$$

$$\text{part } Y = b_0$$

$$H_0 : b_1 = b_2 = \dots = b_k = 0. \quad H_1 : \text{zωta} \text{ vna} \neq 0$$

$$\text{To partial model eδw } Y = b_0 \Rightarrow \hat{b}_0 = \bar{Y}$$

$$SSE_{\text{partial}} = \sum (Y_i - \bar{Y})^2 = SST$$

$$SSE_{\text{par}} - SSE_{\text{full}} = SST - SSE_{\text{full}} = SSR_{\text{full}}$$

$$df_{er, par} - df_{er, full} = k.$$

$$F = \frac{\frac{SSR_{full}}{df_{mod, full}}}{MSE_{full}} = \frac{MSR_{full}}{MSE_{full}} = F$$

↓  
για το  
test των  
των μετεξων

## Είδει απιντων

full  $Y = b_0 + b_1 X_1 + \dots + b_p X_p + b_{p+1} X_{p+1}$

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

$$H_0 : b_{p+1} = 0$$

$$H_1 : b_{p+1} \neq 0$$



t-test  $b_{p+1}$

Ο ελεγχος t για  $b_j = 0$  η  $b_j \neq 0$

ελέγξει τα συντελεστωτα των  $X_j$

Προϋποθέτων όλων των υπόλοιπων μεταβλητών

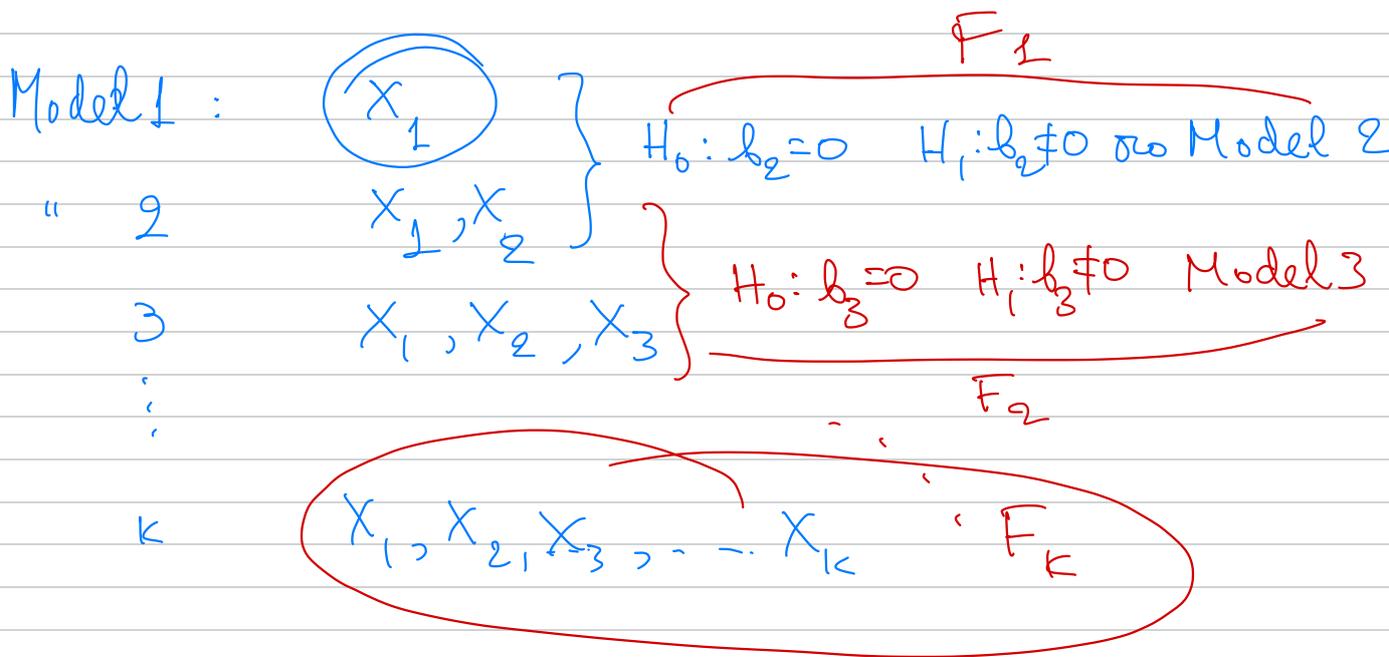
(δηλ. όλων η  $X_j$  περαιτε ζεφεωζατα  
παρονηα όλων των αλληων).

3 Έστω ότι στο μοντέλο

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

α μεταβλητές μεινώνων με τα όρια

$$X_1, X_2, \dots, X_k$$



$F_1, F_2, \dots, F_k$  : Ftests type I  
 (sequential Ftests)  
εξαρτώνται από τα όρια  
εισαγωγής

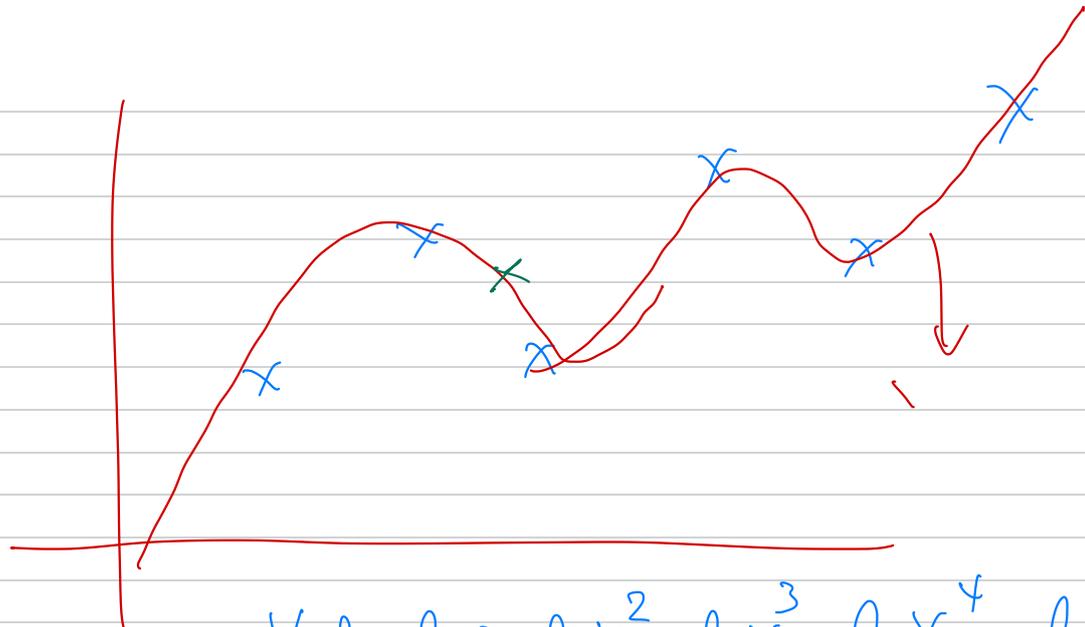
4 Ftests type III (partial)

$$H_0: b_j = 0 \quad H_1: b_j \neq 0 \text{ στο Model } k$$

(H  $X_j$  αξιοποιείται)  $\Leftrightarrow$  t-tests

δεν εξαρτώνται από τα όρια εισαγωγής

n=6.



$$Y = b_0 + b_1 X + b_2 X^2 + b_3 X^3 + b_4 X^4 + b_5 X^5$$

$$SSE = 0 \Rightarrow R^2 = 1$$

$$df_{er} = 0$$

$$MSE = \frac{0}{0} \quad ?$$

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

