

30-3-2023

# Διαγνωστικοί Έλεγχοι - Έλεγχοι Κατανομών

$$Y = b_0 + b_1 X_1 + \varepsilon$$

$$Y = b_0 + b_1 X_1 + \dots + b_k X_k \quad \left| \begin{array}{l} \text{αυ } n-1 \\ \text{dfer} = n-k-1 \end{array} \right. \quad \left| \begin{array}{l} \text{df} = n-1-k-1 \end{array} \right.$$

Δείγμα  $\begin{pmatrix} x_1, y_1 \\ \vdots \\ x_n, y_n \end{pmatrix}$

$$Y_j = b_0 + b_1 X_j + \varepsilon_j$$

$$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n \text{ iid } \sim \mathcal{N}(0, \sigma^2)$$

τυχαίες  
ανεξαρτητές

1)  $\varepsilon_j \sim \mathcal{N}$  κανονικά

2)  $\text{Var}(\varepsilon_j) = \sigma^2$  ομοσκεδαστικά

3) Ανεξαρτητά

πως ελέγχω  
αν ισχύουν οι  
υποθέσεις;

Αν έχω ένα δείγμα  $(z_1, z_2, \dots, z_n)$  από

μια κατανομή  $\Rightarrow \exists$  έλεγχοι  $\left\{ \begin{array}{l} \text{κανονικότητα} \\ \text{ανεξαρτησία} \end{array} \right.$

Όμως

εδώ

$$\varepsilon_j = y_j - b_0 - b_1 x_j$$

$\uparrow$   $\uparrow$   $\uparrow$   
αγνωστα!!  $\uparrow$   $\uparrow$   $\uparrow$   
αγνωστα!!

Αν αντί των αγνωστων  $b_0, b_1$  πάρουμε  $\hat{b}_0, \hat{b}_1$

$$\hat{\varepsilon}_j = e_j = y_j - (\hat{b}_0 + \hat{b}_1 x_j), \quad j=1, \dots, n$$

$$= y_j - \hat{y}_j = \underline{\text{κατάλοιπα}}$$

$$\sum (y_i - \hat{y}_i)^2 = SSE = \sum e_j^2$$

① Έστω Δυσκολίες  $s_j^2 = \text{Var}(e_j)$   $s_j^2$  δεν είναι ίσα

② Πάντα  $\sum e_j = 0 \iff$  δεν είναι ανεξάρτητα

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

?  $\exists$  μετασχηματισμός των  $e_j$   
που καθιστούν ως ιδέμετες

### Μεθόδοι Καθαρισμού

1) Τυποποιημένα καθαρίσματα (standardized)

$$\tilde{e}_j = \frac{e_j}{s_j} \Rightarrow \text{Var}(\tilde{e}_j) = 1$$

συνεπώς να υπάρχει πρόβλημα  
με ανεξαρτησία.

2) Studentized residuals

$$r_j = \frac{\tilde{e}_j}{\sqrt{1-h_j}} = \frac{e_j}{s_j \sqrt{1-h_j}}$$

$h_j = \text{leverage}$   
(μάρτυρας)  
ως παραμετρική

(H: hat matrix  $H = X^T (X^T X)^{-1} X$ )

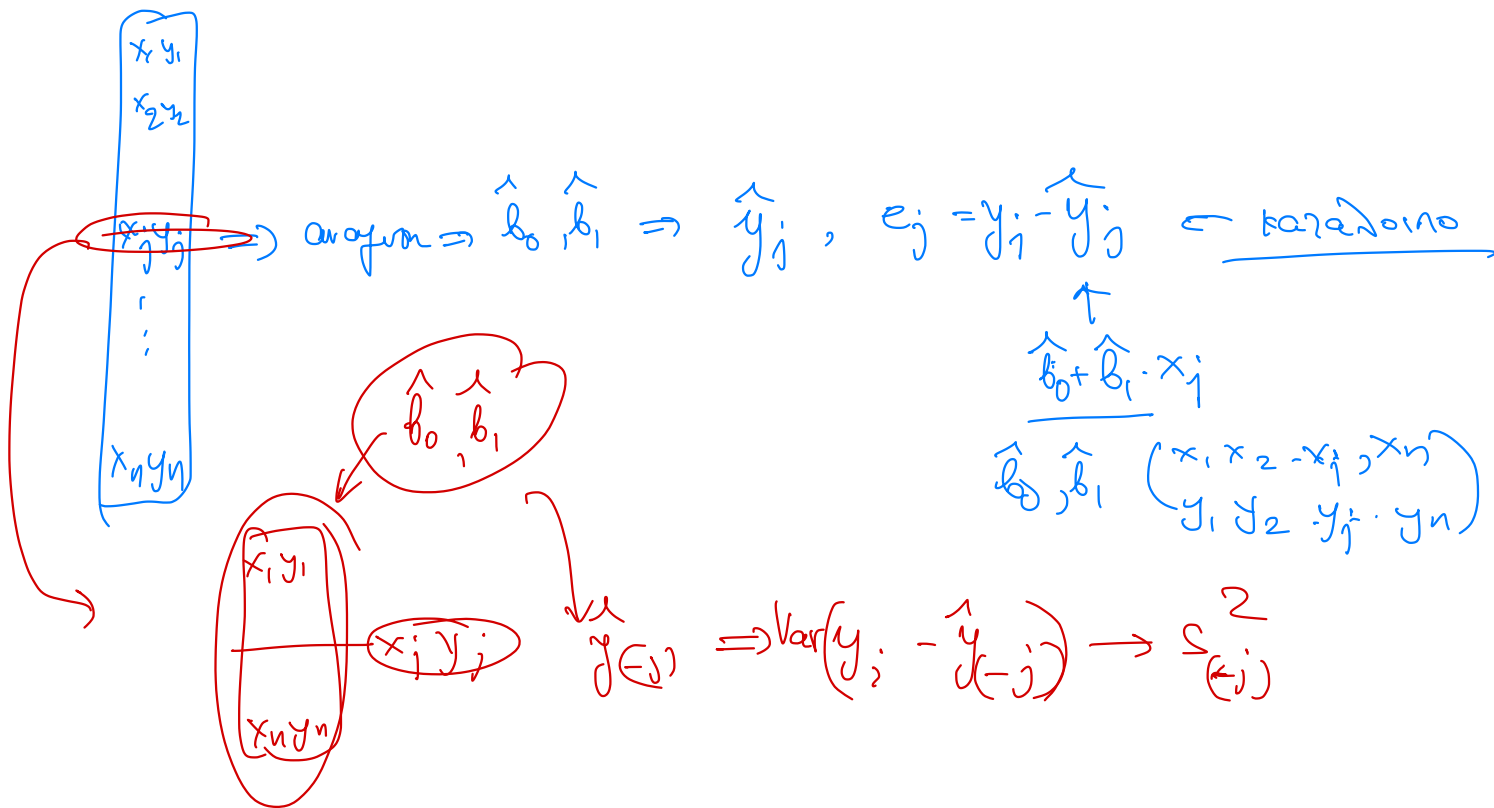
$h_j$  διαγώνια στοιχεία του H.

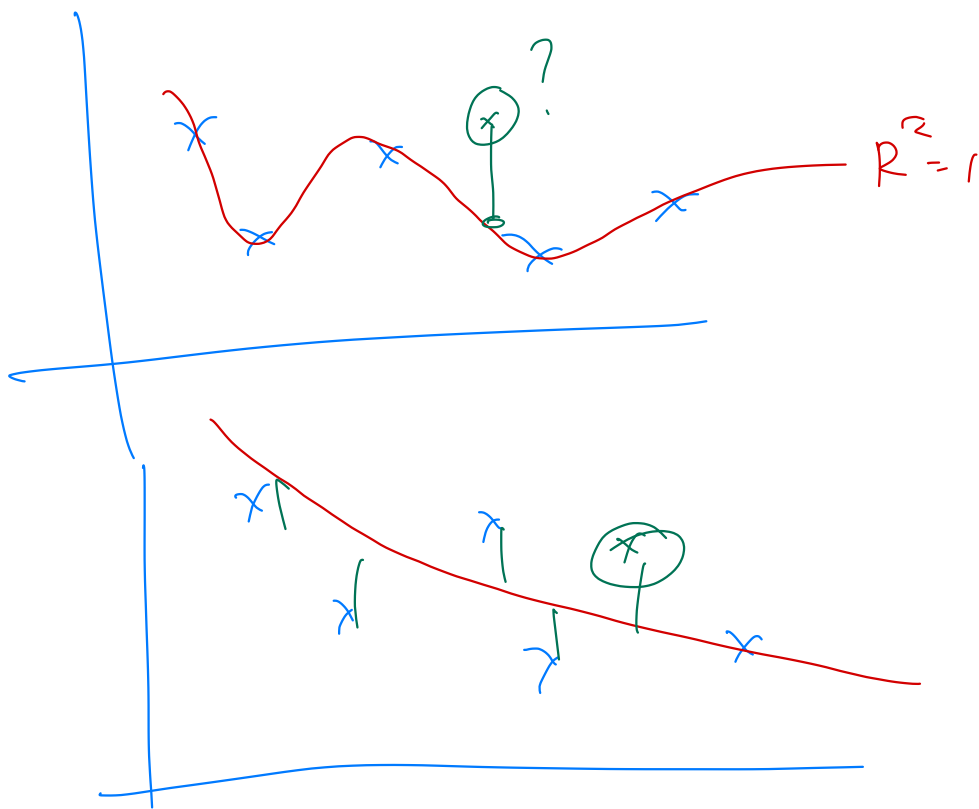
Αποδεικνύεται  $t_j, j=1, \dots, n$  ποσότητες ανεξάρτητες z-τε.  
 $\sim t$

### 3) Jackknife residuals

$$\sqrt{2} t_j = \frac{e_j}{S_{(-j)} \sqrt{1-n_j}} \left\{ \begin{array}{l} \text{αυσχέραια} \\ \sim t_{n-k-2} \left( \begin{array}{l} \text{όταν } n-k-2 \gg \\ t_{n-k-2} \sim N \end{array} \right) \end{array} \right.$$

$S_{(-j)}$  = ζη. απόκτη του  $\hat{y}_j$  παρακείμενου  $j = y_j - \hat{y}_{(-j)}$   
 όπως όπως σου αναφέρει ότι έχει  
 ομοιομορφία με παρακείμενο  $j$ .





Τρεις διαφορετικοί έλεγχοι χρησιμοποιούνται  
 είτε ως studentized  
 " " jackknife residuals

Στο Stata } → studentized → "standard"  
 κ' όσο R } jackknife → "student"

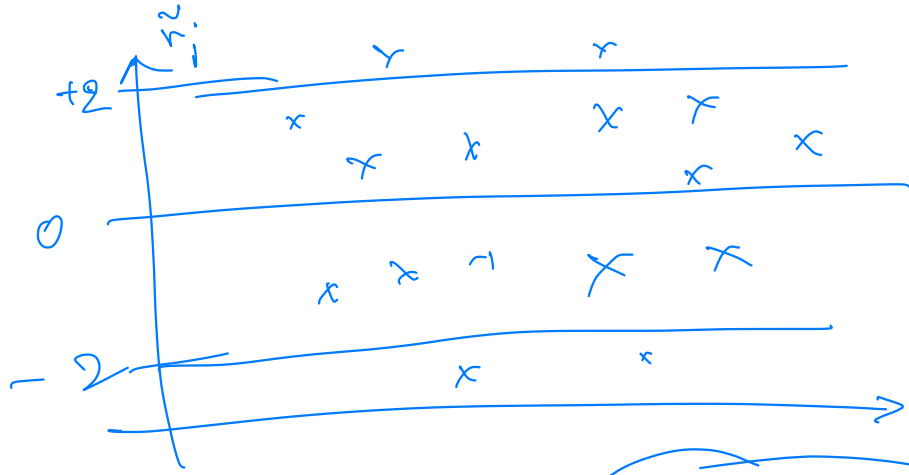
## Διαγνωτικοί Έλεγχοι

Έλεγχοι των υποθέσεων των γραμμικών  
 μοντέλων με βάση τα κατάλοιπα (student  
 " jackknife)

- 1) Τραβήματα
- 2) Έλεγχοι υποθέσεων.

Συνήθως το κατατάξω σε 2 ομάδες

κατάλογος  
προς  
πρόβλ. υπεS



θεωρία

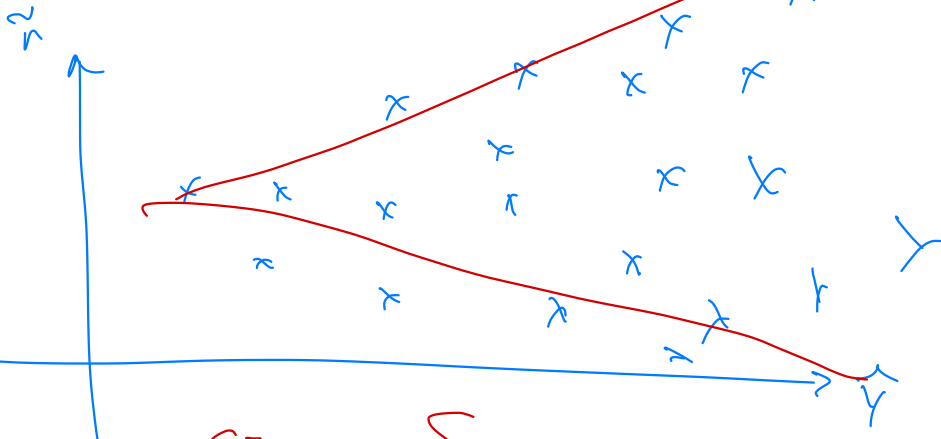
$Cov(e_j, \hat{y}_j) = 0$

Παράδειγμα υποθέσεων

1

επιτυχία  
εξεροσκέδαση

επιτυχία  
Bartlett



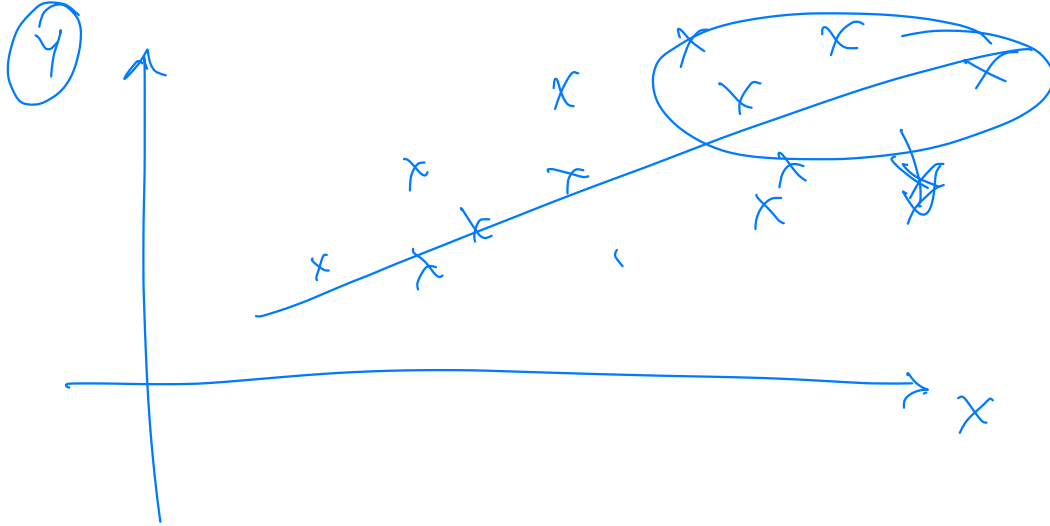
χοάνη

εξεροσκέδαση

Διορθωτικά μέτρα

μετασχηματισμοί ως  $Y$ .

π.χ.  $\tilde{Y} = \sqrt{Y}$  μετριάσει ως προς ως  $Y$



Προσοχή

$Y = b_0 + b_1 X + \varepsilon$ :  $b_1$ : μεταβολή της  $EY$  όταν  $X$  αυξηθεί κατά 1

αν  $\sqrt{Y} = \gamma_0 + \gamma_1 X + \varepsilon$

$\tilde{Y} = \gamma_0 + \gamma_1 X + \varepsilon$

$\gamma_1$ : μεταβολή της  $E(\tilde{Y}) = E(\sqrt{Y})$

$E(\sqrt{Y}) \neq \sqrt{EY}$

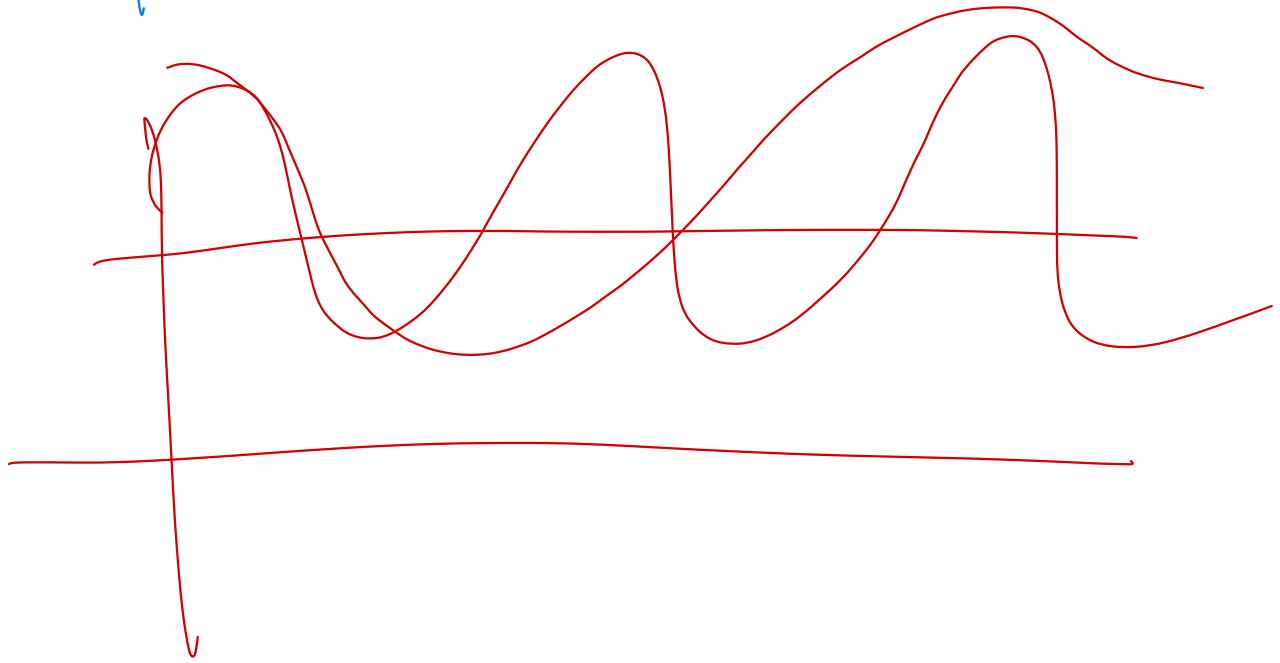
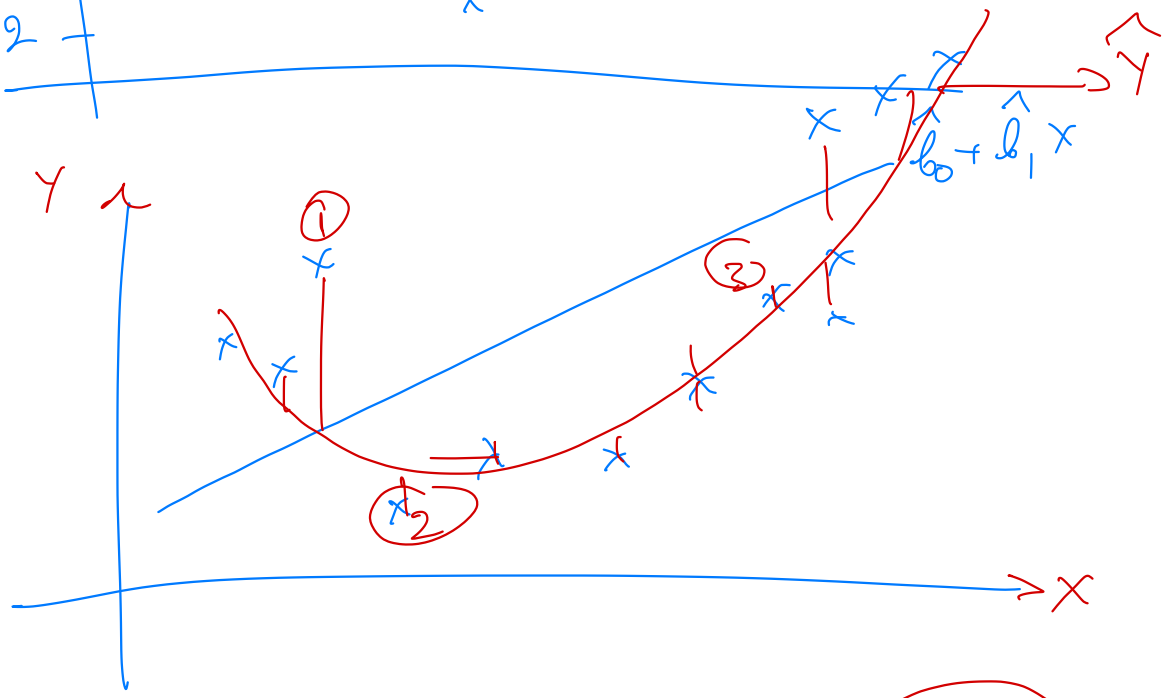
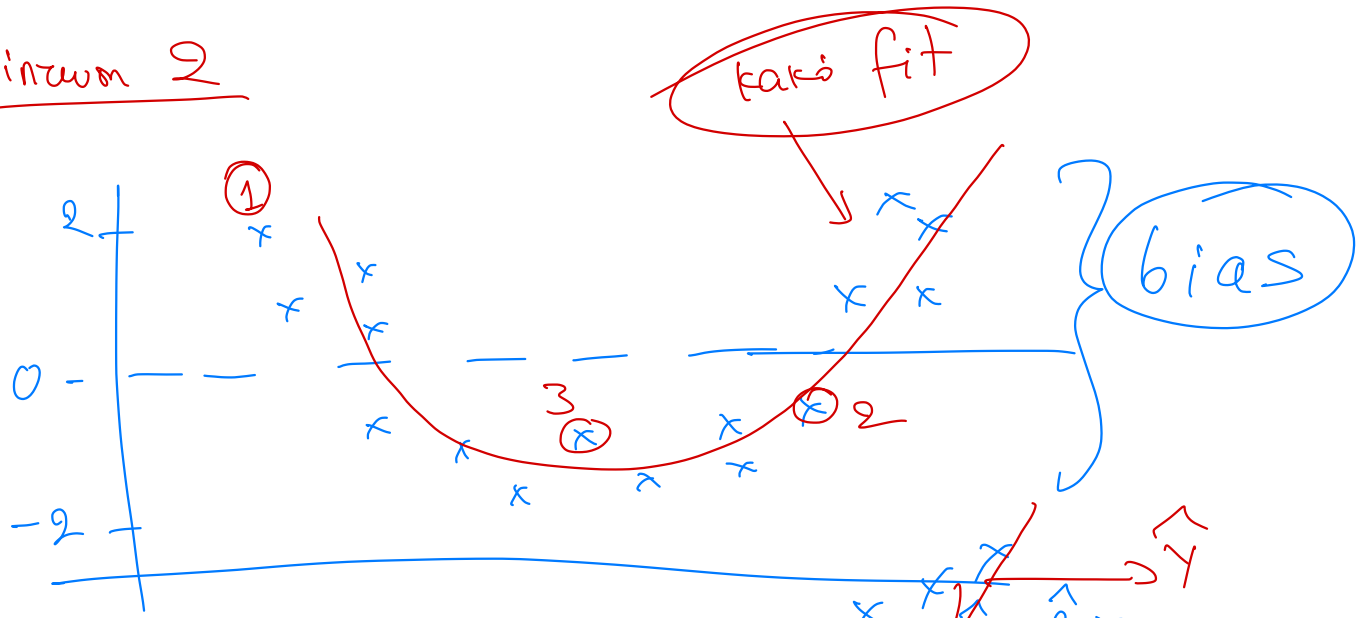
Οι κλίσεις δεν έχουν σημασία ως μεταβολής της  $E(Y)$

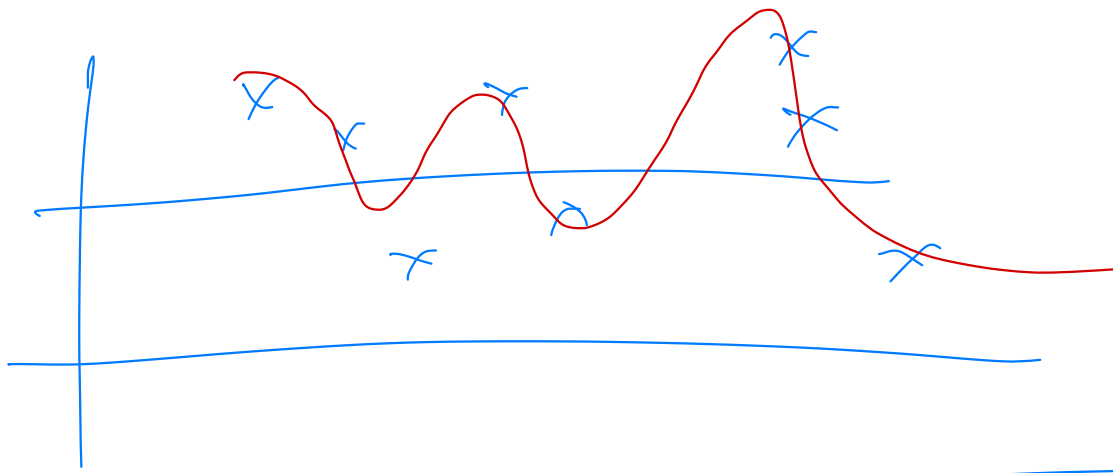
Άλλοι μετασφ.  $\sqrt[3]{Y}$ ,  $\log(Y)$

Τι γίνεται αν στο δείγμα υπάρχουν  $Y < 0$ ?

$\tilde{Y}_j = Y_j - \min(Y) > 0 \leftarrow$  μετασφ.

# Пример 2





### ③ Κανονικότητα

plot γραφικά

έλεγχοι { KS

Shapiro Wilk.

### ④ Αυτοσυσχέτιση Καρακρίνου.

$\tilde{r}_1, \tilde{r}_2, \tilde{r}_3, \dots, \tilde{r}_n \leftarrow$  ανεξάρτητα

Η σειρά ως παρατήρηση μπορεί να παιξε  
ρόλο όταν το δείγμα προέρχεται από  
χρονοσειρά !!

π.χ. χρονοσειρά θερμοκρασίας  $x_j$  }  $j=1, 2, \dots$   
 + σχετ. όφους  $y_j$  } αλλά Σάββατο

Χρειαζόμαστε μεθόδους χρονοσειρών



# Ελεγχος αυτοσυσχέτισης

## Durbin-Watson test

$d$  : durbin-watson statistic

$$d \in [0, 4]$$

$d \approx 2$  =  $\cancel{A}$  αυτοσυσχέτιση

$d \approx 0$  αρνητική αυτοσυσχ

$d \approx 4$  θετική "

Αποδοτική  
Ελεγχος

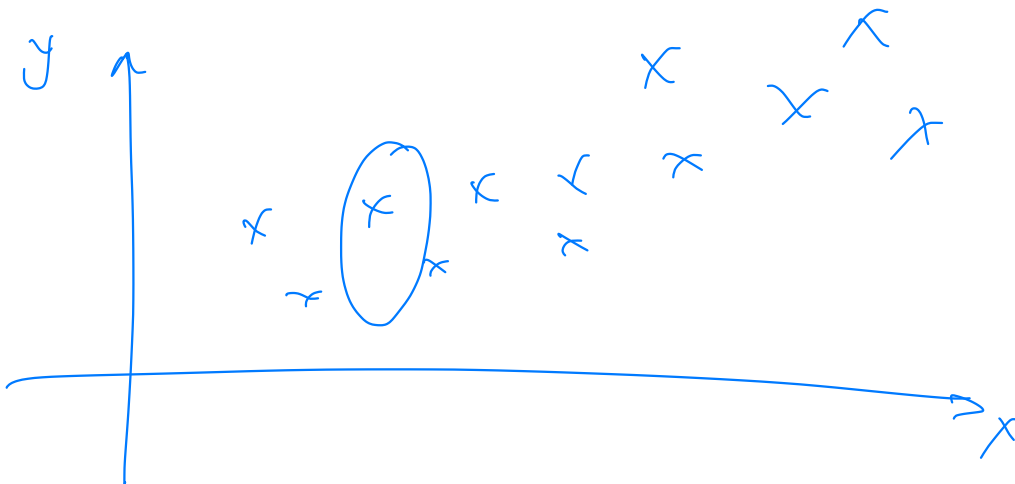
p-value  
προσέγγιση

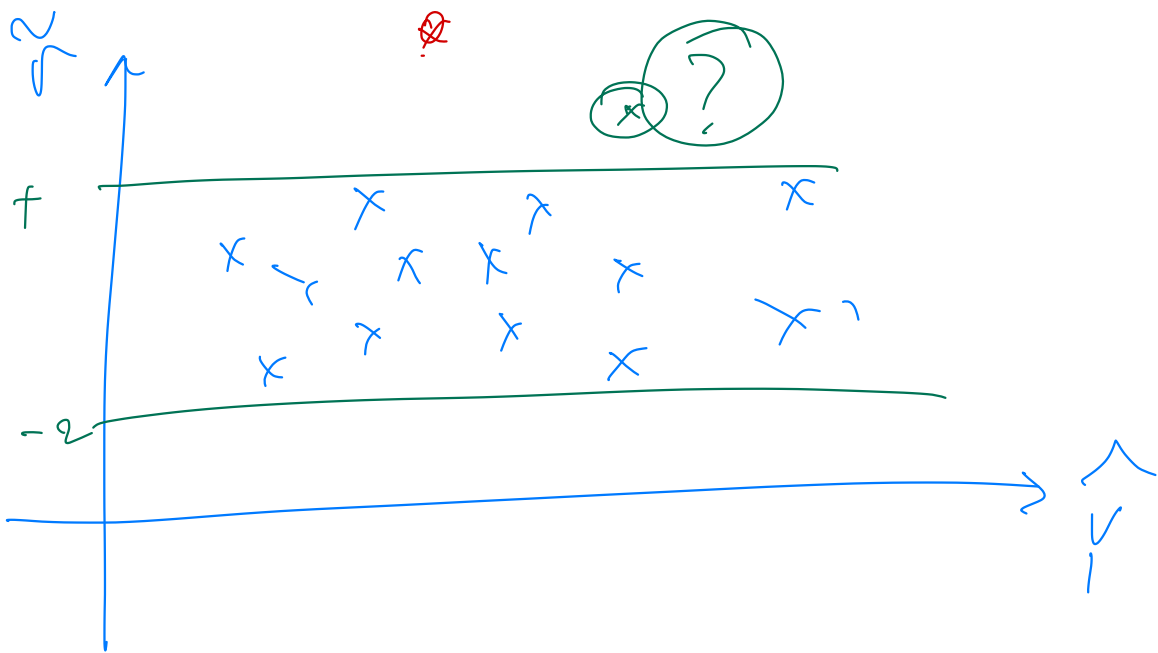
## Εκτροπές - Ενδραστηρικές Παρατηρήσεις

1

(Outliers - Influential Observation)

outlier  $\bar{x} \leftarrow 2\sigma$





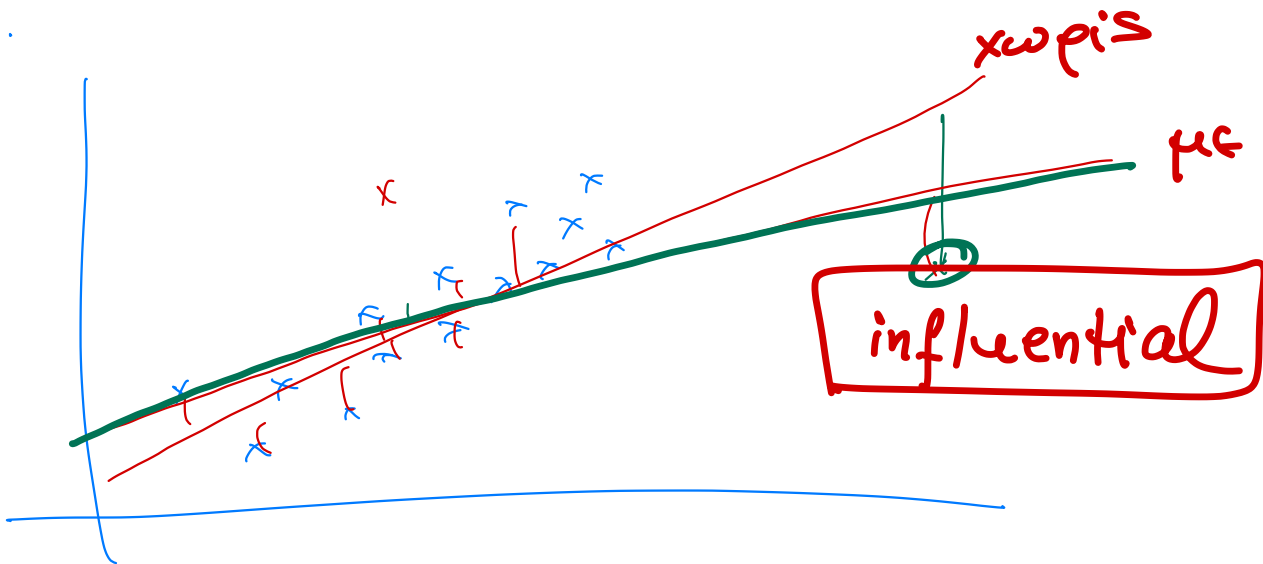
① λαιος καταραφι

② " μέρησων

✓  
??

Στατιστικοί έλεγχοι για outliers \*

②



Outlier

$\tilde{r}_j$   $\begin{matrix} \text{ναός} & \text{μεγάλη} & \text{ή} \\ \text{ναός} & \text{μικρή} & \text{απει} \end{matrix}$  κρίσιμα

Επιπτώσεις  $\left\{ \begin{array}{l} \text{leverage } h_j \quad \underline{0 \leq h_j \leq 1} \\ \text{Cook's distance } d_j \end{array} \right.$

Πίνακες κρίσεων απειών για  $\tilde{r}_j, h, d.$

Πολλαπλοί έλεγχοι

$P(\text{ουστό}) = 1 - \alpha$

$P(\text{η έχω ουστό}) \approx (1 - \alpha)^n$   
 $\approx 1 - n\alpha$   
 $= 1 - b$

$\left[ \begin{array}{l} \alpha \approx 0 \\ (1 - \alpha)^n \approx 1 - n\alpha \end{array} \right]$

$b = n\alpha \Rightarrow \alpha = \frac{b}{n}$

↑ εν. σημ.

προσέγγιση  
διόρθωση  
Bonferroni

Πίνακς κρίσεων απειών  $\tilde{r}$  έχω  
 $h$  διόρθω  
 $d$  Bonferroni

# Επιλογή Μοντέλου

## Βήματα Επιλογής Μοντέλου

- ① Μέγιστο Μοντέλο (full model)
- ② κριτήρια σύγκρισης μοντέλων
- ③ Στατιστικές επιλογής μεταβλητών
- ④ Ανάλυση Τεχνικού Μοντέλου
- ⑤ Ερμηνείες - Προβλέψεις

## ① Μέγιστο Μοντέλο (full model)

Περιέχει όλες τις υποτιμώμενες μεταβλητές  
ε' όσον τους επιμαέον όρους

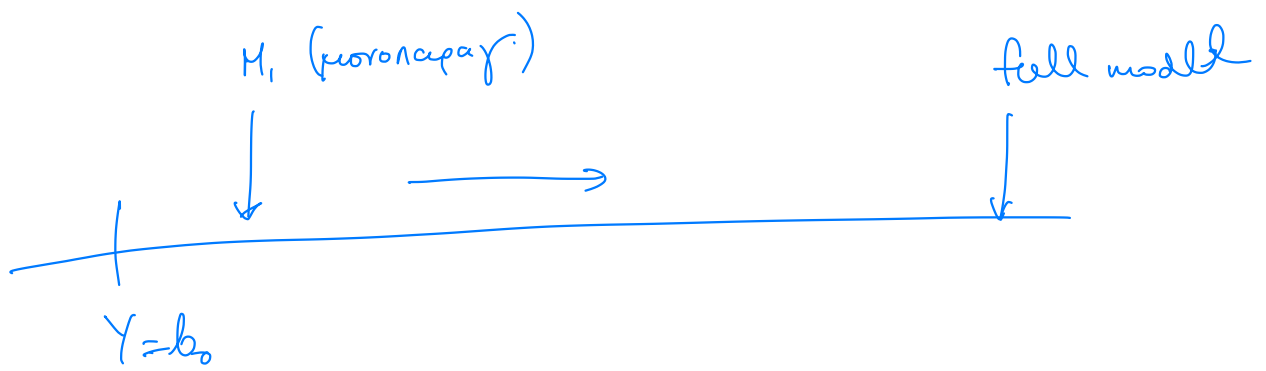
(π.χ. μεγαλύτερη τάξη, αλληλεπίδραση κτλ)

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

df<sub>er</sub> : "μεγάλο"

$$df_{er} \geq 20 \text{ ή } 30$$

Πιο ασφαλές κανόνας  $n \geq 5k$  ή  $n \geq 10k$ .



## 2) Κριτήρια Σύγκρισης Μοντέλων

### 1) Για nested models

$$\begin{array}{ll}
 \text{(full)} & Y = b_0 + b_1 X_1 + \dots + b_k X_k \quad \text{(full)} \\
 \text{(part)} & Y = b_0 + b_1 X_1 + \dots + b_p X_p \quad \text{(Υποσύνολο)}
 \end{array}$$

← nested στο full

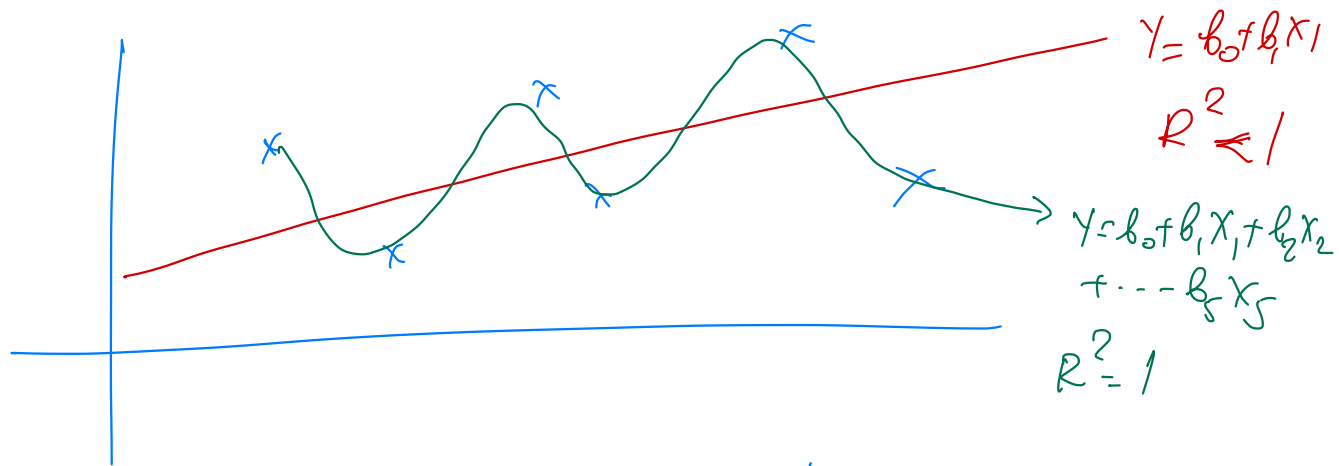
### a) Partial F-test

full  $H_0: b_{p+1} = b_{p+2} = \dots = b_k = 0$   $H_1: \text{zωσταx evx} \neq 0$

### b) Adjusted $R^2$

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

Για nested models  $R^2_{\text{full}} \geq R^2_{\text{part}}$



$$\text{adj-}R^2 = 1 - (1 - R^2) \cdot \frac{n-1}{\underbrace{n-p-1}_{df_{er}}} \quad p = \text{ap. μεταβλητών.}$$

Όταν  $n-1 \approx df_{er} \quad (p \ll n) \quad \frac{n-1}{n-p-1} \approx 1$

$$\text{adj } R^2 \approx 1 - (1 - R^2) = R^2$$

Όταν  $p \rightarrow n$   $\frac{n-1}{n-p-1}$  μεγάλο  $\Rightarrow \underline{\text{adj } R^2 < R^2}$

Πιο αυστηρά κριτήριο για το % μεταβλητών που εξηγεί, λαμβάνοντας υπόψη κι το μέγεθος των ποσών

## Ⓡ Στατιστικά $C_p$ -Mallows

Full model ( $k$ )

Partial " ( $p$ )

$$C_p = \frac{\text{SSE}(p)}{\text{MSE}(k)} - [n - 2(p+1)]$$

$$\text{Ou } \text{Moyéto}(p) \approx \text{Moyéto}(k)$$

$$\text{MSE}(p) \approx \text{MSE}(k)$$

$$\text{Tóze } \frac{\text{SSE}(p)}{\underbrace{n-(p+1)}_{df_{er, (p)}}} \approx \text{MSE}(k) \Rightarrow$$

$$\frac{\text{SSE}(p)}{\text{MSE}(k)} \approx n-(p+1) \Rightarrow$$

$$\Rightarrow C_p \approx n-(p+1) - [n-2(p+1)] = p+1$$

$$\text{Ozan } \text{Model}(p) \approx \text{Model}(k) \Rightarrow \text{C}_p \approx \underline{p+1}$$

Ou Moyéto(p) xepótepo ano Moyéto(k)

$$\Rightarrow C_p > p+1$$

② Κριτήρια Πληροφόρησης (oxi αναπαίρω για nested models)

AIC = Akaike Information Criterion

BIC = Bayesian " "

$$\text{AIC} = -2 \log(L) + 2k$$

L : πιθανοφάνεια δείγματος  $\hat{\beta}$  από τις εκτιμήσεις των  $\hat{\beta}$ .

k : μέγεθος μοτέλου

μικρότερη τιμή AIC  $\Rightarrow$  καλύτερο μοντέλο

## Στατιστική Επιλογή Μεταβλητών

Έστω ένα full model με  $k$  ανεξ. μεταβ.

Πόσα δυνατά υπομόντέλα μπορεί να ορίσω;  $2^k$  !!

$$4 \left\{ \begin{array}{l} Y = b_0 \\ Y = b_0 + b_1 X_1 \\ Y = b_0 + b_1 X_2 \\ Y = b_0 + b_1 X_1 + b_2 X_2 \end{array} \right.$$

## Μεθόδους stepwise regression

① Forward  $(X_1, \dots, X_k)$  υποψήφιες μεταβλητές

	R	State
ⓐ) Όλα τα παρατηρήσιμα	AIC	p-value
$Y = b_0 + b_1 X_1$	AIC	p-value
$Y = b_0 + b_1 X_2$	⋮	⋮
$Y = b_0 + b_1 X_k$	AIC	

Επιλέγουμε από με το min p εστω  $X_j$

$$b_0 + b_1 X_j \left\{ \begin{array}{l} b_0 + b_1 X_j + b_2 X_1 \\ b_0 + b_1 X_j + b_2 X_2 \\ \vdots \\ b_0 + b_1 X_j + b_2 X_k \end{array} \right. \min(p)$$

↑  
εκτός της  $X_j$   
p-value της νέας στο μοντέλο

Αν  $\min(p) > p_{\text{enter}} \Rightarrow$  σταματάμε



Για να γίνει μια νέα μεταβλητή πρέπει να έχει το μικρότερο  $p$  αντίοφες που είναι ερωτά  $t'$  το  $P < P_{enter}$ .

## 2) Backward

Ξεκινάμε με full model

$$Y = b_0 + b_1 X_1 + \dots + b_j X_j + \dots + b_p X_p$$

$\uparrow$              $\uparrow$              $\uparrow$              $\rightarrow \max(P)$   
 $P$              $P$              $P$

$P_{max} > P_{remove}$

$X_j$  φεύγει

↓  
 ∴  
 συνέχεια μέχρι για  $P < P_{remove}$

## 3) Stepwise

Κάθε φορά που προστίθεται μια μεταβλητή με  $P < P_{enter}$

εξετάζονται όφες οι μεταβλητές του νέου μοντέλου  $t'$  αφαιρείται αυτή

με  $\max p$  αν αυτό είναι  $P > P_{remove}$ .

Προσέγγιση: Αν  $P_{\text{enter}} = 0,07$

$P_{\text{remove}} = 0,05$

μπορεί να βρούμε σε κάποιο βήμα μια

$X_i$  έτσι που αν προσέχτι έχει  $p=0,06$

$P_{\text{enter}} < P_{\text{remove}}$ .