

11-10-2021

# Seminar R / Rstudio

R Markdown  $\rightarrow$  πακέτο (knitr) του R.

Τεμάρι: Εφαρμογές (σε Εργασίες)

Κάτω από τον υπόθετον

$$Y|X=x = x^T b + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$\varepsilon_1, \dots, \varepsilon_N$  : ανεξάρτητα

$$\hat{b} = G y$$

$$E \hat{b} = b, \quad \text{Var}(\hat{b}) = (X^T X)^{-1} \sigma^2$$

$$y \sim \mathcal{N}(x^T b, \sigma^2 I_N)$$

$$\Rightarrow \hat{b} \sim \mathcal{N}(b, (X^T X)^{-1} \sigma^2)$$

$\rightarrow$  Διαστ. Γενικ. b  
Εγκυροί t-test

$$(N-p-1) \hat{\sigma}^2 \sim \sigma^2 \chi^2_{N-p-1}$$

$$\hat{b}, \hat{\sigma}^2 \text{ ανεξάρτητα}$$

Θεώρημα Gauss-Markov:  $\hat{b}$  : ελάχιστη διασπορά  
ανάμεσα στις γραμμικές-αφροδύναμες εκτιμήσεις  
του b.

# Επιλογή Μεταβλητών / Μέθοδοι Στατιστικής Μεταβλητών

Αν  $p$  'μεγάλο' σε σχέση με το  $N$

$$df_{\text{er}} = \underbrace{N - p - 1}_{\text{5p} \approx 10p} \quad \text{Εμπειρικός}$$

$$\hat{\sigma}^2 = \frac{\sum (y_j - \hat{\mu})^2}{N - 1}$$

$$Y_j \sim N(\mu, \sigma^2)$$

---

αντί του άγνωστου  $\sum_{j=1}^N x_j^T b$

$$\hat{\sigma}^2 = \frac{\sum_{j=1}^N (y_j - \sum_{j=1}^p x_j^T \hat{b})^2}{\underbrace{N - (p+1)}_{\rightarrow df_{\text{error}}}} \quad b \in \mathbb{R}^{p+1}$$

---

Αν έχουμε λογικές ανεξ. μεταβλητές.

ελέγχος υποθέσεων :  $H_0: b_j = 0, \quad H_1: b_j \neq 0$

Το αποτέλεσμα του ελέγχου εξαρτάται από όλο το μοντέλο κ' όχι μόνο από τη μεταβλητή  $X_j$

λόγω λογιστικής συσχρηματικότητας

συσχρηματικότητα ανάμεσα στις ανεξάρτητες μεταβλητές.

# Επιλογή Μεταβλητών

$X_1, \dots, X_p$  ανεξάρτητες μεταβλητές

$2^p$  δυνατά μοντέλα

⊗

$Y = b_0$

$Y = b_0 + b_1 X_1$

$Y = b_0 + b_1 X_2$

$\vdots$

$Y = b_0 + b_1 X_p$

$Y = b_0 + b_1 X_1 + b_2 X_2$

$Y = b_0 + b_1 X_1 + b_2 X_3$

$\vdots$

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

Δάμπερς

## Best subset selection Methods

Συγκρίνουμε ομά τα δυνατά μοντέλα  
ε' επιλέγουμε το καλύτερο με βάση  
κάποιο κριτήριο.

## Προσεγγιστικές Μέθοδοι Επιλογής Μεταβλητών

### Stepwise Regression

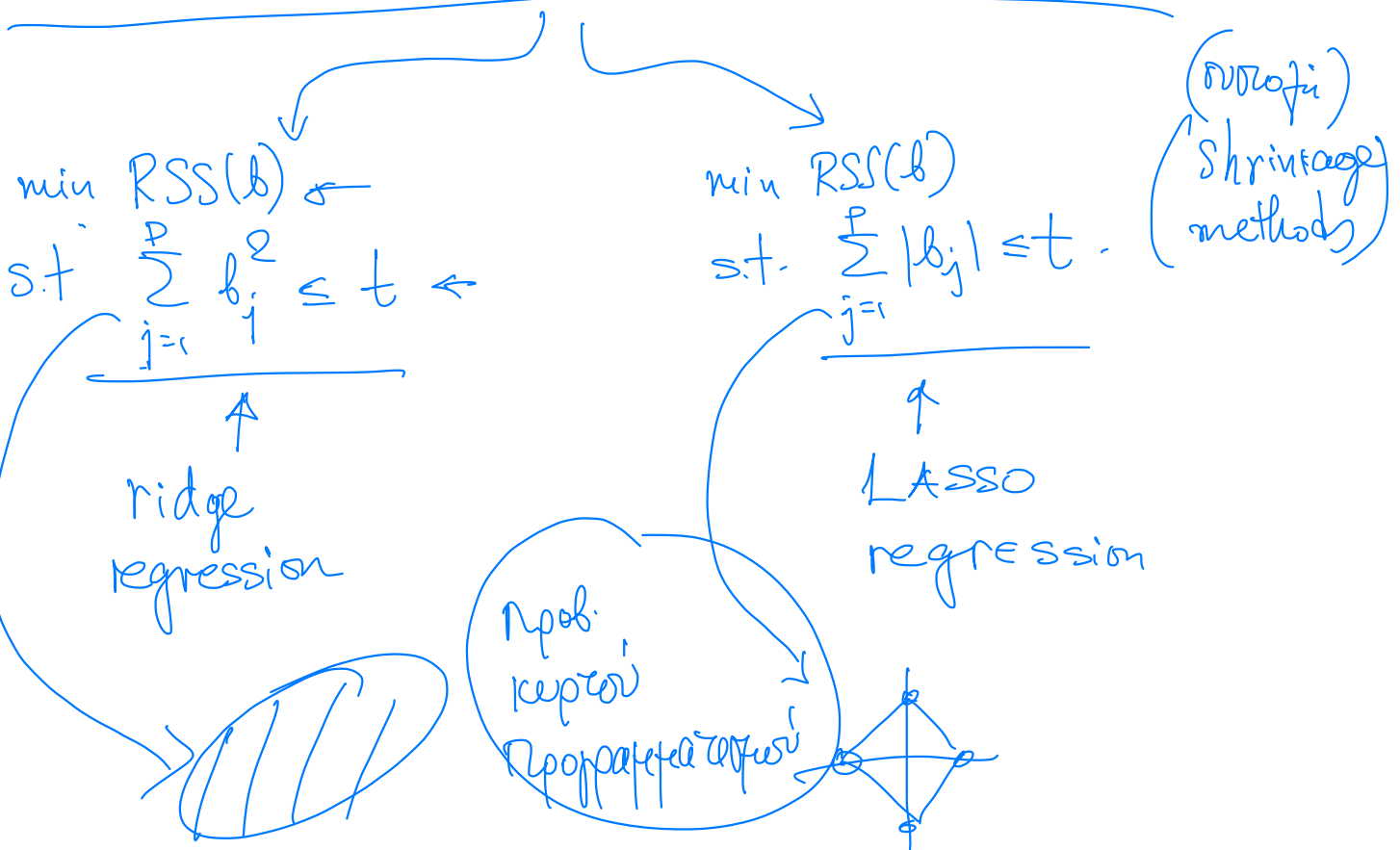
↙  
forward

Προσίδεται μια  
μεταβλητή σε κάθε  
βήμα

# Αλμ Πρoσέγγιση

$$\left\{ \begin{array}{l} \min \text{RSS}(b) \\ \text{s.t. } \underline{\text{αρ. περ. στο ποσείο}} \leq \underline{k} \end{array} \right\}$$

$$\begin{array}{l} \min \text{RSS}(b) \\ \sum_{j=1}^p 1(b_j \neq 0) \leq k. \end{array}$$



Ridge

