

Γραμμικά Μοντέλα  
Το Γραμμικό Μοντέλο με Κανονικά Κατανεμημένα  
Σφάλματα

Διδάσκουσα: Λουκία Μελιγκοτσίδου  
Εθνικό και Καποδιστριακό Πανεπιστήμιο Αθηνών  
Τμήμα Μαθηματικών

April 2, 2020

## Εκτίμηση με τη Μέθοδο Μέγιστης Πιθανοφάνειας

Το απλό γραμμικό μοντέλο, κάτω από την υπόθεση της κανονικότητας για τους τυχαίους όρους, γράφεται ως

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n.$$

**Παρατήρηση:** Η υπόθεση των ασυσχέτιστων σφαλμάτων,  $Cov(\varepsilon_i, \varepsilon_j) = 0, i \neq j$ , κάτω από την υπόθεση της κανονικότητας αντιστοιχεί σε υπόθεση ανεξαρτησίας (ανεξάρτητα σφάλματα).

Επομένως, για τις απαντητικές μεταβλητές μεταβλητές έχουμε υποθέσει  $Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$ ,  $i = 1, \dots, n$ . Η υπόθεση συγκεκριμένης κατανομής για της μεταβλητές ενδιαφέροντες επιτρέπει την εκτίμηση των παραμέτρων του μοντέλου με τη μέθοδο μέγιστης πιθανοφάνειας.

Συνάρτηση πιθανοφάνειας (*likelihood function*):

$$\begin{aligned} L(\beta_0, \beta_1, \sigma^2) &= \prod_1^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(Y_i - \beta_0 - \beta_1 X_i)^2\right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2}(Y_i - \beta_0 - \beta_1 X_i)^2\right\} \end{aligned}$$

$$\begin{aligned} \text{log-likelihood: } \log L &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum (Y_i - \beta_0 - \beta_1 X_i)^2 \\ \left. \begin{aligned} \frac{\partial \log L}{\partial \beta_0} &= \frac{1}{\sigma^2} \sum (Y_i - \beta_0 - \beta_1 X_i) = 0 \\ \frac{\partial \log L}{\partial \beta_1} &= \frac{1}{\sigma^2} \sum X_i(Y_i - \beta_0 - \beta_1 X_i) = 0 \end{aligned} \right\} \text{κανονικές εξισώσεις} \end{aligned}$$

$\Rightarrow \hat{\beta}_0, \hat{\beta}_1$  ίδιες με εκτιμήτριες ελαχίστων τετραγώνων.

$$\begin{aligned} \frac{\partial \log L}{\partial \sigma^2} &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (Y_i - \beta_0 - \beta_1 X_i)^2 = 0 \\ \Rightarrow \hat{\sigma}^2 &= \frac{1}{n} \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2 \end{aligned}$$

**Προσοχή:** Η εκτιμήτρια μέγιστης πιθανοφάνειας της διασποράς δεν είναι αμερόληπτη.

### Παρατηρήσεις

Οι εκτιμήτριες  $\hat{\beta}_0$  και  $\hat{\beta}_1$  ως εκτιμήτριες ελαχίστων τετραγώνων (EET)

1) είναι αμερόληπτες (Unbiased)

2) έχουν ελάχιστη διασπορά μεταξύ των α.ε. που είναι γραμμικοί συνδυασμοί των  $Y_i$  (Best Linear Unbiased Estimators - BLUE).

Επίσης ως εκτιμήτριες μέγιστης πιθανοφάνειας (ΕΜΠ) είναι

(1) συνεπείς

(2) επαρκείς

(3) αμερόληπτες εκτιμήτριες ελάχιστης διασποράς (έχουν ελάχιστη διασπορά μεταξύ όλων των α.ε.)

### Κατανομές των Εκτιμητών των Συντελεστών της Παλινδρόμησης

#### Κατανομή της $\hat{\beta}_1$

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, i = 1, \dots, \varepsilon_i \sim N(0, \sigma^2), \text{ανεξάρτητα} \\ \Rightarrow Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2).$$

$$\text{Είναι } \hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \sum k_i Y_i.$$

Άρα η  $\hat{\beta}_1$  ακολουθεί την κανονική κατανομή ως γραμμικός συνδυασμός κανονικών τυχαίων μεταβλητών.

'Εχουμε δείξει ότι

$$E(\hat{\beta}_1) = \beta_1$$

$$\sigma^2(\hat{\beta}_1) = V(\sum k_i Y_i) = \sum k_i^2 V(Y_i) = \sigma^2 \sum k_i^2 = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$'Αρα η  $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2})$$$

**Παρατήρηση:** Η διασπορά  $\sigma^2$  είναι άγνωστη. Μπορεί όμως να εκτιμηθεί από την α.ε. της

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2.$$

Οπότε η εκτιμούμενη διασπορά της  $\hat{\beta}_1$  είναι

$$S^2(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}.$$

Κατανομή του  $\hat{\beta}_0$

Είναι  $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$ , όπου τα  $Y_i$  είναι κανονικές τ.μ. και το  $\hat{\beta}_1$  επίσης κανονική τ.μ. Άρα το  $\hat{\beta}_0$  ακολουθεί την κανονική κατανομή ως γραμμικός συνδυασμός κανονικών τ.μ.

Εχουμε δείξει ότι

$$\begin{aligned} E(\hat{\beta}_0) &= \beta_0 \\ \sigma^2(\hat{\beta}_0) &= V(\bar{Y} - \hat{\beta}_1 \bar{X}) = V(\bar{Y}) + V(\hat{\beta}_1 \bar{X}) - 2Cov(\bar{Y}, \hat{\beta}_1 \bar{X}) \\ &= \frac{\sigma^2}{n} + \bar{X}^2 V(\hat{\beta}_1) - 2\bar{X} \underbrace{Cov(\bar{Y}, \hat{\beta}_1)}_0 = \frac{\sigma^2}{n} + \bar{X}^2 \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \sigma^2 \left[ \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]. \end{aligned}$$

Εκτιμούμενη διασπορά:  $S^2(\hat{\beta}_0) = \hat{\sigma} \left[ \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]$ .

**Πρόταση.** Ισχύει

$$Cov(\hat{\beta}_1, \bar{Y}) = 0$$

**Απόδειξη.**

$$\begin{aligned} \text{Ιδιότητες συνδιακύμανσης: } Cov(X+Y, Z) &= Cov(X, Z) + Cov(Y, Z) \\ Cov(X+Y, Z+W) &= Cov(X, Z) + Cov(X, W) + Cov(Y, Z) + Cov(Y, W) \end{aligned}$$

$$\begin{aligned}
Cov(\hat{\beta}_1, \bar{Y}) &= Cov\left(\sum_{i=1}^n k_i Y_i, \frac{\sum_{j=1}^n Y_j}{n}\right) = \sum_{i=1}^n \sum_{j=1}^n Cov(k_i Y_i, \frac{Y_j}{n}) = \sum_{i=1}^n \sum_{j=1}^n \frac{k_i}{n} Cov(Y_i, Y_j) \\
&= \sum_{i=1}^n \frac{k_i}{n} Cov(Y_i, Y_i) = \frac{1}{n} \sum_{i=1}^n k_i V(Y_i) = \frac{\sigma^2}{n} \underbrace{\sum_{i=1}^n k_i}_{0} = 0
\end{aligned}$$

αφού  $Cov(Y_i, Y_j) = 0, i \neq j$ .

**Συμπερασματολογία για τους Συντελεστές του Γραμμικού Μοντέλου**

**Πρόταση.** Η τυχαία μεταβλητή  $\frac{\hat{\beta}_1 - \beta_1}{s(\hat{\beta}_1)} \sim t_{(n-2)}$

**Απόδειξη.**

Ισχύει: Αν  $Z \sim N(0, 1)$ ,  $U \sim X_{(r)}^2$  και  $Z, U$  ανεξάρτητες, τότε

$$T = \frac{Z}{\sqrt{\frac{U}{r}}} \sim t_{(r)}$$

Έχουμε  $\hat{\beta}_1 \sim N(\beta_1, \sigma^2(\hat{\beta}_1)) \Rightarrow \frac{\hat{\beta}_1 - \beta_1}{\sigma(\hat{\beta}_1)} \sim N(0, 1)$ .

Επίσης ισχύει ότι ( $\theta$ α αποδειχθεί στη συνέχεια):

$$\frac{\sum \hat{\varepsilon}_i^2}{\sigma^2} = \frac{\sum (Y_i - \hat{Y}_i)^2}{\sigma^2} = \frac{\sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2}{\sigma^2} \sim X^2(n-2).$$

$$\text{Και } \varepsilon\pi\varepsilon\delta\hat{\sigma}^2 = \frac{1}{n-2} \sum \hat{\varepsilon}_i^2, \text{ έχουμε } \frac{\hat{\sigma}^2(n-2)}{\sigma^2} \sim X^2(n-2).$$

Τώρα

$$\left. \begin{array}{l} S^2(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ \sigma^2(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \end{array} \right\} \Rightarrow \frac{S^2(\hat{\beta}_1)}{\sigma^2(\hat{\beta}_1)} = \frac{\hat{\sigma}^2}{\sigma^2}.$$

$$\text{Άρα } \frac{S^2(\hat{\beta}_1)(n-2)}{\sigma^2(\hat{\beta}_1)} \sim X^2(n-2).$$

Επομένως

$$\frac{\frac{\hat{\beta}_1 - \beta_1}{\sigma(\hat{\beta}_1)}}{\sqrt{\frac{S^2(\hat{\beta}_1)(n-2)}{\sigma^2(\hat{\beta}_1)(n-2)}}} = \frac{\hat{\beta}_1 - \beta_1}{s(\hat{\beta}_1)} \sim t_{(n-2)}.$$

**Παρατήρηση:** Η πρόταση αυτή είναι ιδιαίτερα χρήσιμη για στατιστική συμπερασματολογία, καθώς στο αποτέλεσμα αυτό στηρίζεται η κατασκευή διαστημάτων εμπιστοσύνης και η διενέργεια ελέγχων υποθέσεων για τους συντελεστές του γραμμικού μοντέλου.

### Κατασκευή Διαστημάτων Εμπιστοσύνης

Θα κατασκευάσουμε διάστημα εμπιστοσύνης συντελεστή 1- $\alpha$  για το  $\beta_1$ .

$$\text{Είναι } P\left(-t_{\frac{\alpha}{2}}(n-2) \leq \frac{\hat{\beta}_1 - \beta_1}{s(\hat{\beta}_1)} \leq t_{\frac{\alpha}{2}}(n-2)\right) = 1 - \alpha.$$

Άρα το ζητούμενο δ.ε. για το  $\beta_1$  είναι

$$\hat{\beta}_1 \pm t_{\frac{\alpha}{2}}(n-2)s(\hat{\beta}_1).$$

Ομοίως με παραπάνω είναι  $\frac{\hat{\beta}_0 - \beta_0}{s(\hat{\beta}_0)} \sim t(n-2)$ , άρα το διάστημα εεμπιστοσύνης με συντελεστή 1- $\alpha$  για το  $\beta_0$  είναι

$$\hat{\beta}_0 \pm t_{\frac{\alpha}{2}}(n-2)s(\hat{\beta}_0).$$

### Έλεγχοι Υποθέσεων

Χρησιμοποιώντας την κατανομή  $t(n - 2)$  μπορούμε να κάνουμε και ελέγχους υποθέσεων για τις παραμέτρους  $\beta_0$  και  $\beta_1$ .

Η λογική των ελέγχων υποθέσεων είναι να βρεθεί στατιστική συνάρτηση (συνάρτηση του δείγματος που δεν εμπλέκει άγνωστες ποσότητες) με γνωστή κατανομή (κατανομή που δεν έχει άγνωστες παραμέτρους) και για την οποία ακραίες τιμές της αντιστοιχούν σε ένδειξη ότι η  $H_0$  πρέπει να απορριφθεί έναντι της συγκεκριμένης εναλλακτικής ως προς την οποία ελέγχεται.

Οι έλεγχοι

$$\begin{array}{ll} H_0 : \beta_0 = 0 & H_0 : \beta_1 = 0 \\ H_1 : \beta_0 \neq 0 & H_1 : \beta_1 \neq 0 \end{array}$$

ονομάζονται έλεγχοι στατιστικής σημαντικότητας των αντίστοιχων παραμέτρων. Αν σε έναν έλεγχο στατιστικής σημαντικότητας δεν μπορεί να απορριφθεί η αρχική υπόθεση αυτό αντιστοιχεί στην παραδοχή ότι δεν υφίσταται η συγκεκριμένη παράμετρος στο μοντέλο. Στην περίπτωση που ο έλεγχος γίνεται για τον συντελεστή κλίσης,  $\beta_1$ , μη απόρριψη της  $H_0$  οδηγεί στο συμπέρασμα ότι δεν υφίσταται γραμμική σχέση ανάμεσα στην απαντητική μεταβλητή,  $Y$ , και την ερμηνευτική μεταβλητή,  $X$ . Πράγματι, αν  $\beta_1 = 0$ , το μοντέλο υποθέτει ότι  $Y_i \sim N(\beta_0, \sigma^2)$ ,  $i = 1, \dots, n$ , δηλαδή τα  $Y_i$  είναι ανεξάρτητα και ισόνομα.

Για τον έλεγχο στατιστικής σημαντικότητας του  $\beta_1$  χρησιμοποιείται η στατιστική συνάρτηση  $T = \frac{\hat{\beta}_1 - \beta_1}{s(\hat{\beta}_1)}$  (η οποία όπως έχουμε πει ακολουθεί πάντα  $t(n-2)$  κατανομή).

Κάτω από την ισχύ της  $H_0$ , η στατιστική συνάρτηση ελέγχου (ελεγχοσυνάρτηση)  $T = \frac{\hat{\beta}_1}{s(\hat{\beta}_1)} \sim t(n-2)$ .

Για τα δεδομένα του δείγματός μας υπολογίζουμε την παρατηρούμενη τιμή της ελεγχοσυνάρτησης στο δείγμα, έστω  $t^*$ . Απορρίπτουμε την  $H_0$ , σε επίπεδο στατιστικής σημαντικότητας  $\alpha$ , αν  $|t^*| > t_{\alpha/2}(n-2)$  (ισοδύναμα για  $t^* > t_{\alpha/2}(n-2)$  ή  $t^* < -t_{\alpha/2}(n-2)$ ). Δηλαδή, απορρίπτουμε την  $H_0 : \beta_1 = 0$  έναντι της αμφίπλευρης εναλλακτικής, αν η παρατηρούμενη τιμή της ελεγχοσυνάρτησης υπερβαίνει το  $\alpha/2$  άνω ποσοστιαίο σημείο της  $t(n-2)$  κατανομής.

Για τον έλεγχο της  $H_0 : \beta_1 = 3$  έναντι της αμφίπλευρης εναλλακτικής  $H_1 : \beta_1 \neq 3$  η διαδικασία είναι αντίστοιχη. Η ελεγχοσυνάρτηση που χρησιμοποιείται

τώρα είναι η  $T = \frac{\hat{\beta}_1 - 3}{s(\hat{\beta}_1)}$ , η οποία, κάτω από την ισχύ της  $H_0$  ακολουθεί  $t(n - 2)$  κατανομή. Για τα δεδομένα του δείγματός μας υπολογίζουμε την παρατηρούμενη τιμή αυτής της ελεγχοσυνάρτησης στο δείγμα, έστω  $t^*$ . Απορρίπτουμε την  $H_0$ , σε επίπεδο στατιστικής σημαντικότητας  $\alpha$ , αν  $|t^*| > t_{\alpha/2}(n - 2)$ .

Για τον έλεγχο της  $H_0 : \beta_1 = 3$  έναντι της μονόπλευρης εναλλακτικής  $H_1 : \beta_1 > 3$  η ελεγχοσυνάρτηση που χρησιμοποιείται είναι πάλι η  $T = \frac{\hat{\beta}_1 - 3}{s(\hat{\beta}_1)}$ , η οποία, κάτω από την ισχύ της  $H_0$  ακολουθεί  $t(n - 2)$  κατανομή. Απορρίπτουμε την  $H_0$ , σε επίπεδο στατιστικής σημαντικότητας  $\alpha$ , έναντι της συγκεκριμένης εναλλακτικής, αν  $t^* > t_{\alpha}(n - 2)$ .

Για τον έλεγχο της  $H_0 : \beta_1 = 3$  έναντι της μονόπλευρης εναλλακτικής  $H_1 : \beta_1 < 3$  η ελεγχοσυνάρτηση που χρησιμοποιείται είναι πάλι η  $T = \frac{\hat{\beta}_1 - 3}{s(\hat{\beta}_1)}$ , η οποία, κάτω από την ισχύ της  $H_0$  ακολουθεί  $t(n - 2)$  κατανομή. Απορρίπτουμε την  $H_0$ , σε επίπεδο στατιστικής σημαντικότητας  $\alpha$ , έναντι της συγκεκριμένης εναλλακτικής, αν  $t^* < -t_{\alpha}(n - 2)$ .

#### Παρατηρούμενο Επίπεδο Σημαντικότητας (p-value)

Το παρατηρούμενο επίπεδο σημαντικότητας (observed level of significance ή p-value) ορίζεται ως η πιθανότητα να πάρει η ελεγχοσυνάρτηση,  $T$ , τιμή τόσο ή περισσότερο ακραία από την παρατηρούμενη στο δείγμα μας,  $t^*$ .

**Παρατήρηση:** 'Όταν λέμε ακραία στον ορισμό του p-value, εννοούμε ακραία ως προς την συγκεκριμένη εναλλακτική του ελέγχου. 'Ετσι, για αμφίπλευρο έλεγχο,  $p\text{-value} = P(|T| > |t^*|)$ , ενώ για τους μονόπλευρους ελέγχους  $p\text{-value} = P(T > t^*)$  ή  $p\text{-value} = P(T < t^*)$ , αντίστοιχα.

Σε οποιονδήποτε έλεγχο υποθέσεων, απορρίπτουμε την  $H_0$  αν το παρατηρούμενο επίπεδο σημαντικότητας (p-value) είναι χαμηλότερο από το προκαθορισμένο επίπεδο σημαντικότητας,  $\alpha$ . Μάλιστα, όσο μικρότερη είναι η τιμή του p-value τόσο ισχυρότερη ένδειξη αποτελεί αυτό εναντίον της  $H_0$ .

### Άσκηση

Έστω το εναλλακτικό απλό γραμμικό μοντέλο

$$Y_i = \beta_0^* + \beta_1(X_i - \bar{X}) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad \text{ανεξάρτητα.}$$

- (a) Να βρεθούν οι εκτιμητές μέγιστης πιθανοφάνειας των παραμέτρων του μοντέλου.
- (b) Να βρεθεί η κατανομή του  $\hat{\beta}_0^*$ .
- (c) Να κατασκευαστεί διάστημα εμπιστοσύνης συντελεστή  $1 - \alpha$  για το  $\beta_0^*$ .
- (d) Να πραγματοποιηθεί ο έλεγχος στατιστικής σημαντικότητας του συντελεστή  $\beta_0^*$  σε ε.σ.σ  $\alpha$ .

### Λύση του (b)

Το εναλλακτικό γραμμικό μοντέλο της άσκησης είναι ισοδύναμο με το αρχικό απλό γραμμικό μοντέλο.

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, = \beta_0 + \beta_1 \bar{X} + \beta_1(X_i - \bar{X}) + \varepsilon_i = \beta_0^* + \beta_1(X_i - \bar{X}) + \varepsilon_i \quad \text{όπου } \beta_0^* = \beta_0 + \beta_1 \bar{X}.$$

$$\text{Άρα } \hat{\beta}_0^* = \hat{\beta}_0 + \hat{\beta}_1 \bar{X} = \bar{Y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 \bar{X} = \bar{Y}$$

Επειδή τα  $Y_i$  είναι κανονικές τ.μ., το  $\hat{\beta}_0^* = \bar{Y}$  είναι επίσης κανονική τ.μ. ως γραμμικός συνδυασμός κανονικών.

$$E(\hat{\beta}_0^*) = E(\hat{\beta}_0 + \hat{\beta}_1 \bar{X}) = E(\hat{\beta}_0) + \bar{X} E(\hat{\beta}_1) = \beta_0 + \beta_1 \bar{X} = \beta_0^*.$$

Εναλλακτικά

$$\begin{aligned} E(\hat{\beta}_0^*) &= E(\bar{Y}) = E\left(\frac{\sum Y_i}{n}\right) = \frac{1}{n} \sum E(Y_i) = \frac{1}{n} \sum (\beta_0^* + \beta_1(X_i - \bar{X})) \\ &= \frac{1}{n} n \beta_0^* + \frac{1}{n} \underbrace{\sum (X_i - \bar{X})}_0 = \beta_0^*. \end{aligned}$$

Δηλαδή  $\hat{\beta}_0^*$  αμερόληπτη εκτιμήτρια.

$$V(\hat{\beta}_0^*) = \sigma^2(\hat{\beta}_0^*) = V(\bar{Y}) = V\left(\frac{\sum Y_i}{n}\right) = \frac{1}{n^2} \sum V(Y_i) = \frac{1}{n^2} \sum \sigma^2 = \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n}.$$

$$\text{Άρα } \hat{\beta}_0^* \sim N(\beta_0^*, \frac{\sigma^2}{n}).$$

## Κατανομές Τετραγωνικών Μορφών

**Ορισμός:** Ένα ομογενές πολυώνυμο 2ου βαθμού σε η μεταβλητές καλείται τετραγωνική μορφή σε αυτές τις μεταβλητές. Εάν και οι συντελεστές και οι μεταβλητές είναι πραγματικοί αριθμοί τότε καλείται πραγματική τετραγωνική μορφή.

**Σημείωση:** Ομογενές πολυώνυμο είναι αυτό που έχει ίδια δύναμη σε όλους τους όρους του.

## Παραδείγματα

(1)  $X_1^2 + X_1X_2 + X_2^2$  τετραγωνική μορφή ως προς  $X_1, X_2$

(2)  $X_1^2 + X_2^2 + X_3^2 - 2X_1X_2$  τετραγωνική μορφή ως προς  $X_1, X_2, X_3$

$$(3) X_1^2 + X_2^2 - 2X_1 - 4X_2 + 5 \text{ δεν είναι τετραγωνική μορφή ως προς } X_1, X_2 \\ (X_1 - 1)^2 + (X_2 - 2)^2 \quad \text{αλλά είναι ως προς } (X_1 - 1), (X_2 - 2)$$

**Πρόταση.** Έστω  $X_1, X_2, \dots, X_n$  τ.δ. από την τ.μ.  $X$  και  $\bar{X}, S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$  ο δειγματικός μέσος και η δειγματική διασπορά αντίστοιχα. Η τ.μ.  $(n-1)S^2$  είναι μια τετραγωνική μορφή στις  $n$  τ.μ.  $X_1, X_2, \dots, X_n$ .

Απόδειξη.

$$\begin{aligned}
(n-1)S^2 &= \sum(X_i - \bar{X})^2 = \sum(X_i - \frac{\sum X_i}{n})^2 \\
&= \sum X_i^2 + \frac{n(\sum X_i)^2}{n^2} - 2 \sum X_i \frac{\sum X_i}{n} = \\
&= \sum X_i^2 + \frac{(\sum X_i)^2}{n} - \frac{2}{n} (\sum X_i)^2 = \\
&= \sum X_i^2 - \frac{(\sum X_i)^2}{n} = \sum X_i^2 - \frac{1}{n} \left\{ \sum_{\substack{i,j=1 \\ i < j}}^n X_i X_j \right\} =
\end{aligned}$$

$$= \frac{n-1}{n} \sum X_i^2 - \frac{2}{n} \sum_{\substack{i,j=1 \\ i < j}}^n X_i X_j \quad \text{τετραγωνική μορφή}$$

**Παρατήρηση:** Γνωρίζουμε ότι εάν  $X \sim N(\mu, \sigma^2)$  και  $X_1, X_2, \dots, X_n$  τ.δ. τότε  $\frac{(n-1)S^2}{\sigma^2} \sim X_{(n-1)}^2$ .

**Παρατήρηση:** Γνωρίζουμε ότι εάν  $Q_1, Q_2, \dots, Q_k$  είναι ανεξάρτητες  $X^2$  τ.μ. με  $r_i, i = 1, \dots, k$  β.ε. αντίστοιχα, τότε η τ.μ.  $Q = \sum_{i=1}^k Q_i$  είναι  $X_{(r_1+r_2+\dots+r_k)}^2$ .

**Θεώρημα.** 'Εστω  $Q = \sum_{i=1}^k Q_i$  όπου  $Q_1, Q_2, \dots, Q_k$  είναι  $k$  πραγματικές τετραγωνικές μορφές σε  $n$  ανεξάρτητες τ.μ. που ακολουθούν την κανονική κατανομή με μέσους  $\mu_1, \mu_2, \dots, \mu_n$  αντίστοιχα και την ίδια διασπορά  $\sigma^2$ .

'Εστω επίσης ότι (i)  $\frac{Q}{\sigma^2}, \frac{Q_1}{\sigma^2}, \frac{Q_2}{\sigma^2}, \dots, \frac{Q_{k-1}}{\sigma^2}$  έχουν την κατανομή  $X^2$  με  $r, r_1, r_2, \dots, r_{k-1}$  β.ε. αντίστοιχα, και (ii)  $Q_k$  μη αρνητική. Τότε

(1)  $Q_1, Q_2, \dots, Q_k$  είναι ανεξάρτητες και

$$(2) \frac{Q_k}{\sigma^2} \sim X_{(r_k)}^2, \text{ όπου } r_k = r - (r_1 + r_2 + \dots + r_{k-1}).$$

**Πρόταση.** 'Εστω το γραμμικό μοντέλο

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i,$$

ή

$$Y_i = \beta_0^* + \beta_1(X_i - \bar{X}) + \varepsilon_i \text{ με } \beta_0^* = \beta_0 + \beta_1 \bar{X}.$$

Εάν  $\varepsilon_i \sim N(0, \sigma^2)$ , τότε  $Y_i \stackrel{\text{ανεξ.}}{\sim} N(\beta_0 + \beta_1 X_i, \sigma^2)$   
ή  $Y_i \sim N(\beta_0^* + \beta_1(X_i - \bar{X}), \sigma^2)$ .

Ισχύει ότι  $\frac{\sum \hat{\varepsilon}_i^2}{\sigma^2} = \frac{\sum(Y_i - \hat{Y}_i)^2}{\sigma^2} = \frac{\sum(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2}{\sigma^2} \sim X_{(n-2)}^2$   
ή αλλιώς  $\frac{\sum(Y_i - \hat{\beta}_0^* - \hat{\beta}_1(X_i - \bar{X}))^2}{\sigma^2} \sim X_{(n-2)}^2$ .

**Απόδειξη.**

$$\begin{aligned}
 & \sum(Y_i - \beta_0 - \beta_1 X_i)^2 = \sum(Y_i - \beta_0^* - \beta_1(X_i - \bar{X}))^2 = \\
 & = \sum \left[ (\hat{\beta}_0^* - \beta_0^*) + (\hat{\beta}_1 - \beta_1)(X_i - \bar{X}) + [Y_i - \hat{\beta}_0^* - \hat{\beta}_1(X_i - \bar{X})] \right]^2 \\
 & \quad (\text{όπου } \epsilon \text{ έχουμε προσθαφαιρέσει } \hat{Y}_i = \hat{\beta}_0^* + \hat{\beta}_1(X_i - \bar{X})) \\
 & = n(\hat{\beta}_0^* - \beta_0^*)^2 + (\hat{\beta}_1 - \beta_1)^2 \sum(X_i - \bar{X})^2 + \sum \left[ (Y_i - \hat{\beta}_0^* - \hat{\beta}_1(X_i - \bar{X})) \right]^2 + 0 \\
 & \quad (\text{αθροίσματα διπλασίων γινομένων} = 0)
 \end{aligned}$$

Συμβολίζουμε με

$$\begin{aligned}
 Q &= \sum(Y_i - \beta_0 - \beta_1 X_i)^2 = \sum \left( Y_i - \beta_0^* - \beta_1(X_i - \bar{X}) \right)^2 \\
 Q_1 &= n(\hat{\beta}_0^* - \beta_0^*)^2 \\
 Q_2 &= (\hat{\beta}_1 - \beta_1)^2 \sum(X_i - \bar{X})^2 \\
 Q_3 &= \sum \left[ Y_i - \hat{\beta}_0^* - \hat{\beta}_1(X_i - \bar{X}) \right]^2
 \end{aligned}$$

Είναι  $Q = Q_1 + Q_2 + Q_3$ .

Τα  $Q, Q_1, Q_2, Q_3$  είναι τετραγωνικές μορφές ως προς κανονικές τ.μ.

$$\text{Επίσης } \epsilon \text{ έχουμε } Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2) \Rightarrow \frac{(Y_i - \beta_0 - \beta_1 X_i)^2}{\sigma^2} \sim X_{(1)}^2.$$

$$\text{Οι } Y_i \text{ είναι ανεξάρτητες άρα } \sum_{i=1}^n \frac{(Y_i - \beta_0 - \beta_1 X_i)^2}{\sigma^2} \sim X_{(n)}^2,$$

$$\delta \eta \lambda \cdot \frac{Q}{\sigma^2} \sim X_{(n)}^2.$$

$$\hat{\beta}_0^* \sim N(\beta_0, \frac{\sigma^2}{n}) \Rightarrow \frac{n(\hat{\beta}_0^* - \beta_0^*)^2}{\sigma^2} \sim X_{(1)}^2, \delta \eta \lambda \cdot \frac{Q_1}{\sigma^2} \sim X_{(1)}^2.$$

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}) \Rightarrow \frac{\sum_{i=1}^n (X_i - \bar{X})^2 (\hat{\beta}_1 - \beta_1)^2}{\sigma^2} = \frac{Q_2}{\sigma^2} \sim X_{(1)}^2.$$

To  $Q_3$  είναι μη αρνητικό, άρα με βάση το θεώρημα

$$\frac{Q_3}{\sigma^2} = \frac{\sum \left[ Y_i - \hat{\beta}_0^* - \hat{\beta}_1(X_i - \bar{X}) \right]^2}{\sigma^2} = \frac{\sum \hat{\epsilon}_i^2}{\sigma^2} \sim X_{(n-1-1)}^2 \equiv X_{(n-2)}^2.$$

Μένει να δείξουμε ότι τα αθροίσματα διπλασίων γινομένων είναι μηδενικά.

$$\begin{aligned}
 (i) \quad & 2 \sum (\hat{\beta}_0^* - \beta_0^*)(\hat{\beta}_1 - \beta_1)(X_i - \bar{X}) = 2(\hat{\beta}_0^* - \beta_0^*)(\hat{\beta}_1 - \beta_1) \sum (X_i - \bar{X}) = 0 \\
 (ii) \quad & 2 \sum (\hat{\beta}_0^* - \beta_0^*) \left( Y_i - \hat{\beta}_0^* - \hat{\beta}_1(X_i - \bar{X}) \right) = \\
 & = 2(\hat{\beta}_0^* - \beta_0^*) \sum (Y_i - \bar{Y}) - 2(\hat{\beta}_0^* - \beta_0^*) \hat{\beta}_1 \sum (X_i - \bar{X}) = 0 \\
 (iii) \quad & 2 \sum (\hat{\beta}_1 - \beta_1)(X_i - \bar{X}) \left( Y_i - \hat{\beta}_0^* - \hat{\beta}_1(X_i - \bar{X}) \right) = \\
 & = 2(\hat{\beta}_1 - \beta_1) \sum (X_i - \bar{X})(Y_i - \bar{Y}) - 2(\hat{\beta}_1 - \beta_1) \hat{\beta}_1 \sum (X_i - \bar{X})^2 = \\
 & = 2(\hat{\beta}_1 - \beta_1) \left[ \sum (X_i - \bar{X})(Y_i - \bar{Y}) - \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2} \sum (X_i - \bar{X})^2 \right] = 0.
 \end{aligned}$$

### Ασκηση

Δίνεται το γραμμικό μοντέλο

$$Y_i = \beta(X_i - \bar{X}) + \varepsilon_i, i = 1, \dots, n, \varepsilon_i \sim N(0, \sigma^2) \text{ ασυσχέτιστα.}$$

- (i) Να βρεθεί η ε.ε.τ.  $\hat{\beta}$  της  $\beta$  και η κατανομή της.
- (ii) Να βρεθεί η κατανομή της στατιστικής συνάρτησης  $\frac{\hat{\beta} - \beta}{s(\hat{\beta})}$  όπου  $s^2(\hat{\beta})$  είναι α.ε. της διασποράς του  $\hat{\beta}$ .
- (iii) Για δεδομένη τιμή  $X_0$  να βρεθεί δ.ε. συντελεστή 1- $\alpha$  για το  $\beta(X_0 - \bar{X})$ .

### Λύση

- (i) Εκτίμηση ελαχίστων τετραγώνων: ελαχιστοποιούμε την ποσότητα  $Q = \sum \varepsilon_i^2 = \sum (Y_i - \beta(X_i - \bar{X}))^2$  ως προς  $\beta$ .

$$\frac{\partial Q}{\partial \beta} = -2 \sum (X_i - \bar{X})[Y_i - \beta(X_i - \bar{X})] = 0$$

$$\Rightarrow - \sum (X_i - \bar{X})Y_i + \beta \sum (X_i - \bar{X})^2 = 0 \Rightarrow \hat{\beta} = \frac{\sum (X_i - \bar{X})Y_i}{\sum (X_i - \bar{X})^2}$$

To  $\hat{\beta}$  είναι γραμμικός συνδυασμός των  $Y_i$ , αφού

$$\hat{\beta} = \sum k_i Y_i, \text{ όπου } k_i = \frac{X_i - \bar{X}}{\sum (X_i - \bar{X})^2} \text{ σταθερές.}$$

Επειδή  $\varepsilon_i \sim N(0, \sigma^2)$  ασυσχέτιστες άρα και ανεξάρτητες τ.μ., ισχύει ότι

$$Y_i \sim N(\beta(X_i - \bar{X}), \sigma^2).$$

Άρα το  $\hat{\beta}$  ακολουθεί κανονική κατανομή σαν γραμμικός συνδυασμός ανεξάρτητων κανονικών τ.μ.

Είναι

$$\begin{aligned} E(\hat{\beta}) &= E\left[\frac{\sum(X_i - \bar{X})Y_i}{\sum(X_i - \bar{X})^2}\right] = \frac{\sum(X_i - \bar{X})E(Y_i)}{\sum(X_i - \bar{X})^2} \\ &= \frac{\sum(X_i - \bar{X})\beta(X_i - \bar{X})}{\sum(X_i - \bar{X})^2} = \frac{\beta \sum(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} = \beta \end{aligned}$$

Δηλαδή είναι α.ε. του  $\beta$ .

$$V(\hat{\beta}) = V\left[\frac{\sum(X_i - \bar{X})Y_i}{\sum(X_i - \bar{X})^2}\right] = \frac{\sum(X_i - \bar{X})^2V(Y_i)}{\left(\sum(X_i - \bar{X})^2\right)^2} = \frac{\sigma^2 \sum(X_i - \bar{X})^2}{\left(\sum(X_i - \bar{X})^2\right)^2} = \frac{\sigma^2}{\sum(X_i - \bar{X})^2}$$

Άρα

$$\hat{\beta} \sim N(\beta, \frac{\sigma^2}{\sum(X_i - \bar{X})^2})$$

$$(ii) \text{ Θα δείξουμε πρώτα ότι το } \frac{\sum(Y_i - \hat{\beta}(X_i - \bar{X}))^2}{\sigma^2} \sim X_{(n-1)}^2.$$

$$\begin{aligned} \sum(Y_i - \beta(X_i - \bar{X}))^2 &= \sum(Y_i - \hat{\beta}(X_i - \bar{X}) + \hat{\beta}(X_i - \bar{X}) - \beta(X_i - \bar{X}))^2 = \\ &= \sum\left[(\hat{\beta} - \beta)(X_i - \bar{X}) + (Y_i - \hat{\beta}(X_i - \bar{X}))\right]^2 = \\ &= (\hat{\beta} - \beta)^2 \sum(X_i - \bar{X})^2 + \sum(Y_i - \hat{\beta}(X_i - \bar{X}))^2 + 2(\hat{\beta} - \beta) \sum(X_i - \bar{X})(Y_i - \hat{\beta}(X_i - \bar{X})) = \\ &= (\hat{\beta} - \beta)^2 \sum(X_i - \bar{X})^2 + \sum(Y_i - \hat{\beta}(X_i - \bar{X}))^2 \\ \text{γιατί} \quad \sum(X_i - \bar{X})(Y_i - \hat{\beta}(X_i - \bar{X})) &= \sum(X_i - \bar{X})Y_i - \hat{\beta} \sum(X_i - \bar{X})^2 \\ &= \sum(X_i - \bar{X})Y_i - \frac{\sum(X_i - \bar{X})Y_i}{\sum(X_i - \bar{X})^2} \sum(X_i - \bar{X})^2 = 0 \end{aligned}$$

'Αρα έχουμε

$$Q = Q_1 + Q_2, \text{ όπου}$$

$$\begin{aligned} Q &= \sum \left( Y_i - \beta(X_i - \bar{X}) \right)^2 && \text{πραγματικές} \\ Q_1 &= (\hat{\beta} - \beta)^2 \sum (X_i - \bar{X})^2 && \text{τετραγωνικές μορφές} \\ Q_2 &= \sum (Y_i - \hat{\beta}(X_i - \bar{X}))^2 \end{aligned}$$

Επειδή  $Y_i \sim N(\beta(X_i - \bar{X}), \sigma^2)$  έχουμε

$$\frac{[Y_i - \beta(X_i - \bar{X})]^2}{\sigma^2} \sim X_{(1)}^2$$

$$\Rightarrow \frac{Q}{\sigma^2} = \frac{\sum [Y_i - \beta(X_i - \bar{X})]^2}{\sigma^2} \sim X_{(n)}^2$$

$$\hat{\beta} \sim N(\beta, \frac{\sigma^2}{\sum (X_i - \bar{X})^2}) \Rightarrow \frac{(\hat{\beta} - \beta)^2 \sum (X_i - \bar{X})^2}{\sigma^2} = \frac{Q_1}{\sigma^2} \sim X_{(1)}^2$$

Επειδή  $Q_2$  είναι μη αρνητικό, από το θεώρημα τετραγωνικών μορφών έχουμε

$$\frac{Q_2}{\sigma^2} = \frac{\sum (Y_i - \hat{\beta}(X_i - \bar{X}))^2}{\sigma^2} \sim X_{(n-1)}^2$$

Ξέρουμε ότι αν  $W \sim X_{(n-1)}^2$ , τότε  $E(W) = n - 1$ .

$$'Αρα  $E\left[\frac{\sum(Y_i - \hat{\beta}(X_i - \bar{X}))^2}{\sigma^2}\right] = n - 1 \Rightarrow E\left[\frac{\sum(Y_i - \hat{\beta}(X_i - \bar{X}))^2}{n - 1}\right] = \sigma^2$ .$$

$$\begin{aligned} \text{'Αρα } \eta \hat{\sigma}^2 &= \frac{\sum(Y_i - \hat{\beta}(X_i - \bar{X}))^2}{n - 1} \text{ είναι α.ε. τού } \sigma^2 \\ \text{και } \eta S^2(\hat{\beta}) &= \frac{\hat{\sigma}^2}{\sum(X_i - \bar{X})^2} \text{ α.ε. τού } \sigma^2(\hat{\beta}) = \frac{\sigma^2}{\sum(X_i - \bar{X})^2}. \end{aligned}$$

$$\begin{aligned} \text{'Εχουμε } \beta \text{ ότι } \frac{\sum[Y_i - \beta(X_i - \bar{X})]^2}{\sigma^2} &= \frac{(n - 1)\hat{\sigma}^2}{\sigma^2} \sim X_{(n-1)}^2 \\ \text{'Αρα } \frac{(n - 1)S^2(\hat{\beta})}{\sigma^2(\hat{\beta})} &\sim X_{(n-1)}^2 \end{aligned}$$

και η κατανομή της είναι ανεξάρτητη του  $\hat{\beta}$ .

'Αρα

$$T = \frac{\frac{\hat{\beta} - \beta}{\sigma(\hat{\beta})}}{\sqrt{\frac{(n-1)S^2(\hat{\beta})}{\sigma^2(\hat{\beta})(n-1)}}} = \frac{\hat{\beta} - \beta}{s(\hat{\beta})} \sim t(n-1)$$

(iii) Ζητάμε Δ.Ε. για το  $\beta(X_0 - \bar{X})$ .

$$\text{Είναι } P\left(-t_{\frac{\alpha}{2}}(n-1) < \frac{\hat{\beta} - \beta}{s(\hat{\beta})} < t_{\frac{\alpha}{2}}(n-1)\right) = 1 - \alpha$$

'Αρα το ζητούμενο Δ.Ε. για το  $\beta$  με συντελεστή  $1-\alpha$  είναι

$$\hat{\beta} \pm t_{\frac{\alpha}{2}}(n-1)s(\hat{\beta})$$

οπότε το Δ.Ε. για το  $\beta(X_0 - \bar{X})$  είναι

$$\hat{\beta}(X_0 - \bar{X}) \pm t_{\frac{\alpha}{2}}(n-1)s(\hat{\beta})(X_0 - \bar{X})$$