

ΕΚΤΙΜΗΣΗ, Δ.Ε. και
ΕΛΕΓΧΟΙ ΥΠΟΘΕΣΕΩΝ για την $E(Y) = \beta_0 + \beta_1 \cdot X$

Εκτιμήτρια της $E(Y)$ είναι η $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 \cdot X$.

Ιδιότητες της εκτιμήτριας: (Υποθέτω ότι έχω κανονικά και ανεξάρτητα σφάλματα)

$$E(\hat{Y}) = E(\hat{\beta}_0 + \hat{\beta}_1 \cdot X) = E(\hat{\beta}_0) + E(\hat{\beta}_1) \cdot X = \beta_0 + \beta_1 \cdot X$$

$$Var(\hat{Y}) = Var(\hat{\beta}_0 + \hat{\beta}_1 \cdot X) = Var(\hat{\beta}_0) + X^2 \cdot Var(\hat{\beta}_1) + 2 \cdot X \cdot Cov(\hat{\beta}_0, \hat{\beta}_1) =$$

$$= \sigma^2 \cdot \frac{\sum_{i=1}^n X_i^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} + X^2 \cdot \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} - 2 \cdot X \cdot \sigma^2 \cdot \frac{\bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} =$$

$$= \frac{\sigma^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} \left(\sum_{i=1}^n X_i^2 + n \cdot X^2 - 2 \cdot n \cdot X \cdot \bar{X} \right) =$$

$$= \frac{\sigma^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} \left(\sum_{i=1}^n X_i^2 + n \cdot X^2 - 2 \cdot n \cdot X \cdot \bar{X} + n \cdot \bar{X}^2 - n \cdot \bar{X}^2 \right) =$$

$$= \frac{\sigma^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} \left[\left(\sum_{i=1}^n X_i^2 - n \cdot \bar{X}^2 \right) + \left(n \cdot X^2 - 2 \cdot n \cdot X \cdot \bar{X} + n \cdot \bar{X}^2 \right) \right] =$$

$$= \frac{\sigma^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} \left[\sum_{i=1}^n (X_i - \bar{X})^2 + n \cdot (X - \bar{X})^2 \right] = \sigma^2 \cdot \left[\frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]$$

Αφού λόγω της υπόθεσης της κανονικότητας, τα $\hat{\beta}_0, \hat{\beta}_1$ είναι γραμμικοί συνδυασμοί κανονικών, τότε και η \hat{Y} θα ακολουθεί την κανονική κατανομή σαν γραμμικός συνδυασμός των $\hat{\beta}_0, \hat{\beta}_1$, δηλαδή,

$$\hat{Y} \sim N \left(\beta_0 + \beta_1 \cdot X, \left[\frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2 \right)$$

Οπότε, έχουμε σαν αντιστρεπτή ποσότητα την $\frac{\hat{Y} - (\beta_0 + \beta_1 \cdot X)}{\sqrt{\hat{Var}(\hat{Y})}} \sim t_{n-2}$,

$$\text{όπου } \hat{Var}(\hat{Y}) = \left[\frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \hat{\sigma}^2.$$

Αρα ένα $(1-\alpha)\%$ Δ.Ε. για την $E(Y)$ είναι:

$$\left(\hat{Y} - t_{n-2, \alpha/2} \cdot \sqrt{\hat{Var}(\hat{Y})}, \hat{Y} + t_{n-2, \alpha/2} \cdot \sqrt{\hat{Var}(\hat{Y})} \right)$$

Για τον έλεγχο της υπόθεσης $H_0 : \beta_0 + \beta_1 \cdot X = v$ προς $H_1 : \beta_0 + \beta_1 \cdot X \neq v$, σε επίπεδο σημαντικότητας α , απορρίπτουμε την $H_0 : \beta_0 + \beta_1 \cdot X = v$ όταν $|t| > t_{n-2, \alpha/2}$,

$$\text{όπου } t = \frac{\hat{Y} - v}{\sqrt{\hat{Var}(\hat{Y})}}.$$

ΔΙΑΣΤΗΜΑ ΠΡΟΒΛΕΨΗΣ για την $E(Y) = \beta_0 + \beta_1 \cdot X$

Σε αυτό το σημείο θέλουμε να κάνουμε πρόβλεψη της τιμής Y για κάποιο μελλοντικό X .

Χρησιμοποιούμε το εκτιμώμενο μοντέλο $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 \cdot X$.

$$\text{Με βάση αυτό το μοντέλο, έχουμε ότι } \hat{Y} \sim N \left(\beta_0 + \beta_1 \cdot X, \left[\frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2 \right).$$

Σκοπός μας είναι να ελαχιστοποιήσουμε το σφάλμα στην πρόβλεψη.

Από τη διακύμανση του εκτιμητή, μπορεί να φανεί και το μέγεθος του σφάλματος:

(α) Αν η τιμή της X είναι κοντά στο \bar{X} , τότε αναμένουμε να έχουμε μικρό σφάλμα στην πρόβλεψη.

(β) Αν το εύρος των X_i είναι μεγάλο, τότε και το $\sum_{i=1}^n (X_i - \bar{X})^2$ θα είναι μεγάλο, άρα αναμένεται μικρό σφάλμα στην πρόβλεψη.

(γ) Όταν το πλήθος των παρατηρήσεων είναι μεγάλο, τότε πάλι αναμένουμε μικρό σφάλμα στην πρόβλεψη.

Βασικός μας στόχος είναι να δημιουργήσουμε ένα διάστημα πρόβλεψης (Δ.Π.) για την πραγματική τιμή της Y .

Έχουμε ότι $Y - \hat{Y} = Y - \hat{\beta}_0 - \hat{\beta}_1 \cdot X$.

Επίσης, $E(Y - \hat{Y}) = E(Y) - E(\hat{Y}) = \beta_0 + \beta_1 \cdot X - \hat{\beta}_0 - \hat{\beta}_1 \cdot X = 0$ και

$$\begin{aligned} \text{Var}(Y - \hat{Y}) &= \text{Var}(Y) + \text{Var}(\hat{Y}) - 2 \cdot \text{Cov}(Y, \hat{Y}) = \sigma^2 + \sigma^2 \cdot \left[\frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] - 0 = \\ &= \sigma^2 \cdot \left[1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \end{aligned}$$

Σημείωση: Γνωρίζουμε ότι $\text{Cov}(Y_i - Y_j) = 0, i \neq j$. Η \hat{Y} εκτιμήθηκε βάσει τις $Y_i, i = 1, 2, \dots, n$ παρατηρήσεις και η Y είναι μια καινούρια παρατήρηση. Άρα η Y είναι ανεξάρτητη από τις $Y_i, i = 1, 2, \dots, n$ παρατηρήσεις, άρα $\text{Cov}(Y, \hat{Y}) = 0$.

Επομένως, $Y - \hat{Y} \sim N \left(0, \sigma^2 \cdot \left[1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \right)$, σαν γραμμικός συνδυασμός κανονικών.

Άρα έχουμε σαν αντιστρεπτή ποσότητα την $t = \frac{Y - \hat{Y}}{\hat{\sigma} \cdot \sqrt{1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}} \sim t_{n-2}$.

Οπότε ισχύει ότι, $P \left(-t_{n-2, \alpha/2} < \frac{Y - \hat{Y}}{\hat{\sigma} \cdot \sqrt{1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}} < t_{n-2, \alpha/2} \right) = 1 - \alpha$

και άρα ένα $(1-\alpha)\%$ Δ.Π. για την τιμή της Y είναι:

$$\left(\hat{Y} - t_{n-2, \alpha/2} \cdot \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}, \hat{Y} + t_{n-2, \alpha/2} \cdot \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}} \right)$$

Έλεγχος έλλειψης προσαρμογής:

Όταν υιοθετούμε το μοντέλο $Y_i = \beta_0 + \beta_1 \cdot X_i + \varepsilon_i, i=1,2,\dots,n$, δεν σημαίνει ότι η σχέση μεταξύ X και Y είναι γραμμική.

Για παράδειγμα, θα μπορούσε η πραγματική σχέση μεταξύ X και Y να ήταν $Y_i = \beta_0^* + \beta_1^* \cdot X_i + \beta_2^* \cdot X_i^2 + \varepsilon_i, i=1,2,\dots,n$.

Σε ορισμένες περιπτώσεις μπορούμε να ελέγξουμε κατά πόσο το μοντέλο είναι σωστό.

Συνέπειες από την υιοθέτηση λανθασμένου μοντέλου:

(α) Οι υποθέσεις για τα σφάλματα παραμένουν οι ίδιες.

(β) Οι υποθέσεις για τα Y_i γίνονται ως εξής:

$$E(Y_i) \neq \beta_0 + \beta_1 \cdot X_i$$

$$Cov(Y_i, Y_j) = 0, i \neq j$$

$$Var(Y_i) = \sigma^2$$

$$Y_i \sim N(E(Y_i), \sigma^2)$$

(γ) Το εκτιμώμενο μοντέλο θα παραμείνει το ίδιο: $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot X_i$

(δ) Γνωρίζουμε ότι για το σωστό μοντέλο ισχύει: $E(SS_{res}) = (n-2) \cdot \sigma^2$

Για το λανθασμένο μοντέλο ισχύει:

$$Y_i - \hat{Y}_i = [(Y_i - \hat{Y}_i) - E(Y_i - \hat{Y}_i)] + E(Y_i - \hat{Y}_i) = q_i + B_i, \text{ όπου}$$

$$q_i = [(Y_i - \hat{Y}_i) - E(Y_i - \hat{Y}_i)] \text{ και } B_i = E(Y_i - \hat{Y}_i).$$

$$\text{Άρα, } SS_{res} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (q_i + B_i)^2 = \sum_{i=1}^n q_i^2 + \sum_{i=1}^n B_i^2 + 2 \cdot \sum_{i=1}^n q_i \cdot B_i.$$

$$\text{Άρα, } E(SS_{res}) = \sum_{i=1}^n E(q_i^2) + \sum_{i=1}^n B_i^2 + 2 \cdot \sum_{i=1}^n B_i \cdot E(q_i) = \sum_{i=1}^n E(q_i^2) + \sum_{i=1}^n B_i^2,$$

αφού $E(q_i) = 0$.

$$\text{Αλλά, } E(q_i^2) = Var(Y_i - \hat{Y}_i) = Var(Y_i) + Var(\hat{Y}_i) - 2 \cdot Cov(Y_i, \hat{Y}_i) =$$

$$\begin{aligned}
&= \sigma^2 + \left[\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2 - 2 \cdot \text{Cov}(Y_i, \bar{Y} + \hat{\beta}_1 \cdot (X_i - \bar{X})) = \\
&= \sigma^2 + \left[\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2 - 2 \cdot \text{Cov}(Y_i, \bar{Y}) - 2 \cdot (X_i - \bar{X}) \cdot \text{Cov}(Y_i, \hat{\beta}_1) = \\
&= \sigma^2 + \left[\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2 - 2 \cdot \frac{\sigma^2}{n} - 2 \cdot (X_i - \bar{X}) \cdot \text{Cov}\left(Y_i, \sum_{i=1}^n K_i \cdot Y_i\right) = \\
&= \sigma^2 + \left[\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2 - 2 \cdot \frac{\sigma^2}{n} - 2 \cdot \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \cdot \sigma^2 = \\
&= \left[1 - \frac{1}{n} - \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2
\end{aligned}$$

Άρα τελικά,

$$\begin{aligned}
E(SS_{res}) &= \sum_{i=1}^n \left[1 - \frac{1}{n} - \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \cdot \sigma^2 + \sum_{i=1}^n B_i^2 = (n-1) \cdot \sigma^2 + \sum_{i=1}^n B_i^2 \Rightarrow \\
\Rightarrow E(MS_{res}) &= \sigma^2 + \frac{1}{n-2} \cdot \sum_{i=1}^n B_i^2
\end{aligned}$$

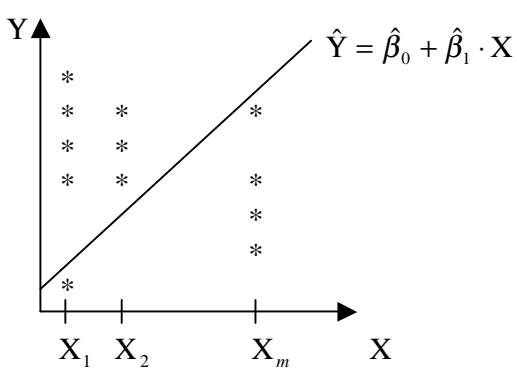
Άρα για το λανθασμένο μοντέλο, η ποσότητα MS_{res} δεν εκτιμά την σ^2 , αλλά μια ποσότητα μεγαλύτερη.

Αυτό επηρεάζει τις διακυμάνσεις των $\hat{\beta}_0, \hat{\beta}_1$, αφού $Var(\hat{\beta}_1) = \frac{MS_{res}}{\sum_{i=1}^n (X_i - \bar{X})^2}$.

Γι' αυτό το λόγο χρησιμοποιούμε τον παρακάτω έλεγχο ορθότητας υποθέτοντας ότι έχουμε επαναλαμβανόμενες παρατηρήσεις:

$$\begin{array}{llll}
X = X_1 & Y_{11}, Y_{12}, \dots, Y_{1n_1} & (\text{τ.δ.}) & \bar{Y}_1 \\
X = X_2 & Y_{21}, Y_{22}, \dots, Y_{2n_2} & (\text{τ.δ.}) & \bar{Y}_2
\end{array}$$

$$\begin{array}{ccc}
 \vdots & \vdots & \vdots \\
 \vdots & \vdots & \vdots \\
 X = X_m & Y_{m1}, Y_{m2}, \dots, Y_{mm_m} & (\text{τ.δ.}) & \bar{Y}_m
 \end{array}$$



Θεωρώ το άθροισμα των υπολοίπων:

$$\begin{aligned}
 SS_{res} &= \sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{ij} - \hat{Y}_i)^2 = \sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i + \bar{Y}_i - \hat{Y}_i)^2 = \\
 &= \sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 + \sum_{i=1}^m \sum_{j=1}^{n_i} (\bar{Y}_i - \hat{Y}_i)^2 + 2 \cdot \sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)(\bar{Y}_i - \hat{Y}_i) = \\
 &= \sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 + \sum_{i=1}^m n_i \cdot (\bar{Y}_i - \hat{Y}_i)^2 + 2 \cdot \sum_{i=1}^m (\bar{Y}_i - \hat{Y}_i) \cdot \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i) = \\
 &= \sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 + \sum_{i=1}^m n_i \cdot (\bar{Y}_i - \hat{Y}_i)^2
 \end{aligned}$$

Παρατηρώντας ότι $\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$ εκφράζει την μεταβλητότητα του i δείγματος,

έχουμε ότι $\sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$ είναι το άθροισμα των τετραγώνων που οφείλεται σε γνήσιο σφάλμα (SS of pure error, SS_{pe}).

Για τον δεύτερο όρο στο παραπάνω άθροισμα, έχω ότι $E(\bar{Y}_i) = \beta_0 + \beta_1 \cdot X$ όταν ισχύει το μοντέλο, αλλά $E(\hat{Y}_i) = E(\hat{\beta}_0 + \hat{\beta}_1 \cdot X_i) = \beta_0 + \beta_1 \cdot X_i$.

Άρα για το σωστό μοντέλο τα \bar{Y}_i είναι κοντά στα \hat{Y}_i .

Επομένως για το σωστό μοντέλο, αν το $\sum_{i=1}^m n_i \cdot (\bar{Y}_i - \hat{Y}_i)^2$ είναι μικρό, τότε έχω ένδειξη για την ορθότητα του μοντέλου. Το παραπάνω άθροισμα ονομάζεται άθροισμα τετραγώνων που οφείλεται στην έλλειψη προσαρμογής (SS of lack of fit, SS_{lof}).

Γενικά έχω ότι: $SS_{res} = SS_{pe} + SS_{lof}$

$$(n-2) \quad (n-m) \quad (m-2) \quad (\text{βαθμοί ελευθερίας})$$

(Όπου οι β.ε. $(n-m) = \sum_{i=1}^m (n_i - 1)$, αφού $n = \sum_{i=1}^m n_i$.)

Όταν ισχύει η κανονικότητα, μπορεί εύκολα να δειχθεί ότι $\frac{SS_{lof}}{\sigma^2} \sim X_{m-2}^2$.

Άρα απορρίπτω την H_0 : "το μοντέλο είναι σωστό", αν $\frac{SS_{lof}}{\sigma^2} > c$.

Επειδή όμως η παράμετρος σ^2 είναι άγνωστη, την εκτιμούμε με $\hat{\sigma}^2 = MS_{pe} = \frac{SS_{pe}}{n-m}$,

οπότε θεωρούμε την ελεγχοσυνάρτηση $F = \frac{MS_{lof}}{MS_{pe}} = \frac{SS_{lof}/m-2}{SS_{pe}/n-m} \sim F_{m-2, n-m}$.

Άρα απορρίπτεται η H_0 , σε επίπεστο σημαντικότητας α , αν $F > F_{m-2, n-m, \alpha}$.

Πρόταση: $E(MS_{pe}) = \sigma^2$, δηλαδή MS_{pe} είναι αμερόληπτη εκτιμήτρια της σ^2 .

Απόδειξη: Παραλείπεται.

ΠΙΝΑΚΑΣ ΑΝΑΔΙΑ

Πηγή Μεταβλητότητας	Άθροισμα Τετραγώνων	Βαθμοί Ελευθερίας	Μέσα Τετράγωνα	F - test
Παλινδρόμηση	SS_{reg}	1	$MS_{reg} = \frac{SS_{reg}}{1}$	$F = \frac{MS_{reg}}{MS_{res}}$
Υπόλοιπα	SS_{res}	$n-2$	$MS_{res} = \frac{SS_{res}}{n-2}$	
Έλλειψη Προσαρμογής	SS_{lof}	$m-2$	$MS_{lof} = \frac{SS_{lof}}{m-2}$	$F = \frac{MS_{lof}}{MS_{pe}}$
Γνήσιο Σφάλμα	SS_{pe}	$n-m$	$MS_{pe} = \frac{SS_{pe}}{n-m}$	
Ολική Μεταβλητότητα	SS_{tot}	$n-1$		

Παρατηρήσεις:

1) Πρώτα εκτελώ το δεύτερο F - test του πίνακα ΑΝΑΔΙΑ, $F = \frac{MS_{lof}}{MS_{pe}}$, για τον έλεγχο ορθότητας του μοντέλου. Αν δεχτώ αυτό το F - test, τότε εκτελώ και το

πρώτο F -test του πίνακα ANADIA, $F = \frac{MS_{reg}}{MS_{res}}$, για τον έλεγχο ύπαρξης ανεξαρτησίας μεταξύ των X και των παρατηρήσεων Y .

2) Η παραπάνω ανάλυση είναι δυνατή μόνο αν έχω επαναλαμβανόμενες μετρήσεις.

3) Η ποσότητα MS_{pe} είναι πάντα αμερόληπτη εκτιμήτρια για τη σ^2 , ανεξαρτήτως του μοντέλου. Αυτό ισχύει γιατί οι επαναλαμβανόμενες μετρήσεις, επιτρέπουν την εύρεση εκτιμήτριας που να μην εξαρτάται από το μοντέλο.

4) Υπάρχει αναλογία μεταξύ: (α) $E(SS_{res}) = (n-2) \cdot \sigma^2 + \sum_{i=1}^n B_i^2$
και (β) $SS_{res} = SS_{pe} + SS_{lof}$.