

## ΑΠΛΗ ΠΑΛΙΝΔΡΟΜΗΣΗ

Θεωρώ ότι έχω μόνο μία ανεξάρτητη μεταβλητή  $X$  και μία εξαρτημένη μεταβλητή  $Y$ .

Αν χρησιμοποιούμε περισσότερες από μία ανεξάρτητες μεταβλητές  $(X_1, X_2, \dots, X_n)$ , τότε έχουμε πολλαπλή παλινδρόμηση. (Θα μελετηθεί σε επόμενα μαθήματα.)

Υποθέτουμε ότι: Η  $Y$  είναι τυχαία μεταβλητή.  
Η  $X$  παίρνει προκαθορισμένες τιμές, δηλαδή δεν είναι τυχαία μεταβλητή. (Αν θεωρήσουμε την  $X$  ως τυχαία μεταβλητή, τότε θεωρούμε πως την έχουμε σταθεροποιήσει σε συγκεκριμένες τιμές.)

Διαθέτουμε δεδομένα  $(X_i, Y_i)$ ,  $i = 1, 2, \dots, n$ , έτσι ώστε  $Y_i$  να είναι η παρατήρηση που αντιστοιχεί στην τιμή  $X = x_i$ .

Πρόβλημα: Θέλουμε να προσδιορίσουμε τη βέλτιστη σχέση που συνδέει την  $Y$  με τη  $X$ . (Για την ακρίβεια, τη σχέση μεταξύ της  $Y$ , δεδομένου της  $X$ , δοθέντος.)

Η συνάρτηση που ψάχνουμε θα είναι της μορφής  $g(Y/X = x) = f(x)$ . Η βέλτιστη αυτή σχέση μπορεί να βρεθεί ελαχιστοποιώντας το μέσο τετραγωνικό σφάλμα μεταξύ  $Y$  και  $f(x)$ .

$$\text{Άρα } \min_f E[Y - f(x)] \Rightarrow f(x) = E(Y/X = x).$$

Με βάση τα δεδομένα μας, μπορούμε να εκτιμήσουμε την  $f(x)$  με  $\hat{f}(x)$ , δηλαδή εκτιμήσαμε την  $E(Y/X = x) = f(x)$  με  $f(x) = \hat{E}(Y/X = x) = \hat{f}(x)$ .

Θέτοντας  $Y = E(Y/X = x) = f(x)$ , καταλήγουμε στο συμπέρασμα ότι κατά μέσο όρο, η  $Y$  έχει τιμή  $\hat{f}(x)$ .

Για την εκτίμηση της  $f(x)$  κάνουμε γραφική παράσταση των  $(X_i, Y_i)$ ,  $i = 1, 2, \dots, n$ , με βάση την οποία υιοθετούμε κάποιο μοντέλο για την  $E(Y/X)$ .

$$\begin{aligned} \text{π.χ. } E(Y/X) &= \beta_0 + \beta_1 \cdot X && (\text{γραμμικό μοντέλο}) \\ E(Y/X) &= \exp(\beta_0 + \beta_1 \cdot X) && (\text{εκθετικό μοντέλο}) \end{aligned}$$

( Σημείωση: Για ευκολία, θα χρησιμοποιούμε  $Y$  αντί  $Y/X$  και  $E(Y)$  αντί  $E(Y/X)$ .)

ΑΠΛΟ ΓΡΑΜΜΙΚΟ ΜΟΝΤΕΛΟ ή  
ΑΠΛΗ ΓΡΑΜΜΙΚΗ ΠΑΛΙΝΔΡΟΜΗΣΗ

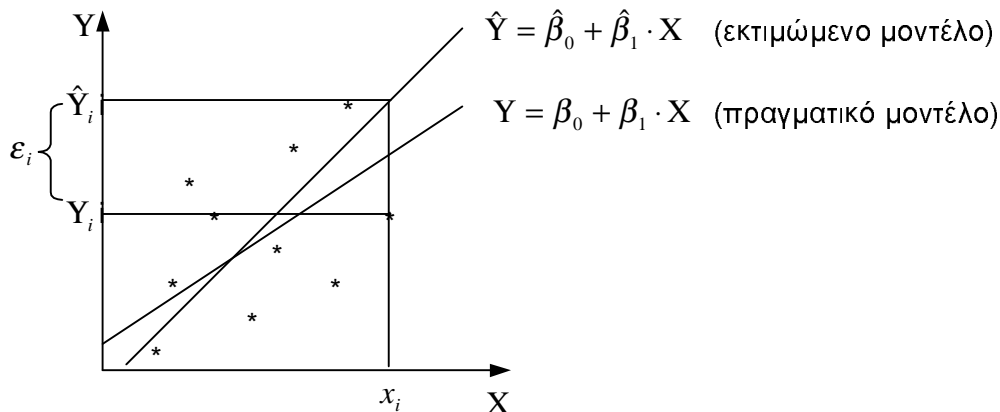
Είναι το μοντέλο με εξίσωση παλινδρόμησης  $E(Y) = \beta_0 + \beta_1 \cdot X$  (1), όπου  $\beta_0, \beta_1$  είναι άγνωστοι παράμετροι, οι οποίοι ονομάζονται συντελεστές παλινδρόμησης.

Το μοντέλο ονομάζεται απλό γιατί έχει μόνο μία ανεξάρτητη μεταβλητή  $X$  και γραμμικό ως προς τις παραμέτρους  $\beta_0, \beta_1$  και ως προς την ανεξάρτητη μεταβλητή  $X$ .

Δοθέντος του δείγματος  $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ , αναζητώ εκτιμήτριες των παραμέτρων  $\beta_0, \beta_1$ .

Ορίζουμε  $\varepsilon = Y - E(Y) = Y - \beta_0 + \beta_1 \cdot X \Rightarrow Y = \beta_0 + \beta_1 \cdot X + \varepsilon$  και από την (1) έχουμε ότι  $E(\varepsilon) = 0$ .

Από αυτό το μοντέλο, έχουμε για κάθε μέτρηση  $Y_i = \beta_0 + \beta_1 \cdot X_i + \varepsilon_i, i = 1, 2, \dots, n$ , με  $E(\varepsilon_i) = 0$ .



Εκτιμούμε τα  $\beta_0, \beta_1$  με τη μέθοδο ελαχίστων τετραγώνων, δηλαδή ελαχιστοποιούμε την ολική απόσταση των τετραγώνων των σφαλμάτων,  $\sum_{i=1}^n \varepsilon_i^2$ .

Συμβολικά:

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n \varepsilon_i^2 = \min_{\beta_0, \beta_1} \sum_{i=1}^n [Y_i - E(Y_i)]^2 = \min_{\beta_0, \beta_1} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 \cdot X_i)^2 = \min_{\beta_0, \beta_1} S(\beta_0, \beta_1).$$

Όμως,

$$\frac{\partial S}{\partial \beta_0} = 0 \Rightarrow -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 \cdot X_i) = 0 \Rightarrow \hat{\beta}_0 + \hat{\beta}_1 \cdot \bar{X} = \bar{Y} \quad \text{και}$$

$$\frac{\partial S}{\partial \beta_1} = 0 \Rightarrow -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 \cdot X_i) \cdot X_i = 0 \Rightarrow \hat{\beta}_0 \cdot \sum_{i=1}^n X_i + \hat{\beta}_1 \cdot \sum_{i=1}^n X_i^2 = \sum_{i=1}^n X_i \cdot Y_i.$$

Άρα,

$$\hat{\beta}_0 + \hat{\beta}_1 \cdot \bar{X} = \bar{Y} \quad \text{και}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X}) \cdot (Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

Σημείωση: Οι  $\frac{\partial S}{\partial \beta_0} = 0$  και  $\frac{\partial S}{\partial \beta_1} = 0$ , λέγονται πρώτη και δεύτερη κανονική εξίσωση αντίστοιχα.

Τώρα, από την (1) και απο τις πιο πάνω εξισώσεις, έχω:

$$\left. \begin{array}{l} E(Y) = \beta_0 + \beta_1 \cdot X \\ \hat{Y} = \hat{E}(Y) \end{array} \right\} \Rightarrow \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 \cdot X \quad (2)$$

Για  $X = X_i$ , η τιμή της  $\hat{Y}$  από την (2) συμβολίζεται με  $\hat{Y}_i$  και  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot X_i$  (εκτιμώμενη τιμή της  $Y$  για  $X = X_i$ ).

Παρατηρήσεις:

(1) Ένας άλλος τρόπος εκτίμησης των  $\beta_0, \beta_1$  είναι η ελαχιστοποίηση των απολύτων

τιμών των  $\varepsilon_i$ , δηλαδή,  $\min_{\beta_0, \beta_1} \sum_{i=1}^n |\varepsilon_i| = \min_{\beta_0, \beta_1} \sum_{i=1}^n |Y_i - \beta_0 - \beta_1 \cdot X_i|$ .

Η μέθοδος των ελαχίστων τετραγώνων πλεονεκτεί στο ότι παραγωγίζεται παντού, ενώ οι απόλυτες τιμές δεν παραγωγίζονται στο 0.

(2) Εναλλακτική μορφή του  $\hat{\beta}_1$ :

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X}) \cdot (Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\text{Αλλά, } \sum_{i=1}^n (X_i - \bar{X}) \cdot (Y_i - \bar{Y}) = \sum_{i=1}^n (X_i - \bar{X}) \cdot Y_i - \bar{Y} \cdot \sum_{i=1}^n (X_i - \bar{X}) = \sum_{i=1}^n (X_i - \bar{X}) \cdot Y_i .$$

$$\text{Άρα, } \hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X}) \cdot Y_i}{\sum_{i=1}^n (X_i - \bar{X})^2}, \text{ αλλά θέτοντας } \kappa_i = \frac{(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2}, i = 1, 2, \dots, n,$$

έχουμε ότι:  $\hat{\beta}_1 = \sum_{i=1}^n \kappa_i \cdot Y_i$ , δηλαδή βλέπουμε ότι το  $\hat{\beta}_1$  είναι γραμμικός συνδυασμός των παρατηρήσεων  $Y_i$ .

Ισχύει ότι:

$$\sum_{i=1}^n \kappa_i = \sum_{i=1}^n \left[ \frac{(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] = \frac{\sum_{i=1}^n (X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n X_i - n \cdot \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} = 0$$

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

(3) Εναλλακτική μορφή του  $\hat{\beta}_0$ :

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \cdot \bar{X} = \frac{1}{n} \cdot \sum_{i=1}^n Y_i - \left( \sum_{i=1}^n \kappa_i \cdot Y_i \right) \cdot \bar{X} = \sum_{i=1}^n \left( \frac{1}{n} - \kappa_i \cdot \bar{X} \right) \cdot Y_i ,$$

$$\text{όπου } \kappa_i = \frac{(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2}, i = 1, 2, \dots, n .$$

Παρατηρούμε ότι και το  $\hat{\beta}_0$  γράφεται σαν γραμμικός συνδυασμός των παρατηρήσεων  $Y_i$  και επιπλέον βλέπουμε ότι δεν εξαρτάται από το  $\hat{\beta}_1$  (όπως στην αρχική του μορφή).

(4) Εναλλακτική μορφή για το εκτιμώμενο μοντέλο:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 \cdot X = (\bar{Y} - \hat{\beta}_1 \cdot \bar{X}) + \hat{\beta}_1 \cdot X = \bar{Y} + \hat{\beta}_1 \cdot (X - \bar{X}), \text{ δηλαδή για } X = X_i \text{ θα έχω}$$

$$\hat{Y}_i = \bar{Y} + \hat{\beta}_1 \cdot (X_i - \bar{X}).$$

Άρα παρατηρούμε ότι δεν χρειάζεται να γνωρίζουμε το  $\hat{\beta}_0$  για να εκτιμήσουμε το  $\hat{Y}$ .

(5) Για την εκτίμηση των  $\hat{\beta}_0, \hat{\beta}_1$  δεν χρειαστήκαμε καμία πιθανοθεωρητική ιδιότητα για τα  $Y_i, i = 1, 2, \dots, n$ , δηλαδή, η κατανομή των  $Y_i$  δεν επηράζει καθόλου την εκτίμηση των  $\hat{\beta}_0, \hat{\beta}_1$ .

ΣΤΑΤΙΣΤΙΚΕΣ ΙΔΙΟΤΗΤΕΣ  
ΕΚΤΙΜΗΤΡΙΩΝ ΕΛΑΧΙΣΤΩΝ ΤΕΤΡΑΓΩΝΩΝ

Χρειαζόμαστε υποθέσεις για να μελετήσουμε τις ιδιότητες των  $\hat{\beta}_0, \hat{\beta}_1$ . Οι υποθέσεις μπορούν να εκφραστούν είτε μέσω των σφαλμάτων  $\varepsilon_i$  είτε μέσω των παρατηρήσεων  $Y_i, i = 1, 2, \dots, n$ .

A) Μέσω των σφαλμάτων  $\varepsilon_i, i = 1, 2, \dots, n$ :

(1) Υποθέτουμε ότι  $E(\varepsilon_i) = 0, i = 1, 2, \dots, n$ .

(2) Υποθέτουμε ότι  $Cov(\varepsilon_i, \varepsilon_j) = 0, i \neq j, i, j = 1, 2, \dots, n$  (ασυσχέτιστα σφάλματα) και  $Var(\varepsilon_i) = \sigma^2, i = 1, 2, \dots, n$  (σταθερή). ( $Var(\varepsilon_i) = \sigma^2$  είναι ισχυρή υπόθεση)

(3) Πιο ειδικά μπορώ να θεωρήσω ότι  $\varepsilon_i \sim N(0, \sigma^2), i = 1, 2, \dots, n$ , τυχαίο δείγμα (ανεξάρτητα).

B) Μέσω των παρατηρήσεων  $Y_i, i = 1, 2, \dots, n$ :

(1)'  $E(Y_i) = \beta_0 + \beta_1 \cdot X_i, i = 1, 2, \dots, n$  (αφού  $E(\varepsilon_i) = 0$ ).

(2)'  $Var(Y_i) = \sigma^2, i = 1, 2, \dots, n$  και  $Cov(Y_i, Y_j) = 0, i \neq j, i, j = 1, 2, \dots, n$ .

(3)'  $Y_i \sim N(\beta_0 + \beta_1 \cdot X_i, \sigma^2), i = 1, 2, \dots, n$  (ανεξάρτητες παρατηρήσεις, αλλά όχι ισόνομες).

Παρατήρηση:

Αν ισχύει η (3), ή ισοδύναμα η (3)', τότε ισχύουν και οι (1) και (2), ή ισοδύναμα οι (1)' και (2)'. Το αντίθετο δεν ισχύει.

Πρόταση 1:

Αν ισχύει η (1), ή ισοδύναμα η (1)', τότε  $\hat{\beta}_0$  και  $\hat{\beta}_1$  είναι αμερόληπτες εκτιμήτριες των  $\beta_0$  και  $\beta_1$  αντίστοιχα.

Απόδειξη:

Για το  $\hat{\beta}_1$  έχουμε:

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\sum_{i=1}^n \kappa_i \cdot Y_i\right) = \sum_{i=1}^n \kappa_i \cdot E(Y_i) = \sum_{i=1}^n \kappa_i \cdot (\beta_0 + \beta_1 \cdot X_i) = \beta_0 \cdot \sum_{i=1}^n \kappa_i + \beta_1 \cdot \sum_{i=1}^n \kappa_i \cdot X_i = \\ &= 0 + \beta_1 \cdot \frac{\sum_{i=1}^n (X_i - \bar{X}) \cdot X_i}{\sum_{i=1}^n (X_i - \bar{X})^2} = \beta_1 \cdot \frac{\sum_{i=1}^n (X_i - \bar{X}) \cdot (X_i + \bar{X} - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \\ &= \beta_1 \cdot \frac{\sum_{i=1}^n (X_i - \bar{X})^2 + \bar{X} \cdot \sum_{i=1}^n (X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \beta_1 \cdot \frac{\sum_{i=1}^n (X_i - \bar{X})^2 + 0}{\sum_{i=1}^n (X_i - \bar{X})^2} = \beta_1 \end{aligned}$$

Επίσης για το  $\hat{\beta}_0$  :

$$\begin{aligned} E(\hat{\beta}_0) &= E(\bar{Y} - \hat{\beta}_1 \cdot \bar{X}) = E(\bar{Y}) - \bar{X} \cdot E(\hat{\beta}_1) = \frac{1}{n} \cdot \sum_{i=1}^n E(Y_i) - \beta_1 \cdot \bar{X} = \\ &= \frac{1}{n} \cdot \sum_{i=1}^n E(\beta_0 + \beta_1 \cdot X_i) - \beta_1 \cdot \bar{X} = \beta_0 + \beta_1 \cdot \bar{X} - \beta_1 \cdot \bar{X} = \beta_0 \end{aligned}$$

Λήμμα:

Αν ισχύει η (2), ή ισοδύναμα η (2)', τότε  $Cov(\bar{Y}, \hat{\beta}_1) = 0$  και  $Var(\bar{Y}) = \frac{\sigma^2}{n}$ .

Απόδειξη:

$$\text{Γενικά έχω ότι } Cov\left(\sum_{i=1}^n \alpha_i \cdot Y_i, \sum_{j=1}^m \beta_j \cdot Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \cdot \beta_j \cdot Cov(Y_i, Y_j).$$

Άρα:

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

Είναι επίσης προφανές ότι:  $Var(\bar{Y}) = \frac{1}{n^2} \cdot \sum_{i=1}^n Var(Y_i) = \frac{1}{n^2} \cdot n \cdot \sigma^2 = \frac{\sigma^2}{n}$

Πρόταση 2:

Αν ισχύει η (2), ή ισοδύναμα η (2)', τότε  $Var(\hat{\beta}_0) = \sigma^2 \cdot \frac{\sum_{i=1}^n X_i^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2}$  και

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

Απόδειξη:

$$Var(\hat{\beta}_1) = Var\left(\sum_{i=1}^n \kappa_i \cdot Y_i\right) = \sum_{i=1}^n \kappa_i^2 \cdot Var(Y_i) = \sigma^2 \cdot \sum_{i=1}^n \kappa_i^2 = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\begin{aligned}
\text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{Y} - \hat{\beta}_1 \cdot \bar{X}) = \text{Var}(\bar{Y}) + \bar{X}^2 \cdot \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{n} + \bar{X}^2 \cdot \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = \\
&= \sigma^2 \cdot \frac{\sum_{i=1}^n (X_i - \bar{X})^2 + n \cdot \bar{X}^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} = \sigma^2 \cdot \frac{\sum_{i=1}^n X_i^2 - n \cdot \bar{X}^2 + n \cdot \bar{X}^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} = \sigma^2 \cdot \frac{\sum_{i=1}^n X_i^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2}
\end{aligned}$$

Πρόταση 3:

Αν ισχύει η (3), ή ισοδύναμα η (3)', τότε:

$$\hat{\beta}_0 \sim N \left( \beta_0, \sigma^2 \cdot \frac{\sum_{i=1}^n X_i^2}{n \cdot \sum_{i=1}^n (X_i - \bar{X})^2} \right)$$

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

Απόδειξη:

Προφανές, αφού τα  $\hat{\beta}_0$  και  $\hat{\beta}_1$  είναι γραμμικοί συνδυασμοί των παρατηρήσεων, οι οποίες ακολουθούν την κανονική κατανομή και ισχύουν συγχρόνως οι (1) και (2).