

Προβλήματα

Φροντιστήριο #2

Απλή γραμμική παλινδρόμηση

Ενδιαφερόμαστε να μελετήσουμε τη σχέση μεταξύ δύο μεταβλητών X και Y .

- Αν για κάθε τιμή x της X καθορίζεται μονοσήμαντα η τιμή y της Y τότε ορίζεται μία συνάρτηση $y = f(x)$, όπου $f(x)$ είναι ο μαθηματικός τύπος υπολογισμού της y μέσω της x .

- Για κάθε τιμή x της X η ποσότητα Y μπορεί να είναι μία $z.p.$

Στην πρώτη περίπτωση λέμε ότι υπάρχει μία συναρτησιακή ή ντετερμινιστική (deterministic, functional) σχέση μεταξύ των μεταβλητών X και Y . Στη δεύτερη περίπτωση έχουμε μία στοχαστική ή στατιστική ή μη ντετερμινιστική (stochastic, statistic, non-deterministic) σχέση.

Για τις σχέσεις της δεύτερης κοφής, η μεταβλητή X θα καλείται ανεξάρτητη ή προβλέπουσα ή ελεγχόμενη ενώ η Y θα καλείται εξαρτημένη ή μεταβλητή απόκρισης.

Συνήθως οι ανεξάρτητες μεταβλητές είναι αυτές στις οποίες μπορούμε να δώσουμε μία τιμή (π.χ. τιμή πώλησης ενός προϊόντος, ποσότητα λιπάσματος για έναν αγρό κλπ) και οι εξαρτημένες μεταβλητές (π.χ. αριθμός πωλήσεων παραγωγή του αγρού) εμπνέονται από μεταβολές των ανεξάρτητων μεταβλητών.

Το πιο απλό μοντέλο ανάγωγης παλινδρόμησης είναι -2-
το γραμμικό (στοχαστικό) μοντέλο για το οποίο η μορφή
εξάρτησης των μεταβλητών X και Y περιγράφεται από
τη σχέση:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

όπου X είναι η τιμή που πήρε η μη-τυχαία μεταβλητή
 X , Y είναι η τιμή που περιγράφει τη μεταβλητή από-
κρισης για τη συγκεκριμένη τιμή X της X και ε
είναι μια τ.ρ. που εκφράζει την απόκλιση της Y από
το γραμμικό όρο $\beta_0 + \beta_1 X$.

Έστω x_1, x_2, \dots, x_n οι τιμές που δίνουμε στην X . Καταγρά-
φουμε τις αντίστοιχες τιμές y_1, y_2, \dots, y_n που λαμβάνει η
μεταβλητή Y . Για τα ζεύγη $(x_i, y_i), i=1, \dots, n$ μπορούμε
να γράψουμε:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, i=1, \dots, n$$

Αν παραστήσουμε τα ζεύγη (x_i, y_i) σε ένα σύστημα
ορθογωνίων αξόνων θα πάρουμε το διάγραμμα διασποράς.
Στο διάγραμμα αυτό, η εξίσωση $y = \beta_0 + \beta_1 x$ θα παρι-
στάνει μία ευθεία που προσεγγίζει τα δεδομένα μας.
Η παράμετρος β_0 δίνει τη θέση όπου η ευθεία τέμνει
τον άξονα y/y και η β_1 δίνει το συντελεστή διεύδη-
σης της ευθείας.

Ερώτηση: Με ποιο τρόπο μπορούμε να καθορίσουμε
τις άγνωστες παραμέτρους β_0 και β_1 ώστε η ευθεία
που θα προκύψει να μας δίνει όσο το δυνατόν μαζύ-
τερη περιγραφή της σχέσης εξάρτησης των X
και Y ;

Ο καθορισμός των β_0, β_1 λέγεται επιτήρηση των παραμέτρων ενώ οι τιμές που προκύπτουν για αυτές με την υλοποίηση της διαδικασίας επιτήρησης λέγονται επιτηρήσεις. Η ευθεία που προσεγγίζει καλύτερα τα σημεία καλείται ευθεία παλινδρόμησης της Y στη X .

Για το απλό γραμμικό μοντέλο παλινδρόμησης ισχύουν οι επόμενες υποθέσεις: $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, i = 1, \dots, n$

1. Οι ποσότητες β_0, β_1 είναι οι άγνωστες παράμετροι
2. Τα x_i είναι γνωστοί αριθμοί που συνήθως καθορίζονται από τον ερευνητή. Είναι οι τιμές της ανεξάρτητης μεταβλητής X .

3. Τα Y_i είναι οι τιμές της εξαρτημένης μεταβλητής Y . Το Y_i είναι τιμ. και με y_i συμβολίζουμε την τιμή που λαμβάνει αυτή όταν $X = x_i$. Δηλ. είναι η παρατηρούμενη τιμή της Y όταν $X = x_i$.

4. Τα $\epsilon_i, i = 1, \dots, n$ είναι τυχαία σφάλματα με μέση τιμή 0 και διασπορά σ^2 , δηλ. $E(\epsilon_i) = 0$
 $Var(\epsilon_i) = \sigma^2$.

5. Τα σφάλματα ϵ_i, ϵ_j που αντιστοιχούν σε διαφορετικές μετρήσεις $i, j, i \neq j$ θεωρούνται ασυσχίτητα δηλ. $cov(\epsilon_i, \epsilon_j) = 0 \forall i \neq j$.

δηλ. $E(Y_i) = \beta_0 + \beta_1 x_i, Var(Y_i) = \sigma^2$

και $cov(Y_i, Y_j) = 0$

Μπορούμε να γράψουμε:

$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, Y = \beta_0 + \beta_1 x + \epsilon$ ή

$$E(Y_i) = \beta_0 + \beta_1 x_i \quad \text{ή} \quad E(Y) = \beta_0 + \beta_1 x$$

-4-

Δε γίνεται καμία υπόθεση για το είδος της κατανομής των Y_i ή των ε_i π.χ. αν είναι κανονική ή κάποια άλλη κατανομή πέραν του ότι αυτή έχει τη ίδια διασπορά σ^2 για όλες τις τιμές x_i της ανεξάρτητης μεταβλητής X .

Το πλέον συνήθισμένο κριτήριο για τη ευρίσκηση των παραμέτρων β_0, β_1 είναι η ελαχιστοποίηση του αθροίσματος των τετραγώνων:

$$g(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

όπου $\varepsilon_i = y_i - (\beta_0 + \beta_1 x_i)$, $i=1, 2, \dots, n$ είναι οι αποκλίσεις ή τα σφάλματα (errors, deviations).

Το άθροισμα $g(\beta_0, \beta_1)$ λέγεται άθροισμα τετραγώνων των σφαλμάτων (sum of squares of errors), ενώ οι ποσότητες που προκύπτουν για τις παραμέτρους β_0, β_1 ελαχιστοποιώντας το άθροισμα αυτό καλούνται ευρισκόμενες ελάχιστες τετραγώνων, $\hat{\beta}_0, \hat{\beta}_1$ των β_0, β_1 , αντίστοιχα. Η συνηνευμένη διαδικασία ευρίσκησης των παραμέτρων β_0, β_1 είναι γνωστή ως μέθοδος των ελάχιστων τετραγώνων (Least Squares Method).

Ελαχιστοποιούμε την $g(\beta_0, \beta_1)$ ως προς β_0, β_1 .

$$\frac{\partial g}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial g}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\Rightarrow \begin{cases} \sum_{i=1}^n y_i = n \beta_0 + \beta_1 \sum_{i=1}^n x_i & (1) \\ \sum_{i=1}^n y_i x_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \end{cases}$$

Με μέθοδο ορισμών

$$\hat{\beta}_1 = \frac{\begin{vmatrix} \sum_{i=1}^n y_i & \sum_{i=1}^n y_i x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{vmatrix}}{\begin{vmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{vmatrix}} = \frac{n \sum_{i=1}^n y_i x_i - \left(\sum_{i=1}^n y_i \right)^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \quad (2)$$

$$\text{Από (1)} \Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (3)$$

Άρα $\hat{\beta}_0, \hat{\beta}_1$ δίνονται από (3) & (2).

Παράδειγμα (1) Έστω $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ η ευθεία ελαχίστων τετραγώνων που βρέθηκε με βάση n ζεύγη σημείων $(x_i, y_i), i=1, \dots, n$. Έστω \hat{y}_i οι επιτηρημένες τιμές

$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, i=1, 2, \dots, n$. Να δείχθει ότι ο μέσος όρος των $\hat{y}_i, i=1, \dots, n$ είναι ίσος με \bar{y} .

Λύση Επειδή $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ έχουμε

$$\hat{y}_i = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1 (x_i - \bar{x}), i=1, \dots, n$$

Αθροίζοντας στα αόθε $i=1, 2, \dots, n$ έχουμε:

$$\sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n (\bar{y} + \hat{\beta}_1 (x_i - \bar{x})) = \sum_{i=1}^n \bar{y} + \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x}) \quad -6-$$

Οπως $\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} = \sum_{i=1}^n x_i - n\bar{x}$
 $= n\bar{x} - n\bar{x} = 0$

Αρα $\sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n \bar{y} = n\bar{y} = \sum_{i=1}^n y_i$

Αν διαφύσουμε τη τελευταία σχέση με n , έχουμε:

$$\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n y_i \Leftrightarrow \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y}$$

δηλαδή ο μέσος όρος των $\hat{y}_i, i=1, 2, \dots, n$ είναι ίσος με \bar{y} .

Παρατήρηση: Αν θέσουμε $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}, \bar{y} = \frac{\sum_{i=1}^n y_i}{n}$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2$$

μπορούμε να γράψουμε:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \bar{y} - \frac{S_{xy}}{S_{xx}} \bar{x}$$

Επίσης, μπορούμε να γράψουμε:

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x}) y_i \text{ και για τη ζ.ρ. } Y$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x}) Y_i \text{ 'Αρα } \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sum_{i=1}^n d_i Y_i$$

$$\text{όπου } c_i = \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{x_i - \bar{x}}{S_{xx}}$$

Με αυτό τον μετασχηματισμό για το $\hat{\beta}_1$ μπορούμε να δείξουμε ότι $E(\hat{\beta}_1) = \beta_1$ δηλαδή ότι $\hat{\beta}_1$ είναι αμερόληπτη εκτίμηση της παραμέτρου β_1 διότι,

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\sum_{i=1}^n c_i Y_i\right) = \sum_{i=1}^n c_i E(Y_i) \\ &= \sum_{i=1}^n c_i (\beta_0 + \beta_1 x_i) = \beta_0 \sum_{i=1}^n c_i + \beta_1 \sum_{i=1}^n c_i x_i \end{aligned}$$

$$\text{Όπως } \sum_{i=1}^n c_i = \frac{\sum_{i=1}^n (x_i - \bar{x})}{S_{xx}} = \frac{n\bar{x} - n\bar{x}}{S_{xx}} = 0$$

$$\begin{aligned} \text{Επίσης} \sum_{i=1}^n c_i x_i &= \frac{\sum_{i=1}^n (x_i - \bar{x}) x_i}{S_{xx}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x} + \bar{x})}{S_{xx}} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{S_{xx}} + \bar{x} \frac{\sum_{i=1}^n (x_i - \bar{x})}{S_{xx}} = \frac{S_{xx}}{S_{xx}} + 0 = 1 \end{aligned}$$

Άρα $E(\hat{\beta}_1) = \beta_0 \cdot 0 + \beta_1 \cdot 1 = \beta_1$

Ακόμη, $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$ και $\hat{\beta}_1 = \sum_{i=1}^n c_i Y_i$

$$\text{Άρα } \hat{\beta}_0 = \bar{Y} - \left(\sum_{i=1}^n c_i Y_i\right) \bar{x} = \sum_{i=1}^n \left(\frac{1}{n}\right) Y_i - \sum_{i=1}^n c_i \bar{x} Y_i$$

$$= \sum_{i=1}^n \left(\frac{1}{n} - c_i \bar{x}\right) Y_i = \sum_{i=1}^n d_i Y_i \text{ όπου}$$

$$d_i = \left(\frac{1}{n} - c_i \bar{x}\right) = \frac{1}{n} - \frac{(x_i - \bar{x}) \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Συνεπώς,

$$E(\hat{\beta}_0) = E(\bar{Y} - \hat{\beta}_1 \bar{x}) = E(\bar{Y}) - \bar{x} E(\hat{\beta}_1) = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0$$

Άρα, $\hat{\beta}_0$ αμερόληπτος εκτιμητής της παραμέτρου β_0 . \square

Εφαρμογή Ένα υγιές συσκευάζεται από ένα εργοστάσιο σε μεγάλα κουτιά ο αριθμός των οποίων ποικίλει ανάλογα με την παραγωγή. Ο πίνακας δίνει τον αριθμό των κουτιών που συσκευάστηκαν ώστε να καλυφθούν οι παραγγελίες του εργοστασίου και οι εργατοώρες που χρειάστηκαν για 10 πρόσφατες παραγγελίες που εντελέστηκαν.

(α) Ποια από τις δύο μεταβλητές (αριθμός κουτιών, εργατοώρες) μπορεί να θεωρηθεί ως ανεξάρτητη μεταβλητή X και ποια ως εξαρτημένη Y ;

(β) Να σχεδιασθεί το αντίστοιχο διάγραμμα διασποράς των δεδομένων.

(γ) Να υπολογισθεί με τη μέθοδο ελαχίστων τετραγώνων η ευθεία γραμμικής παλινδρόμησης της Y πάνω στη X και να χαραχθεί στο αντίστοιχο διάγραμμα διασποράς.

(δ) Να δοθεί ερμηνεία της κλίσης $\hat{\beta}_1$ και του σταθερού όρου $\hat{\beta}_0$ της ευθείας παλινδρόμησης.

#κουτιών	60	40	120	160	80	100	100	140	180	70
Εργατοώρες	230	161	365	515	263	335	335	464	587	245

Λύση (α) Ως ανεξάρτητη μεταβλητή (X) θα θεωρήσουμε τον αριθμό των κουτιών ενώ ως εξαρτημένη (Y) τις εργατοώρες που χρειάζονται για τη συσκευασία τους. Προφανώς, στόχος είναι να μπορέσουμε να κατασκευάσουμε ένα μοντέλο με το οποίο θα μπορούμε να προβλέψουμε

τις εργαζώμενες που θα χρειαστούν όταν δίνεται ο αριθμός των κουτιών που πρέπει να συσκευαστούν.

(β) Από το διάγραμμα διασποράς (βλέπε Σχήμα) μπορούμε να πούμε ότι η σχέση των μεταβλητών X και Y είναι κατά προσέγγιση γραμμική.

(δ) Από τα δεδομένα μας, για $X := \#$ κουτιών και $Y :=$ εργαζώμενες, βρίσκουμε:

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} = \frac{10 \cdot 422250 - 1050 \cdot 3500}{10 \cdot 128500 - 1050^2} = 3$$

$$\hat{\beta}_0 = \frac{1}{n} \sum_{i=1}^n y_i - \hat{\beta}_1 \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{10} 3500 - 3 \frac{1}{10} 1050 = 350 - 315 = 35.$$

Άρα, η ευθεία παλινδρόμησης της Y πάνω στη X είναι

$$y = 35 + 3x \quad \text{ή} \quad \hat{Y} = 35 + 3x.$$

Το σχήμα με τη ευθεία παλινδρόμησης μαζί με το διάγραμμα διασποράς δίνεται παρακάτω.

(ε) Ερμηνεία $\hat{\beta}_1$: Για κάθε επιπλέον κουτί σε μια παραγγελία που φτάνει στο εργοστάσιο χρειάζονται 3 επιπλέον εργαζώμενες για τη συσκευασία του.

Ερμηνεία $\hat{\beta}_0$: Χρειάζονται επίσης 35 εργαζώμενες για προπαρασκευαστικές δουλειές πριν αρχίσει η συσκευασία των κουτιών που απαιτούνται για την κάλυψη των παραγγελιών. Αυτό προκύπτει από το γεγονός ότι για παραγγελία $x=0$ κουτιών η ευθεία παλινδρόμησης δίνει $\hat{Y} = \hat{\beta}_0 = 35$.

① Επισημαίνουμε ότι ο όρος "γραμμικό" για το χαρακτηρισμό του μοντέλου αναφέρεται στις παραμέτρους και όχι στις μεταβλητές.

Για παράδειγμα, το μοντέλο $Y = \beta_0 + \beta_1 \frac{1}{X} + \varepsilon$ είναι επίσης γραμμικό και μπορούμε να υπολογίσουμε τους ε.ε.τ.

των β_0, β_1 :
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \frac{y_i}{x_i} - \sum_{i=1}^n y_i \sum_{i=1}^n \frac{1}{x_i}}{\sum_{i=1}^n \frac{1}{x_i^2} - \left(\sum_{i=1}^n \frac{1}{x_i}\right)^2}$$

και
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \sum_{i=1}^n \frac{1}{x_i}$$

Επίσης, το μοντέλο $Y = e^{\beta_1 X} + \varepsilon$ είναι επίσης γραμμικό ως προς την παράμετρο β_1 γιατί μπορεί να γραφεί:

$\ln(Y) = \beta_1 X + \varepsilon$ και μπορούμε και για αυτό το μοντέλο να υπολογισθεί η ε.ε.τ της β_1 .

② Στην απλή γραμμική παλινδρόμηση της Y στην X η ευθεία ελαχίστων τετραγώνων περνά από το σημείο (\bar{x}, \bar{y}) διότι $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ και άρα

$$\hat{\beta}_0 + \hat{\beta}_1 \bar{x} = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 \bar{x} = \bar{y}$$

③ Ισχύει ότι $\sum_{i=1}^n \hat{\varepsilon}_i = 0$ διότι:

$$\sum_{i=1}^n \hat{\varepsilon}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)$$

$$= \sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\beta}_0 - \sum_{i=1}^n \hat{\beta}_1 x_i$$

$$= \sum_{i=1}^n y_i - n \hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i - n(\bar{y} - \hat{\beta}_1 \bar{x}) - \hat{\beta}_1 n \bar{x} = \sum_{i=1}^n y_i - n \bar{y} = n \bar{y} - n \bar{y} = 0$$