

9ο online μάθημα

23/4/2020

2ος τρόπος: Γραμμικός Συντελεστής Συσχέτισης

Έστω X, Y δύο συνεχείς τ.κ. Εισαγεται τότε η έννοια της
covariances των X, Y η οποία συμβολίζεται $\text{Cov}(X, Y)$
και ορίζεται:

↓
Covariance

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

$$\text{όπου } \mu_X = EX, \mu_Y = EY$$

εύκολα προκύπτει

$$\text{Cov}(X, Y) = E(XY) - (EX)(EY)$$

Παρατήρηση 1η: Η covariances είναι ένα μέτρο της συσχέτισης των X, Y . Έχει ως μονάδες μέτρησης το γινόμενο των μονάδων μέτρησης των X, Y . Με στόχο να απαλλαγούμε από τις μονάδες οδηγούμαστε στο συντελεστή συσχέτισης.

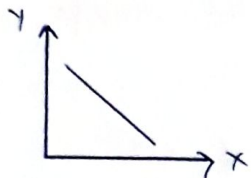
$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var } X} \sqrt{\text{Var } Y}}$$

Θεωρητικός - Πληθυσιακός γραμμικός συντελεστής συσχέτισης ³

- είναι καθαρός αριθμός
- Αποδοικνύεται (μάθημα 531) ότι

α) $-1 \leq P(x, y) \leq 1$

β) $P(x, y) = -1$ έχουμε τέλεια αρνητική γραμμική συσχέτιση



γ) $P(x, y) = 1$ έχουμε τέλεια θετική γραμμική συσχέτιση

δ) $P(x, y) = 0$ τότε δεν υπάρχει γραμμική σχέση

Επομένως για να δούμε αν υπάρχει ή όχι γραμμική σχέση μεταξύ των X, Y αρκεί να ελέγξω: Την:

$H_0: P(x, y) = 0$ • Αν αυτός ο έλεγχος απορρίπτεται ευραίων ότι υπάρχει γραμμική σχέση.

$H_{0a}: P(x, y) \neq 0$ • Αν δεν απορρίπτεται ευραίων ότι δεν υπάρχει γραμμική.

Ο έλεγχος "πρέπει" να γίνει από ένα στατιστικό ανάλογο του $P(x, y)$. Είναι:

$$r = r_{x,y} = r(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

ή ισοδύναμα

$$r = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)}{\sqrt{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \sqrt{n \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i \right)^2}}$$

- Παρατηρήσεις:
- $r(x, y) = r(y, x)$
 - Ορίζεται για ποσοτικές μεταβλητές
 - $r(x, y) > 0$ υποδηλώνει θετική συσχέτιση.
Ενώ $r(x, y) < 0$ αρνητική.
 - Τιμές του $r(x, y)$ "κοντά" στο 0 δηλώνουν μη ύπαρξη γραμμικής σχέσης.
 - ΠΡΟΒΟΛΗ ΕΝΔΕΙΚΤΙΚΩΝ ΓΡΑΦΗΜΑΤΩΝ
 - Η μη ύπαρξη γραμμικής σχέσης δεν συμβαίνει των μη ύπαρξη σχέσης.

Έλεγχος για τον $\rho(x, y)$

$H_0: \rho(x, y) = 0$

έναντι μιας εκ των θετικής γραμμική συσχέτιση ή αρνητική γραμμική συσχέτιση ή $H_1: \rho(x, y) < 0$ ή $H_1: \rho(x, y) \neq 0$ γραμμική συσχέτιση

Υπό την προϋπόθεση ότι (x, y) έχουν διδιάστατη κανονική και δεν υπάρχουν ακέραιες τιμές.

$$t = \frac{r \sqrt{n-2}}{\sqrt{1-r^2}}$$

$H_0 \sim t_{n-2}$ \rightarrow Στηχαστικός συντελεστής συσχέτισης

πλήθος (x_i, y_i)

Άρα περιοχή απόρριψης

$t > t_{n-2, \alpha}$ ή $t < -t_{n-2, \alpha}$ ή $|t| \geq t_{n-2, \alpha/2}$ αντιστοίχα.
 $H_0: \rho > 0$ $H_0: \rho < 0$ $H_0: \rho \neq 0$

Παρατήρηση: Για τον γενικότερο έλεγχο
 $H_0: \rho(x, y) = \rho_0$ ($\rho_0 \in [-1, 1]$), $H_1: \rho > \rho_0$ ή $H_1: \rho < \rho_0$ ή $H_1: \rho \neq \rho_0$

Χρησιμοποιείται το στατιστικό

$$W = \frac{Z - E(Z)}{\sqrt{\text{Var } Z}} \stackrel{H_0}{\sim} N(0,1)$$

όπου

$$Z = \frac{1}{2} \ln \frac{1+r}{1-r}$$

$$E Z = \frac{1}{2} \ln \frac{1+r_0}{1-r_0}$$

$$\text{Var } Z = \frac{1}{n-3} \quad]$$

Εκτός

Μόνο για εγκυκλοπαιδικούς λόγους

Περιοχές απόρριψης:

$$W \geq Z_\alpha \quad \eta \quad W \leq -Z_\alpha \quad \eta \quad |W| \geq Z_{\alpha/2}$$

Αφού διαπιστωθεί η ύπαρξη γραμμικής σχέσης, στόχος είναι ο προσδιορισμός της με απώτερο σκοπό την αξιολόγηση της για τη πρόβλεψη της μιας μεταβλητής όταν η άλλη θεωρείται δεδομένη. Η γραμμική παλινδρόμησης της Y πάνω στην X εκφράζεται

από τη σχέση:

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

με β_0, β_1 σταθερές αλλά άγνωστες ποσότητες και ϵ το τυχαίο σφάλμα που περιβάλλει την απόκλιση της Y από την γραμμική παλινδρόμησης στο βιβλίο X .

Η σχέση προφανώς θα ικανοποιείται από τα διαθέσιμα δεδομένα επί. 1ε

Επομένως $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad i=1, \dots, n$

1^ο Βήμα: Εκτίμηση των παραμέτρων β_0, β_1

Πώς; Προσδιορισμός των β_0, β_1 έτσι ώστε να ελαχιστοποιείται

(6)

Το άθροισμα των τετραγώνων των εφαλμάτων. Η μέθοδος αυτή εκτιμήσεις είναι χρωστή. ως μέθοδος ελάχιστων τετραγώνων θέλουμε να ελαχιστοποιήσουμε την

$$Q(\theta_0, \theta_1) = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2 \quad \text{ως προς } \theta_0, \theta_1$$

Είναι

$$\frac{\partial Q}{\partial \theta_0} = -2 \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)$$

$$\frac{\partial Q}{\partial \theta_1} = -2 \sum_{i=1}^n x_i (y_i - \theta_0 - \theta_1 x_i)$$

Οι εκτιμήσεις, έστω $\hat{\theta}_0, \hat{\theta}_1$ προκύπτουν από την επίλυση του συστήματος

$$\begin{cases} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i) = 0 \\ \sum_{i=1}^n x_i (y_i - \theta_0 - \theta_1 x_i) = 0 \end{cases}$$

Γνωστό ως ΣΥΣΤΗΜΑ
ΚΑΝΟΝΙΚΩΝ ΕΞΙΣΩΣΕΩΝ

ή ισοδύναμα

$$\begin{cases} \sum_{i=1}^n y_i = n \hat{\theta}_0 + \hat{\theta}_1 \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i y_i = \hat{\theta}_0 \sum_{i=1}^n x_i + \hat{\theta}_1 \sum_{i=1}^n x_i^2 \end{cases}$$

Από την πρώτη προκύπτει ότι

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

ενώ έπειτα έχουμε

$$\sum_{i=1}^n x_i y_i = \bar{y} \sum_{i=1}^n x_i - \hat{\theta}_1 \bar{x} \sum_{i=1}^n x_i + \hat{\theta}_1 \sum_{i=1}^n x_i^2 \Rightarrow$$

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

$$\Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n}}{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i \right)^2}{n}}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Παρατήρηση: Κανονικά θα έπρεπε για να είμαστε βέβαιοι ότι είναι μέγιστο ελάχιστου, να ελέγξουμε τις δεύτερες παραχώφους και εβανό πίνακα.

Εκτιμήσεων τιμή \hat{y} :

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

Ερμηνεία $\hat{\beta}_0, \hat{\beta}_1$; ; ;

$\hat{\beta}_0$: τιμή της \hat{y} για $x=0$

$\hat{\beta}_1$: μεταβολή της \hat{y} για μοναδιαία μεταβολή της x

Γιατί;

Έστω (x_1, y_1) τότε

$$\hat{y}_1 = \hat{\beta}_0 + \hat{\beta}_1 x_1$$

Εννι για $x_2 = x_1 + 1$ είναι

$$\hat{y}_2 = \hat{\beta}_0 + \hat{\beta}_1 (x_1 + 1) \Rightarrow \hat{y}_2 = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_1 \Rightarrow \hat{y}_2 = \hat{y}_1 + \hat{\beta}_1$$

Υπόλοιπα: $e_i = y_i - \hat{y}_i \quad i=1, \dots, n$

Παρατηρήσεις: ① Το σημείο (\bar{x}, \bar{y}) ανήκει στην ευθεία παλινδρόμησης.
 Όπως για $x = \bar{x}$ είναι

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

② Τα υπόλοιπα αθροίζονται στο 0.

Όπως

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

Η πιο απλά ικνύει από το σύστημα κανονικών εξισώσεων ότι

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

③ Ικνύει ότι $\sum_{i=1}^n x_i e_i = 0$. Όπως είναι η 2^η σχέση του συστήματος κανονικών εξισώσεων.

④ Ικνύει ότι $\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i$

Όπως

$$\sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n \hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i \Rightarrow \sum_{i=1}^n y_i = n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i$$

Όπως από την 1^η από κανονικές εξισώσεις:

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad \text{Σημείωση}$$

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

Ιδιότητες των εκτιμητών ελαχίστων τετραγώνων

Έχουμε δείξει ότι

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\sum (x_i - \bar{x})(y_i - \bar{y}) = \sum (x_i - \bar{x})y_i - \sum (x_i - \bar{x})\bar{y}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})y_i - \sum_{i=1}^n (x_i - \bar{x})\bar{y}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

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

$$\sum (x_i - \bar{x}) = \sum x_i - n\bar{x} = \sum x_i - n \frac{\sum x_i}{n}$$

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

O.S.O. $E\hat{\beta}_0 = \beta_0$

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

υπό την υπόθεση ότι $E(\epsilon_i) = 0$

είναι

$$E\hat{\beta}_1 = E \left[\frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

$$= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n E[(x_i - \bar{x})y_i] = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x})E y_i$$

όπως $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

άρα $E y_i = \beta_0 + \beta_1 x_i$

Επιδείνωμα

$$E \hat{\beta}_1 = \frac{1}{\sum (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) (\beta_0 + \beta_1 x_i)$$

$$= \frac{1}{\sum (x_i - \bar{x})^2} \left[\beta_0 \sum_{i=1}^n (x_i - \bar{x}) + \beta_1 \sum_{i=1}^n (x_i - \bar{x}) x_i \right]$$

$$= \frac{1}{\sum (x_i - \bar{x})^2} \beta_1 \sum (x_i - \bar{x})^2 = \beta_1$$

Χρησιμοποιήστε ότι

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

Τώρα είναι:

$$E \hat{\beta}_0 = E \bar{y} - \bar{x} E \hat{\beta}_1 \Rightarrow E \hat{\beta}_0 = E \frac{\sum y_i}{n} - \bar{x} \beta_1 \Rightarrow$$

$$E \hat{\beta}_0 = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i) - \bar{x} \beta_1 \Rightarrow E \hat{\beta}_0 = \beta_0$$

Πρόταση: Υπό την υπόθεση ότι $E(\varepsilon_i) = 0$, $\text{Var}(\varepsilon_i) = \sigma^2$, v. s. o.

$$\text{Var} \hat{\beta}_1 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Var} \hat{\beta}_0 = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)$$

$$\text{Var} \hat{\beta}_1 = \text{Var} \left[\frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} \right] \Rightarrow$$

$$\text{Var } \hat{\theta}_1 = \left\{ \frac{1}{\sum (x_i - \bar{x})^2} \right\}^2 \text{Var} \left(\sum_{i=1}^n (x_i - \bar{x}) y_i \right)$$

$$\Rightarrow \text{Var } \hat{\theta}_1 = \left\{ \frac{1}{\sum (x_i - \bar{x})^2} \right\}^2 \sum_{i=1}^n (x_i - \bar{x})^2 \text{Var } y_i$$

$$\rightarrow \text{Var } \hat{\theta}_1 = \left\{ \frac{1}{\sum (x_i - \bar{x})^2} \right\}^2 \sum_{i=1}^n (x_i - \bar{x})^2 \sigma^2$$

και αποδειχθηκε το \int του βρω.

$$\text{Var } \hat{\theta}_0 = \text{Var} (\bar{y} - \hat{\theta}_1 \bar{x})$$

Προσχημα:

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

αρα δεν ειναι ανεξαρτητες τ.κ. για να μπορεσω να πω οτι

π.χ. $\text{Var} (a_1 x_1 + a_2 x_2) = a_1^2 \text{Var } x_1 + a_2^2 \text{Var } x_2$

Ειναι οπως

$$\begin{aligned} \bar{y} - \hat{\theta}_1 \bar{x} &= \sum_{i=1}^n \frac{1}{n} y_i - \frac{\bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) y_i \\ &= \sum_{i=1}^n \left\{ \frac{1}{n} - \frac{\bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} (x_i - \bar{x}) \right\} y_i \end{aligned}$$

Εποκενως:

$$\text{Var } \hat{\theta}_0 = \sum_{i=1}^n \left\{ \frac{1}{n} - \frac{\bar{x} (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right\}^2 \text{Var } y_i$$

$$= 6^2 \left\{ \sum_{i=1}^n \frac{1}{n^2} - 2 \frac{\bar{x} \sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} + \frac{\bar{x}^2 \sum (x_i - \bar{x})^2}{\left\{ \sum (x_i - \bar{x})^2 \right\}^2} \right\}$$

$$= 6^2 \left\{ n \frac{1}{n^2} - 2 \frac{\bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) + \frac{\bar{x}^2}{\left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} \sum (x_i - \bar{x})^2 \right\}$$

$$= 6^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right\}$$

Παρατήρηση: Πιο όμορφες αποδείξεις στο βιβλίο του Γου Εζακίου!!!