

Διάλεξη 2^ηΠαλινδρόμηση και Ανάλυση Διακύμανσης.Σεμινάρια
ΣτατιστικήςΠαρασκευή 18/10/2019
ώρα 12:00 Αιθ. 307α

Απλή Γραμμική Παλινδρόμηση (αχπ)

Τμ.	X	X ₁ , X ₂ , ..., X _n
	Y	Y ₁ , Y ₂ , ..., Y _n

αχπ δηλ

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \quad i = 1, 2, \dots, n$$

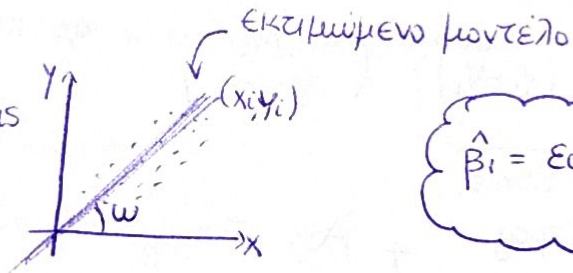
εφαρτημένοι κτρ
σφάλματα
↑
παραμέτροι ανεξάρτητη

ΕΕΤ των β₀, β₁.

$$\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Εκτιμώμενο Μοντέλο: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$

Γραφικά: Διαγράμμα Διασποράς



$$\hat{\beta}_1 = \epsilon\phi\omega$$

Ερμηνεία των $\hat{\beta}_0, \hat{\beta}_1$ Το $\hat{\beta}_0$ εκφράζει την τιμή της Y όταν X=0Το $\hat{\beta}_1$ εκφράζει την μεταβολή της Y σε μοναδιαία μεταβολή της X.Υπολοιπα: Εκφράζουν την απόκλιση του μοντέλου από την πραγματικότητα.

$$e_i = Y_i - \hat{Y}_i$$

↑
πραγματικότητα μοντέλο

Ιδιότητα Υπολοίπων: $\sum_{i=1}^n e_i = 0$

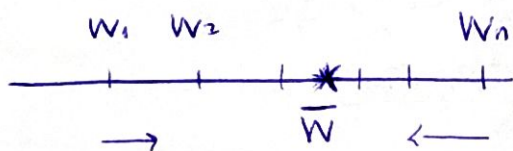
$$\begin{aligned} \text{Αποδ: } \sum_{i=1}^n e_i &= \sum_{i=1}^n (Y_i - \hat{Y}_i) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i) = \sum_{i=1}^n (Y_i - \hat{Y} + \hat{\beta}_1 \bar{X} - \hat{\beta}_1 X_i) \\ &= \sum_{i=1}^n (Y_i - \bar{Y}) - \hat{\beta}_1 \sum_{i=1}^n (X_i - \bar{X}) \\ &= \underbrace{\sum_{i=1}^n Y_i - n\bar{Y}}_{\text{ομοίως}} - n\bar{Y} - n\bar{Y} = 0 \end{aligned}$$

Αρα $\sum_{i=1}^n e_i = 0$.

ΑΝΑΛΥΣΗ ΔΙΑΚΥΜΑΝΣΗΣ στο μοντέλο αgh.

Αν W_1, \dots, W_n είναι τυχαίο δείγμα, η δείγματική διακύμανση

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (W_i - \bar{W})^2.$$



$S^2 \leftarrow$ μέτρο μεταβλητότητας των παρατηρήσεων W_1, \dots, W_n σε σχέση με \bar{W} .

Ανάλογα η μεταβλητότητα των Y_1, \dots, Y_n είναι $\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{SS_{tot}} =$

αν $\pm \hat{Y}_i$ τότε :

$$\underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{SS_{reg}} + \underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{SS_{res}} \quad (\text{Αποδ. Άσκηση})$$

ολική μεταβλητότητα = Αθρ. τετραχ στο μοντέλο αgh + Αθρ. τετραχ Υπόλοιπα

$$SS_{tot} = SS_{reg} + SS_{res}$$



• Αν $SS_{reg} \gg SS_{res}$ τότε το μεγαλύτερο μέρος της ολικής μεταβλητότητας περιγράφεται από το μοντέλο της αgh.

Άρα το μοντέλο της αgh είναι υποσχόμενο

• Αν $SS_{reg} \ll SS_{res}$ τότε το μοντέλο της αgh δεν φαίνεται να είναι υποσχόμενο

ΠΙΝΑΚΑΣ ΑΝΑΔΙΑ μοντέλου αχπ

ΠΗΓΗ ΜΕΤΑΒΛΗΤΟΤΗΤΑΣ	SS	β.ε	MS	F πηλικο
ΜΟΝΤΕΛΟ αχπ	$SS_{reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	1	$MS_{reg} = \frac{SS_{reg}}{1}$	$F = \frac{MS_{reg}}{MS_{res}}$
ΥΠΟΛΟΙΠΑ	$= \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$ $SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$	(n-2)	$MS_{res} = \frac{SS_{res}}{n-2}$	
ΟΛΙΚΗ	$SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2$	n-1		

β.ε : πλήθος ανεξάρτητων πληροφοριών τις οποίες πρέπει να διαθέτουμε ώστε να μπορούμε να υπολογίσουμε το αντίστοιχο άθροισμα τετραγώνων.

Εμπειρικά στα Γραμμικά Μοντέλα

β.ε $SS_{tot} =$ μέγεθος δείγματος $- 1$

β.ε $SS_{μοντέλου} =$ πλήθος ανεξάρτητων μεταβλ.

β.ε $SS_{υπολοίπων} =$ αφαίρεση

Αν $MS_{reg} \gg MS_{res}$ τότε Υποσχόμενο Μοντέλο

Αν $MS_{reg} \ll MS_{res}$ τότε Μη Υποσχόμενο

Οι β.ε του SS_{tot} είναι $n-1$.

$$SS_{tot} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \Rightarrow Y_n + \sum_{i=1}^{n-1} Y_i = n\bar{Y} \Rightarrow Y_n = n\bar{Y} - \sum_{i=1}^{n-1} Y_i$$

$$Y_n - \bar{Y} = n\bar{Y} - \sum_{i=1}^{n-1} Y_i - \bar{Y} = (n-1)\bar{Y} - \sum_{i=1}^{n-1} Y_i \Rightarrow Y_n - \bar{Y} = - \sum_{i=1}^{n-1} (Y_i - \bar{Y})$$

$$\text{Άρα } (Y_n - \bar{Y})^2 = \left[\sum_{i=1}^{n-1} (Y_i - \bar{Y}) \right]^2$$

$$SS_{tot} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^{n-1} (Y_i - \bar{Y})^2 + (Y_n - \bar{Y})^2$$

$$= \sum_{i=1}^{n-1} (Y_i - \bar{Y})^2 + \left[\sum_{i=1}^{n-1} (Y_i - \bar{Y}) \right]^2$$

απαιτούνται $n-1$
ανεξ. πληρωφ. για να
υπολογιστεί.

οι $Y_1 - \bar{Y}, Y_2 - \bar{Y}, \dots, Y_{n-1} - \bar{Y}$.

$$\text{Αποδ } SS_{reg} = \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2$$

⇓

$$SS_{reg} = \sum_{i=1}^n (Y_i - \hat{Y})^2 = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 X_i - \bar{Y})^2$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \quad \sum_{i=1}^n (\bar{Y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 X_i - \bar{Y})^2 = \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2$$

Συντελεστής Προσδιορισμού ή Συντελεστής Προσαρμοστικότητας

$$R^2 = \frac{SS_{reg}}{SS_{tot}}, \quad R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

ΙΔΙΟΤΗΤΕΣ R^2

① $0 \leq R^2 \leq 1$ είναι καθαρός αριθμός (απαλλαγμένος από μονάδες μέτρησης)

② $0 \leq R^2 \leq 1$

Λόγω της ② ο R^2 είναι ποσοστό και εκφράζει το ποσοστό της ολικής μεταβλητότητας της εξαρτημένης μεταβλητής που ερμηνεύεται από το μοντέλο της αχπ.

Τιμές του R^2 κοντά στο 1 σημαίνουν υποσχόμενο μοντέλο

Μικρές τιμές του R^2 σημαίνουν όχι υποσχόμενο μοντέλο.

Για την περαιτέρω ανάπτυξη του μοντέλου της αχπ $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$, $i=1, \dots, n$

απαιτούνται κάποιες πιθανοθεωρητικές υποθέσεις που διατυπώνονται πάνω στα σφάλματα ϵ_i , $i=1, \dots, n$ του μοντέλου.

ΥΠΟΘΕΣΕΙΣ ΓΙΑ ΤΑ ΣΦΑΛΜΑΤΑ.

Τα σφάλματα είναι τ.μ.

① $E(\epsilon_i) = 0$, $\forall i=1, \dots, n$

② Η διακύμανση των σφαλμάτων είναι σταθερή

$Var(\epsilon_i) = \sigma^2$, $\forall i=1, \dots, n$

③ Τα σφάλματα ϵ_i είναι ασυσχέτιστα

$Cov(\epsilon_i, \epsilon_j) = 0$, $i \neq j$, $i, j=1, 2, \dots, n$

④ $\epsilon_i \sim N(0, \sigma^2)$, $i=1, \dots, n$

Οι υποθέσεις για τα σφάλματα αντανακλούν και στα Y_i

① $E(Y_i) = \beta_0 + \beta_1 X_i$, $i=1, \dots, n$

② $Var(Y_i) = \sigma^2$, $i=1, \dots, n$

③ $Cov(Y_i, Y_j) = 0$, $i \neq j$, $i, j=1, \dots, n$ τα Y_i ασυσχέτιστα

④ $Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$

ΥΠΕΝΟΜΗΣΗ

$Cov(Z, W) = E[(Z - E(Z))(W - E(W))]$

$= E(ZW) - (E(Z))(E(W))$

Αν $Z=W$

$Cov(Z, W) = Var(Z)$

ΘΕΩΡΗΜΑ Αν οι υποθέσεις για τα σφάλματα ικανοποιούνται τότε:

$$\alpha) \hat{\beta}_0 \sim N(\beta_0, \text{Var}(\hat{\beta}_0) = \frac{\sum x_i^2}{n \sum (x_i - \bar{x})^2} \sigma^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right])$$

↑
αμερόληπτος

$$\beta) \hat{\beta}_1 \sim N(\beta_1, \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2})$$

Αποδ

$$\beta) \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x}) y_i - \bar{y} \sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

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

ο $\hat{\beta}_1$ είναι γραμμικός συνδυασμός των $y_i, i=1, 2, \dots, n$

Ισχύει ότι αν $W_i \sim \text{Normal}$ τότε κάθε γραμμικός συνδυασμός των $W_i \sim \text{Normal}$

Άρα και ο $\hat{\beta}_1 \sim \text{Normal}$ ως γραμμικός συνδυασμός των $y_i \sim \text{Normal}$ λόγω της (4).

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

$$= \sum \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2} (\beta_0 + \beta_1 x_i)$$

$$= \beta_0 \sum \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2} + \beta_1 \sum \frac{x_i^2 - \bar{x} x_i}{\sum (x_i^2 - 2\bar{x} x_i + \bar{x}^2)} = \beta_0 \frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} + \beta_1 \frac{\sum x_i^2 - \bar{x} \sum x_i}{\sum x_i^2 - 2\bar{x} \sum x_i + n\bar{x}^2}$$

$$= \beta_1 \frac{\sum x_i^2 - \bar{x} n \bar{x}}{\sum x_i^2 - 2\bar{x} n \bar{x} + n \bar{x}^2} = \beta_1 \frac{\sum x_i^2 - n \bar{x}^2}{\sum x_i^2 - n \bar{x}^2} = \beta_1$$

Από ΘΠΣ ισχύει:

$$\text{Αν } W_i \text{ } i=1, \dots, n \text{ είναι τ.μ. τότε } \text{Var}\left(\sum_{i=1}^n \alpha_i W_i\right) = \sum_{i=1}^n \alpha_i^2 \text{Var}(W_i) + \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \text{Cov}(W_i, W_j)$$

Ειδική περίπτωση: Αν W_i αμοιβάτες ($\text{Cov}(W_i, W_j) = 0$) τότε

$$\text{Var}\left(\sum_{i=1}^n \alpha_i W_i\right) = \sum_{i=1}^n \alpha_i^2 \text{Var}(W_i)$$

$$\text{Var}(\hat{\beta}_1) = \text{Var}\left(\sum_{i=1}^n \frac{X_i - \bar{X}}{\sum (X_i - \bar{X})^2} Y_i\right) \quad \begin{array}{l} \text{Υποθέσεις} \\ \text{Υι αμοιβάτες} \end{array} \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sum (X_i - \bar{X})^2}\right)^2 \text{Var}(Y_i)$$

$$\begin{array}{l} \text{Υποθέσεις} \\ \text{Var}(Y_i) = \sigma^2 \end{array} \sigma^2 \sum \frac{(X_i - \bar{X})^2}{[\sum (X_i - \bar{X})^2]^2} = \sigma^2 \frac{\sum (X_i - \bar{X})^2}{[\sum (X_i - \bar{X})^2]^2} = \frac{\sigma^2}{\sum (X_i - \bar{X})^2}$$

Παρατηρώ ότι η $\sigma^2 (= \text{Var}(E_i) = \text{Var}(Y_i))$ είναι αίχμηση

Άρα η σ^2 πρέπει να εκτιμηθεί.

ΘΕΩΡΗΜΑ Υπό τις υποθέσεις για τα σφάλματα του μοντέλου της αιχμής το MS_{res} είναι αμερόληπτος εκτιμητής της σ^2 .

$$\text{δηλ } E(MS_{res}) = \sigma^2.$$

Αποδ σε φυλλάδιο που έδωσε.

Απλοποιώντας αλγόριθμο:

$$E(SS_{res}) = \sigma^2$$

Είναι $MS_{res} = \frac{SS_{res}}{n-2} = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

οπότε $SS_{res} = SS_{tot} - SS_{reg}$

$$E(SS_{tot}) = E\left\{ \sum_{i=1}^n (y_i - \bar{y})^2 \right\} = \sum_{i=1}^n E(y_i - \bar{y})^2 = \sum_{i=1}^n E\{ \beta_0 + \beta_1 x_i + \varepsilon_i - \beta_0 - \beta_1 \bar{x} - \bar{\varepsilon} \}^2 = \sum_{i=1}^n E\{ \beta_0 + \beta_1 x_i + \varepsilon_i - \beta_0 - \beta_1 \bar{x} - \bar{\varepsilon} \}^2$$

Μετα τις πράξεις και λαμβάνοντας υπόψη ότι:

$$E(\varepsilon_i - \bar{\varepsilon}) = E(\varepsilon_i) - \frac{1}{n} \sum_{i=1}^n E(\varepsilon_i) = 0 - 0 = 0$$

παιρνουμε:

$$E(SS_{tot}) = \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n E(\varepsilon_i - \bar{\varepsilon})^2 \quad (1)$$

Αλλά:

$$\begin{aligned} \sum_{i=1}^n E(\varepsilon_i - \bar{\varepsilon})^2 &= \sum_{i=1}^n (E(\varepsilon_i^2) - 2E(\varepsilon_i \bar{\varepsilon}) + E(\bar{\varepsilon}^2)) = \sum_{i=1}^n E(\varepsilon_i^2) - 2E\left(\sum_{i=1}^n \varepsilon_i\right)\bar{\varepsilon} + nE(\bar{\varepsilon}^2) \\ &= \sum_{i=1}^n E(\varepsilon_i^2) - 2E(n\bar{\varepsilon} \cdot \bar{\varepsilon}) + nE(\bar{\varepsilon}^2) = \sum_{i=1}^n E(\varepsilon_i^2) - nE(\bar{\varepsilon}^2) \\ &= \sum_{i=1}^n [Var(\varepsilon_i) + \{E(\varepsilon_i)\}^2] - n[Var(\bar{\varepsilon}) + \{E(\bar{\varepsilon})\}^2] \\ &= \sum_{i=1}^n \sigma^2 - nVar\left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i\right) = n\sigma^2 - n \frac{1}{n^2} \sum_{i=1}^n Var(\varepsilon_i) = n\sigma^2 - \frac{1}{n} \sum_{i=1}^n \sigma^2 \end{aligned}$$

και τελικά: $\sum_{i=1}^n E(\varepsilon_i - \bar{\varepsilon})^2 = (n-1)\sigma^2 \quad (2)$

Από τις (1) και (2):

$$E(SS_{tot}) = \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 + (n-1)\sigma^2 \quad (3)$$

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

και $E(SS_{reg}) = \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 \right\} E(\hat{\beta}_1^2)$

$$E(SS_{\text{reg}}) = \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 \right\} \left[\text{Var}(\hat{\beta}_1) + [E(\hat{\beta}_1)]^2 \right].$$

Λαμβάνοντας υπόψη ότι

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{και} \quad E(\hat{\beta}_1) = \beta_1,$$

Μετα από πράξεις προκύπτει:

$$E(SS_{\text{reg}}) = \sigma^2 + \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 \quad (4)$$

Έτσι από τις (3) και (4) έχουμε:

$$\begin{aligned} E(MS_{\text{res}}) &= \frac{1}{n-2} E(SS_{\text{res}}) = \frac{1}{n-2} \left[E(SS_{\text{tot}}) - E(SS_{\text{reg}}) \right] \\ &= \frac{1}{n-2} \left\{ \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 + (n-1)\sigma^2 - \sigma^2 - \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 \right\} \\ &= \frac{1}{n-2} (n-2)\sigma^2 \\ &= \sigma^2 \quad \blacktriangle \end{aligned}$$

ΘΕΩΡΗΜΑ Οι ΕΕΤ $\hat{\beta}_0$ και $\hat{\beta}_1$ είναι ανεξάρτητα των SS_{res} , MS_{res} .

ΘΕΩΡΗΜΑ Υπό τις υποθέσεις για τα σφάλματα

$$\frac{SS_{res}}{\sigma^2} \sim \chi^2_{n-2}$$

Απόδ

Λόγω των υποθέσεων για τα σφάλματα $\forall i=1,2,\dots,n$ τα $Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$

άρα (τυποποίηση) $\frac{Y_i - (\beta_0 + \beta_1 X_i)}{\sigma} \sim N(0,1) \quad \forall i=1,\dots,n$

$$\Rightarrow \frac{[Y_i - (\beta_0 + \beta_1 X_i)]^2}{\sigma^2} \sim N^2(0,1) \equiv \chi^2_1 \quad \forall i=1,\dots,n$$

για να έφ

$$\Rightarrow \sum \frac{(Y_i - (\beta_0 + \beta_1 X_i))^2}{\sigma^2} \sim \chi^2_{\sum_{i=1}^n 1} \equiv \chi^2_n \quad (*)$$

λόγω του ότι
τα Y_i ασυσχέτιστα
και κανονικά

$$SS_{res} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$(*) \Rightarrow \sum_{i=1}^n \frac{(Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2}{\sigma^2} \sim \chi^2_{n-2}$$

$$\Rightarrow \sum_{i=1}^n \frac{(Y_i - \hat{Y}_i)^2}{\sigma^2} \sim \chi^2_{n-2}$$