

Παλινδρόμηση & Ανάλυση Διακύμανσης

Γενική γραμμική υπόθεση.

Θεωρούμε το μοντέλο π.χ.π $Y = X\beta + \underline{\varepsilon}$, $\underline{\varepsilon} \sim N_n(0, \sigma^2 I_n)$

Έστω ένας πίνακας σταθερών A διαστάσεων $q \times (p+1)$

Γενική Γραμμική Υπόθεση: $H_0: A\beta = \underline{c}$

$\begin{matrix} & \downarrow & \downarrow & \downarrow & \downarrow \\ & q \times (p+1) & (p+1) \times 1 & q \times 1 & q \times 1 \\ & \underbrace{\hspace{2cm}} & & & \\ & q \times 1 & & & \end{matrix}$

Η H_0 περιλαμβάνει ως ειδικές περιπτώσεις τις ακόλουθες ενδιαφέρουσες:

① Αν $A \equiv I_{p+1}$ και $\underline{c} = \underline{0}_{p+1}$ τότε $H_0: A\beta = \underline{c} \rightarrow H_0: \beta_0 = \beta_1 = \dots = \beta_p = 0$
 Αν δεν μπορεί να απορριφθεί \rightarrow ~~μοντέλο π.χ.π~~

② Αν $A = I_{p+1}$ και $\underline{c} = \beta^*$ με β^* γνωστό τότε $H_0: A\beta = \underline{c}$
 $\rightarrow H_0: \beta_0 = \beta_0^*, \dots, \beta_p = \beta_p^*$
 Αν δεν μπορεί να απορριφθεί τότε η μεταβολή της Y για μοναδιαία μεταβολή της X_i όταν οι υπολοίπες X σταθερές είναι ίση με β_i^*

③ Αν $A = \begin{pmatrix} 1 & & & & \\ & 0 & & & \\ & & 1 & & \\ & & & \ddots & \\ 0 & & & & 1 & \\ & & & & & 0 \end{pmatrix}$ και $\underline{c} = \underline{0}$ τότε $H_0: A\beta = \underline{c} \rightarrow H_0: \beta_m = 0$
 όπου β_m διάνυσμα που αποτελείται από κάποιες

μηδενικές συνιστώσες του β

π.χ αν $A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & & 0 & 0 \end{pmatrix}$ τότε $A\beta = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & & 0 & 0 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} = \begin{pmatrix} 0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ 0 \end{pmatrix}$

Άρα $H_0: A\beta = \underline{c} = \underline{0} \rightarrow H_0: \beta_1 = \beta_2 = 0$

◦ Είναι σημαντική η $H_0: \beta_1 = \beta_2 = 0$;

ΝΑΙ παίρα πολύ γιατί αν δεν μπορού να την απορρίψω τότε οι όροι που περιέχουν την X_1 και X_2 πρέπει να διαγραφούν και το μοντέλο θα γίνει: $Y_i = \beta_0 + \beta_3 X_{i3} + \beta_4 X_{i4} + \dots + \beta_p X_{ip} + \epsilon_i$

Περιγραφή της κατασκευής στατιστικού τεστ για τον έλεγχο της $H_0: A\beta = \underline{c}$.

Έστω ΠΜ το πλήρες μοντέλο: $Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i, i=1, \dots, n$

Έστω H_0M το H_0 -μοντέλο που διαμορφώνεται από το ΠΜ αν λάβω υπόψη μου την $H_0: A\beta = \underline{c}$.

πχ Έστω η $H_0: A\beta = \underline{c}$ είναι:

$$H_0: \begin{pmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \dots & 0 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

↓

$$H_0: \beta_1 = 0 \text{ και } \beta_2 = 0$$

Άρα το H_0M είναι $Y_i = \beta_0 + \beta_3 X_{i3} + \dots + \beta_p X_{ip} + \epsilon_i, i=1, \dots, n$

Έστω $\hat{\beta}$ και $\hat{\beta}_{H_0}$ οι ΕΕΤ της παραμέτρου β υπό το ΠΜ και H_0M , αντίστοιχα.

Οι εκτιμώμενες τιμές του Y είναι

$$\underline{\hat{Y}} = X \hat{\beta} \text{ για το ΠΜ}$$

$$\underline{\hat{Y}}_{H_0} = X \hat{\beta}_{H_0} \text{ για το } H_0M$$

Επομένως τα αντίστοιχα αθροίσματα τετραγώνων των υπολοίπων θα είναι:

$$SS_{res}(\Pi M) \cong (\underline{y} - \hat{\underline{y}})' (\underline{y} - \hat{\underline{y}})$$

$$SS_{res}(H_0 M) \cong (\underline{y} - \hat{\underline{y}}_{H_0})' (\underline{y} - \hat{\underline{y}}_{H_0})$$

Αν $SS_{res}(H_0 M) \gg SS_{res}(\Pi M)$ τότε το μοντέλο κάτω από την $H_0: A\beta = \zeta$ ερμηνεύει πολύ μικρότερο κομμάτι από την ολική μεταβλητότητα και επομένως είναι το $H_0 M$ λιγότερο προτιμητέο από το ΠM .
Άρα το $H_0 M$ πρέπει να απορριφθεί

Άρα το τεστ για τον έλεγχο θα πρέπει να βασιστεί στη σύγκριση των $SS_{res}(H_0 M)$ με το $SS_{res}(\Pi M)$.

Αποδεικνύεται (αλγεβρικά):

$$SS_{res}(H_0 M) - SS_{res}(\Pi M) = (A\hat{\beta} - \zeta)' [A(X'X)^{-1}A']^{-1} (A\hat{\beta} - \zeta)$$

Επίσης αποδεικνύεται (πιθανοθεωρητικά με κατανομές τετραγωνικών μορφών).

$$\frac{SS_{res}(H_0 M) - SS_{res}(\Pi M)}{\sigma^2} \sim \chi^2_q \text{ υπό την } H_0: A\beta = \zeta$$

$$\text{Επίσης } \frac{SS_{res}(\Pi M)}{\sigma^2} \sim \chi^2_{n-p-1}$$

Επίσης $\hat{\beta}$ ανεξ του $SS_{res}(\Pi M)$

Άρα και $SS_{res}(H_0 M) - SS_{res}(\Pi M)$ ανεξ του $SS_{res}(\Pi M)$

Έτσι υπό $H_0: A\beta = \zeta$ η στατιστική σωλην:

$$F = \frac{SS_{res}(H_0 M) - SS_{res}(\Pi M) / \sigma^2 \cdot q}{SS_{res}(\Pi M) / \sigma^2 (n-p-1)} \sim F_{q, n-p-1} \text{ υπό την } H_0$$

Συγκεντρωτικά για τον έλεγχο της $H_0: A\beta = \xi$
 χρησιμοποιείται η στατιστική συνάρτηση:

$$F = \frac{(n-p-1)(A\hat{\beta} - \xi)' [A(X'X)^{-1}A']^{-1} (A\hat{\beta} - \xi)}{q \cdot SS_{res}} \sim F_{q, n-p-1} \text{ υπό } H_0$$

και κ.π μεγάλες τιμές του F ($F \geq c$, το c υπολογ από $\alpha = P(A \text{ απόρ } H_0 / H_0 \text{ αληθεία})$)
 $F \geq F_{q, n-p-1, \alpha}$

Διαρεύνηση των Ιδιοτήτων των υπολοίπων *

Ανοίγηση Υπολοίπων

Αξιοποίηση των υπολοίπων στον έλεγχο **

των υποθέσεων για τα σφάλματα.

(αφού τα υπόλοιπα εστιμούν τα σφάλματα του μοντέλου)

* Μοντέλο: $\underline{Y} = X\underline{\beta} + \underline{\varepsilon}$

Υπόλοιπα: $\underline{e} = \underline{Y} - \hat{\underline{Y}} = \underline{Y} - X\hat{\underline{\beta}}$, $\hat{\underline{\beta}} = (X'X)^{-1}X'\underline{Y}$

ΙΔΙΟΤΗΤΕΣ ΤΩΝ ΥΠΟΛΟΙΠΩΝ.

1. Τα υπόλοιπα \underline{e} είναι γραμμική συνάρτηση των $\underline{\varepsilon}$ και \underline{Y}
 Πράγματι,

$$\underline{e} = \underline{Y} - \hat{\underline{Y}} = \underline{Y} - X\hat{\underline{\beta}} = \underline{Y} - X(X'X)^{-1}X'\underline{Y} = (\underbrace{I_n - X(X'X)^{-1}X'}_p)\underline{Y}$$

$$\Rightarrow \underline{e} = (I_n - p)\underline{Y} \text{ όπου } p = X(X'X)^{-1}X'$$

Επίσης,

$$\begin{aligned} \underline{e} &= (I_n - p)\underline{Y} = (I_n - p)(X\underline{\beta} + \underline{\varepsilon}) = (I_n - p)X\underline{\beta} + (I_n - p)\underline{\varepsilon} \\ &= X\underline{\beta} - X\underline{\beta} + (I_n - p)\underline{\varepsilon} = X\underline{\beta} - X(X'X)^{-1}X'X\underline{\beta} + (I_n - p)\underline{\varepsilon} \\ &= X\underline{\beta} - X\underline{\beta} + (I_n - p)\underline{\varepsilon} \Rightarrow \underline{e} = (I_n - p)\underline{\varepsilon} \end{aligned}$$

ΠΑΡΑΤΗΡΗΣΗ: Για τον $P = X(X'X)^{-1}X'$

i) Ο P είναι συμμετρικός $P = P'$

ii) Ο P είναι ταυτοδύναμος, $P^2 = P$.

Αποδ

$$i) P' = [X(X'X)^{-1}X']' \stackrel{(AB)'}{=} \underbrace{X(X'X)^{-1}X'}_{(A^{-1})' = (A')^{-1}} = P$$

$$ii) P^2 = [X(X'X)^{-1}X'] [X(X'X)^{-1}X'] = X(X'X)^{-1} \underbrace{X'X(X'X)^{-1}}_{I_n} X' = X(X'X)^{-1}X' = P.$$

2. $\underline{e} \sim N_n(\underline{0}, (I_n - P)\sigma^2)$

Επειδή $\underline{e} = (I_n - P)\underline{\xi}$ και $\underline{\xi} \sim N_n \Rightarrow \underline{e} \sim N_n$

$$E(\underline{e}) = E[(I_n - P)\underline{\xi}] = (I_n - P)E(\underline{\xi}) \stackrel{\text{από υποθ. σφαλμάτων}}{=} (I_n - P) \cdot \underline{0} = \underline{0}$$

$$\text{Ισχύει } \text{Var}(A\underline{w}) = A \cdot \text{Var}(\underline{w}) \cdot A'$$

$$\text{Var}(\underline{e}) = \text{Var}[(I_n - P)\underline{\xi}] \Rightarrow$$

$$\text{Var}(\underline{e}) = (I_n - P) \text{Var}(\underline{\xi}) (I_n - P)'$$

$$= (I_n - P) \sigma^2 (I_n - P) \quad \underline{P = P'}$$

$$= (I_n - P) \cdot \sigma^2 (I_n - P)$$

$$= \sigma^2 (I_n I_n - P I_n - I_n P + P \cdot P) \quad \underline{P^2 = P}$$

$$= \sigma^2 (I_n - \underbrace{2P}_{I_n}) + P) = \sigma^2 (I_n - P)$$

ΠΑΡΑΤΗΡΗΣΗ: Αφού $\text{Var}(\underline{\xi}) = \sigma^2 (I_n - P)$

$$\text{Var}(e_i) = i\text{-}\text{διαγώνιο στοιχείο του } \sigma^2 (I_n - P)$$

$$= \sigma^2 (1 - p_{ii}) \text{ όπου } p_{ii} \text{ το } i\text{-}\text{διαγώνιο στοιχείο του } P$$

$$\text{Cov}(e_i, e_j) = -\sigma^2 p_{ij}, \text{ όπου } p_{ij} \text{ το } (ij) \text{ στοιχείο του } P$$

3. Τα e_i και \hat{Y}_i ασυσχέτιστα
 δηλ $\text{Cov}(e_i, \hat{Y}_i) = 0$.

Αποδ

$$\text{Cov}(\underline{e}, \underline{\hat{Y}}) = \text{Cov}((I_n - P)\underline{Y}, X\hat{\beta})$$

$$= \text{Cov}((I_n - P)\underline{Y}, X(X'X)^{-1}X'\underline{Y})$$

$$= \text{Cov}((I_n - P)\underline{Y}, P\underline{Y})$$

ΘΠΣ Ισχύει $\text{Cov}(A\underline{Y}, B\underline{Y}) = A\text{Var}(\underline{Y})B'$

$$= \text{Cov}(\underline{e}, \underline{\hat{Y}}) = (I_n - P)\text{Var}(\underline{Y})P'$$

$$= (I_n - P) \cdot (\sigma^2 I_n) P$$

$$= \sigma^2 (I_n - P) P$$

$$= \sigma^2 [I_n - P - P^2]$$

$$= \sigma^2 (P - P)$$

$$= \sigma^2 \cdot 0 \quad \text{πινδενικιστικόνουαοις}$$

$$= 0$$

4. Μαθηματικοποιημένα Υπόλοιπα.

Ορίζονται:

$$t_i = \frac{e_i}{\hat{\sigma} \sqrt{1 - p_{ii}}}$$

$\hat{\sigma} = \sqrt{MS_{res}}$, πη το i -διαγώνιο
στοιχείο του $P = X(X'X)^{-1}X'$

$$= \frac{e_i}{\hat{\text{Var}}(e_i)}$$

Ισχύει $t_i \sim t_{n-p-1}$.

Αποδ $t_i \sim t_{n-p-1}$

Επειδή $\underline{e} \sim N_n(0, \sigma^2(I_n - P)) \rightarrow \omega \ e_i \sim N(0, \sigma^2(1 - p_{ii}))$

$$\Rightarrow \frac{e_i - 0}{\sigma \sqrt{1 - p_{ii}}} \sim N(0, 1) \quad \eta \quad \frac{e_i}{\sigma \sqrt{1 - p_{ii}}} \sim N(0, 1)$$

Επιπλέον $\frac{SS_{res}}{\sigma^2} \sim \chi^2_{n-p-1}$

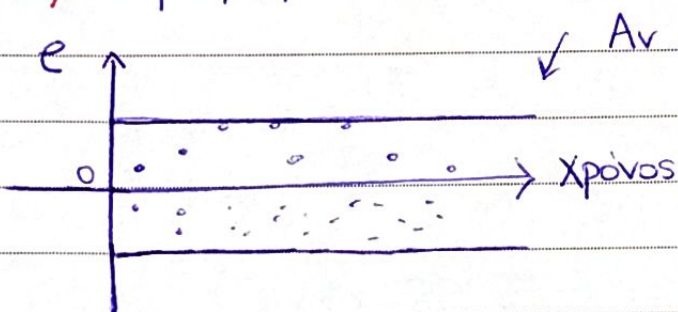
$$t_i \text{ of } \frac{e_i}{\hat{\sigma}\sqrt{1-p_{ii}}} = \frac{e_i/\sqrt{1-p_{ii}}}{\sqrt{MS_{res}}} = \frac{e_i/\sigma\sqrt{1-p_{ii}}}{\sqrt{\frac{SS_{res}}{(n-p-1)\sigma^2}}} \sim \frac{N(0,1)}{\sqrt{\frac{\chi^2_{n-p-1}}{(n-p-1)}}} \sim t_{n-p-1}$$

γιατι e_i και SS_{res} ανεξαρτητα αφού
 $\underline{e} = \underline{y} - \hat{\underline{y}} = \underline{y} - X\hat{\underline{\beta}}$ και $\hat{\underline{\beta}}$ ανεξ SS_{res} .

**

I) Έλεγχος του αουσχέτιστου των σφαλμάτων: $Cov(e_i, e_j) = 0$.

α) Γραφική Μέθοδος



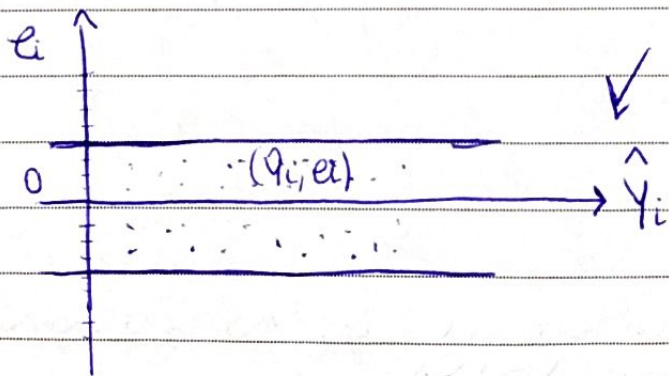
✓ Αν τα υπόλοιπα e_i κατανομονται τυχαία σε μια ζώνη γύρω από το 0, τότε έχω ένδειξη του αουσχέτιστου

β) Τεστ Ραών

γ) Τεστ Durbin-Watson

II) Έλεγχος της υπόθεσης της σταθερής διακύμανσης $Var(e_i) = \sigma^2$
 $i=1, \dots, n$

α) Γραφική Μέθοδος (Γραφική παράσταση των $e_i, \hat{y}_i \quad i=1, \dots, n$)



✓ Αν τα $(\hat{y}_i, e_i) \quad i=1, \dots, n$ κατανομονται τυχαία σε μια ζώνη γύρω από το 0 έχω ισχυρή ένδειξη ότι η διακύμανση των σφαλμάτων είναι κοινή και σταθερή

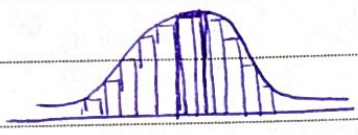
Θεωρητική Δικαιολόγηση

Αποδ. ότι $Var(\epsilon_i) = \sigma^2$ τότε $Cov(\epsilon_i, \hat{Y}_i) = 0$.

β) Τεστ του Levene.

III) Κανονικότητα των σφαλμάτων $\epsilon_i \quad i=1, \dots, n$
δηλ ότι τα $\epsilon_i \sim N$

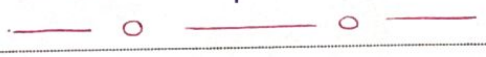
α) Ιστογράμμα $\epsilon_1, \dots, \epsilon_n$



β) Normal probab. Plots ← SPSS κλπ.

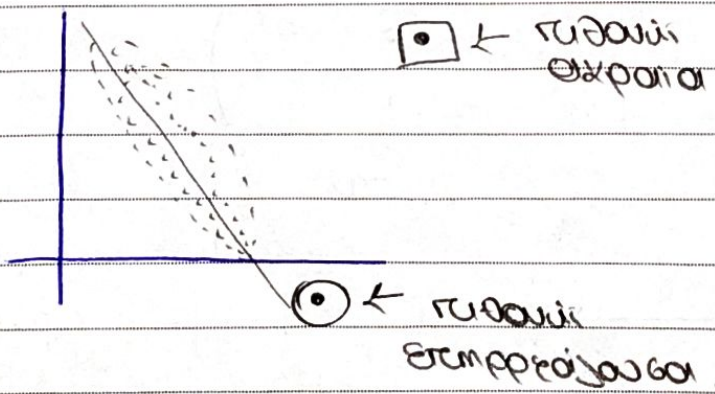
γ) Kolmogorov-Smirnov

δ) Shapiro - Wilks



► Ακραίες - Επιπρεσζαυσες παρατηρήσες

Γενικά παρατηρήσες που διαφοροποιούνται αισθητά από τις υπόλοιπες



• Μια παρατηρήσση Y_k είναι ακραία αν το κ-υπόλοιπο εκ του συστωλικού β' ουτμ είναι $|R_k| > 3$.

• Μια παρατηρήσση Y_k είναι επιπρεσζαυσα

→ ΤΕΣΤ ΒΕΛΓΙΣΤΟΥ ΒΕΜΝ ΟΥΣΤΑΘΜ COOK.