

Θέμα

A1)

Επειδή οι τιμές των ε είναι άγνωστες, η εξέταση των υποθέσεων για τα σφάλματα γίνεται υποχρεωτικά μέσω των υπολοίπων

$$e = y - \hat{y} = y - X\hat{\beta} = y - X(X'X)^{-1}X'y = y - HY = (I - H)y \\ = (I - H)(X\beta + \varepsilon) = X\beta - HX\beta + \varepsilon - H\varepsilon \Rightarrow$$

$$e = (I - H)\varepsilon$$

$$E(e) = E[(I - H)\varepsilon] = (I - H)E(\varepsilon) = 0 \Rightarrow E(e) = 0$$

$$V(e) = V[(I - H)\varepsilon] = (I - H)V(\varepsilon)(I - H)' = \sigma^2(I - H)$$

$$\text{Cov}(e_i, e_j) = \sigma^2 h_{ij} \neq 0$$

$$A2) f(y_i) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left[-\frac{(y_i - \beta'x_i)^2}{2\sigma^2}\right] \Rightarrow$$

$$f(y_i) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left[-\frac{(y_i - x_i'\beta)^2}{2\sigma^2}\right]$$

$$L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - x_i'\beta)^2}{2\sigma^2}\right] \Rightarrow$$

$$L(\beta, \sigma^2) = L = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{(y - X\beta)'(y - X\beta)}{2\sigma^2}\right] \Rightarrow$$

$$\ln L = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} (y - X\beta)'(y - X\beta)$$

$$\frac{\partial \ln L}{\partial \beta} = -\frac{1}{\sigma^2} (X'X\beta - X'y) \Rightarrow \hat{\beta} = (X'X)^{-1}X'y$$

$$\hat{\sigma}^2 = \frac{(y - X\hat{\beta})(y - X\hat{\beta})}{n} = \frac{SSE}{n}$$

Άρα η μεγιστοποιημένη συνάρτηση πιθανοφάνειας θα είναι

$$\hat{l} = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\hat{\sigma}^2) - \frac{n\hat{\sigma}^2}{2\hat{\sigma}^2} \Rightarrow$$

$$\hat{l} = -\frac{n}{2} [\ln(2\pi) + \ln \hat{\sigma}^2 + 1] \Rightarrow$$

$$\hat{l} = -\frac{n}{2} \left[\ln(2\pi \cdot SSE/n) + 1 \right]$$

ii) Το AIC αποτελεί ένα κριτήριο επιλογής των βέλτιστων μοντέλων με το όσο το δυνατόν μικρότερο αριθμό παραμέτρων. Στην γενική περίπτωση ορίζεται από την σχέση

$$AIC = 2d - 2\ln L$$

όπου d το πλήθος των παραμέτρων του μοντέλου και L η μεγιστοποιημένη τιμή της συνάρτησης πιθανοφάνειας.

$$AIC = 2d - 2\ln L = 2d - 2\left(-n\ln\hat{\sigma}^2 - \frac{1}{2}n - \frac{n}{2}\ln(2\pi)\right)$$

$$= 2(p+1) + (2n\ln\hat{\sigma}^2 + n + n\ln(2\pi)) \Rightarrow$$

$$AIC = 2(p+1) + n \left[\ln(2\pi) + 1 + \ln\left(\frac{SSE}{n}\right) \right]$$

A3) Οι έλεγχοι και τα κριτήρια που αφορούν την επιλογή των καταλληλότερων μεταβλητών για την ανάλυση του μοντέλου είναι τα μέτρα καταλληλότητας, οι στατιστικοί έλεγχοι και οι γραφικές παραστάσεις.

Για την αφαίρεση μη επεξηγηματικής μεταβλητής είναι ανάγκη να γίνει ο έλεγχος F για την ουσιαστικότητα των συντελεστών της μεταβλητής αυτής. Έπειτα ~~εξετάζονται~~ εξετάζονται και οι τιμές των R^2 , \bar{R}^2 και AIC του μοντέλου με και χωρίς την μεταβλητή αυτή.

A4) Τα τυποποιημένα υπολοίπα ορίζονται ως :

$$r_i = \frac{e_i}{S\sqrt{1-h_{ii}}} \quad i=1, 2, \dots, n$$

Η τυποποίηση των υπολοίπων αφαιρεί την ετεροσκεδαστικότητα των συνήθων υπολοίπων και έτσι μπορούμε να προβούμε σε ασφαλέστερες αξιολογήσεις για τις υποθέσεις των τυχαιών σφαλμάτων.

Τα deleted υπολοίπα αποτελούν μια τροποποίηση των τυποποιημένων υπολοίπων σύμφωνα με την οποία χρησιμοποιείται μια εναλλακτική εκτίμηση της σ^2 , την οποία αποκτάμε αφού πρωταρχικά το μοντέλο χρησιμοποιώντας όλες τις παρατηρήσεις πλην της i -οστής.

$$r_i^* = \frac{e_i}{S(e)\sqrt{1-h_{ii}}} \quad i=1, 2, \dots, n$$

Τα διαγώνια στοιχεία h_{ii} ονομάζονται μέγεθος. Οι εν λόγω τιμές μετρούν την επίδραση της παρατηρούμενης τιμής της μεταβλητής απόκρισης πάνω στις προβλεπόμενες τιμές.

Οι τιμές των h_{ii} εμφανίζονται από 0 έως 1. Η αναμενόμενη τιμή της είναι $(p+1)/n$. Αν υπάρχουν τιμές h_{ii} 3 φορές μεγαλύτερες από την μέση τιμή τότε έχουμε εκτροπές τιμές στην i -γραφή του πίνακα σχεδιασμού, οι οποίες πιθανότατα διαφοροποιούν έντονα την προσήλωση του ματέλα.

Οι τιμές μοχλίσματος χρησιμοποιούνται επίσης και για τον υπολογισμό των αποστάσεων Cook για να μετράν την συνολική επίρροη της κάθε παρατήρησης στο ματέλο. Τιμές των D_i άνω της μονάδας πρέπει να μας προβληματίσουν.

A5) Ελέγχουμε πρώτα την υπόθεση της απώσιας αλληλεπίδρασης μεταξύ της μεταβλητής, που κατασκευάζουμε, δηλαδή την ισοτιμία των κλίσεων

$$H_0: \beta_3 = 0 \text{ [κατάσταση Β]} \text{ έναντι } H_1: \beta_3 \neq 0 \text{ [κατάσταση Α]}$$

Με τον έλεγχο F:

$$F = \frac{(SSE_B - SSE_A)/1}{SSE_A / (n-4)}$$

Στην περίπτωση που δεν απορριφθεί η μηδενική υπόθεση στο πρώτο βήμα, δηλαδή δεν διαπιστώσαμε κάποια αλληλεπίδραση μεταξύ δύο μεταβλητών, προχωράμε στην εξέταση της τιμής του συντελεστή β_2 .

$$H_0: \beta_2 = 0 \text{ [κατάσταση Γ]} \text{ έναντι } H_1: \beta_2 \neq 0 \text{ [κατάσταση Β]} \\ (\text{αφού αποφασίσαμε ότι } \beta_3 = 0)$$

$$F = \frac{(SSE_\Gamma - SSE_B)/1}{SSE_B / (n-3)}$$

Στη περίπτωση που δεν απορρίψουμε την μηδενική
υπόθεση μπορούμε να πούμε ότι τα δεδομένα προέρχονται
από τον ίδιο ακριβώς πληθυσμό και μπορούν να θεωρηθούν ως
μία ομάδα.