

ΑΝΑΛΥΣΗ ΤΗΣ ΔΙΑΚΥΜΑΝΣΗΣ

Γενικά: Υπάρχει μια εξαρτημένη μεταβλητή Y (ποσοτική). Επίσης υπάρχουν μία ή περισσότερες ανεξάρτητες μεταβλητές. Η πιο συνηθισμένη περίπτωση είναι οι ανεξάρτητες μεταβλητές να είναι ποιοτικές.

Σε αντίθεση με την παλινδρόμηση, δεν μας ενδιαφέρει η σχέση ανεξάρτητων και εξαρτημένης, ούτε και η πρόβλεψη. Κυρίως τα προβλήματα είναι ελέγχου υποθέσεων.

Οι ανεξάρτητες μεταβλητές αναφέρονται σαν παράγοντες (factors).

Ανάλυση της Διακύμανσης κατά ένα παράγοντα

Μαθηματική Διατύπωση:

Έχουμε παρατηρήσεις $Y_{ij} \sim N(\mu_i, \sigma^2)$ $i = 1, \dots, I$, $j = 1, \dots, n_i$ και μπορούν να δοθούν υπό μορφή πίνακα:

$$\begin{array}{cccc} Y_{11} & Y_{12} & \dots & Y_{1n_1} \\ Y_{21} & Y_{22} & \dots & Y_{2n_2} \\ \vdots & & & \\ Y_{I1} & Y_{I2} & \dots & Y_{In_I} \end{array} \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \text{Κάθε γραμμή είναι τυχαίο δείγμα από την } N(\mu_i, \sigma^2)$$

Πρόβλημα: Να ελεγχθεί αν οι μέσες τιμές είναι ίσες.

$$H_0: \mu_1 = \dots = \mu_I$$

Εφαρμογή: Έστω ότι υπάρχουν $I = 4$ ποικιλίες κάποιου αγροτικού προϊόντος. Ενδιαφερόμαστε να εξετάσουμε ποια ποικιλία είναι καλύτερη με την έννοια ότι δίνει καλύτερη απόδοση.

Y_{ij} = απόδοση (παραγωγή) της i ποικιλίας στο j αγρόκτημα.

Θεωρώ τον έλεγχο: $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$ (δεν διαφέρουν οι ποικιλίες)

$$H_1: \mu_\kappa \neq \mu_\lambda \quad (\text{υπάρχουν διαφορές})$$

Εδώ έχουμε ανάλυση της διακύμανσης κατά ένα παράγοντα με 4 επίπεδα (levels).

Έλεγχος: Έστω $n = \sum_{i=1}^I n_i$ (ολικός αριθμός παρατηρήσεων).

$$\underline{\mu}_{n \times 1} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_2 \\ \vdots \\ \mu_I \\ \vdots \\ \mu_I \end{pmatrix} \quad \left. \begin{array}{l} \} n_1 \text{ φορές} \\ \\ \} n_2 \text{ φορές} \\ \\ \} n_I \text{ φορές} \end{array} \right\}$$
$$\underline{Y} = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ Y_{I1} \\ \vdots \\ Y_{In_I} \end{pmatrix} \quad \left. \begin{array}{l} \} n_1 \text{ φορές} \\ \\ \} n_2 \text{ φορές} \\ \\ \} n_I \text{ φορές} \end{array} \right\}$$

$$\underline{\mu}_0 = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_I \end{pmatrix}_{n \times 1}, \text{ όπου } \mu = \frac{n_1\mu_1 + n_2\mu_2 + \dots + n_I\mu_I}{n} = \frac{\sum_{i=1}^I n_i\mu_i}{n}$$

Όταν ισχύει η H_0 τότε έχω ότι $\underline{\mu}_0 = \mu \cdot \underline{1}$, $\underline{1} = (1, \dots, 1)^T$ και μ η κοινή άγνωστη τιμή των μ_i . Άρα, ο έλεγχος παίρνει τη μορφή:

$$H_0: \underline{\mu} = \underline{\mu}_0 \} \Leftrightarrow H_0: \|\underline{\mu} - \underline{\mu}_0\|^2 = 0$$

$$H_1: \underline{\mu} \neq \underline{\mu}_0 \} \quad H_1: \|\underline{\mu} - \underline{\mu}_0\|^2 > 0$$

$$\text{Έστω } \theta = \|\underline{\mu} - \underline{\mu}_0\|^2 = \sum_{i=1}^I n_i (\mu_i - \mu)^2$$

Άρα, έχω τον ισοδύναμο έλεγχο:

$$H_0: \theta = 0 \text{ προς } H_1: \theta > 0.$$

Πρέπει να εκτιμήσω την παράμετρο θ . Για να εκτιμήσω την θ έχω να εκτιμήσω τα μ_i και μ .

$$\text{Εκτιμητήρια του } \mu_i: \hat{\mu}_i = \frac{\sum_{j=1}^{n_i} Y_{ij}}{n_i} = \frac{Y_{i.}}{n_i} = \bar{Y}_{i.}$$

$$\text{Εκτιμητήρια του } \mu: \hat{\mu} = \frac{n_1\bar{Y}_{1.} + \dots + n_I\bar{Y}_{I.}}{n} = \bar{Y}_{..} \quad (= \frac{1}{n} \sum_{i=1}^I \sum_{j=1}^{n_i} Y_{ij}).$$

$$\text{Άρα, η εκτιμητήρια του } \theta \text{ είναι } \hat{\theta} = \sum_{i=1}^I n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2.$$

Άρα, απορρίπτουμε αν $\hat{\theta} > c$, όπου c σταθερά που υπολογίζεται από το επίπεδο σημαντικότητας.

Ισχύει η βασική σχέση:

$$\begin{aligned} \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 &= \sum_{i=1}^I \sum_{j=1}^{n_i} \{(Y_{ij} - \bar{Y}_{i.}) + (\bar{Y}_{i.} - \bar{Y}_{..})\}^2 \\ &= \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 + 2 \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})(\bar{Y}_{i.} - \bar{Y}_{..}) + \sum_{i=1}^I \sum_{j=1}^{n_i} (\bar{Y}_{i.} - \bar{Y}_{..})^2 \\ &= \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 + \sum_{i=1}^I n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2, \text{ αφού } \sum_{i=1}^I (\bar{Y}_{i.} - \bar{Y}_{..}) \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.}) = 0. \end{aligned}$$

$$\begin{aligned} \text{Άρα, } \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 &= \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 + \sum_{i=1}^I n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2. \\ n-1 \text{ β.ε. } SS_{total} &= n-I \text{ β.ε. } SS_{within} + I-1 \text{ β.ε. } SS_{between} \\ &\quad \uparrow \\ &\quad \text{(άθροισμα τετραγώνων} \\ &\quad \text{των υπολοίπων)} \end{aligned}$$

Όμως,
$$\sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 / \sigma^2 \sim \chi_{n-1}^2,$$

γιατί στην περίπτωση που ισχύει η H_0 τα $Y_{ij} \sim N(\mu, \sigma^2)$.

Για τον ίδιο λόγο έχω ότι
$$\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 / \sigma^2 \sim \chi_{n_i-1}^2$$

και άρα
$$\sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 / \sigma^2 \sim \chi_{\sum (n_i-1)}^2 = \chi_{n-I}^2,$$
 λόγω ανεξαρτησίας.

Άρα, αν αποδείξω ότι $\hat{\theta}$ και $\sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2$ είναι ανεξάρτητα, η βασική σχέση δίνει

ότι $\hat{\theta} / \sigma^2 \sim \chi_{I-1}^2.$

Αλλά
$$\text{Cov}(Y_{ij} - \bar{Y}_{i.}, \bar{Y}_{i.} - \bar{Y}_{..}) = \text{Cov}(Y_{ij}, \bar{Y}_{i.}) - \text{Cov}(Y_{ij}, \bar{Y}_{..}) - \text{Cov}(\bar{Y}_{i.}, \bar{Y}_{i.}) + \text{Cov}(\bar{Y}_{i.}, \bar{Y}_{..})$$

$$= \frac{\sigma^2}{n_i} - \frac{\sigma^2}{n} - \frac{\sigma^2}{n_i} + \frac{\sigma^2}{n} = 0$$
 και λόγω κανονικότητας έχω ανεξαρτησία.

Άρα, $\hat{\theta} / \sigma^2 \sim \chi_{I-1}^2.$

Επειδή όμως το σ^2 είναι άγνωστο, από τη σχέση

$$E(SS_{with}) = (n - I)\sigma^2 \Rightarrow E\left(\frac{SS_{with}}{n - I}\right) = \sigma^2.$$

Οπότε, ο έλεγχος γίνεται

$$F = \frac{\hat{\theta} / (I - 1)}{SS_{with} / (n - I)} = \frac{MS_{bet}}{MS_{with}} \sim F_{I-1, n-I}.$$

Απορρίπτω αν $F > F_{I-1, n-I, \alpha}.$