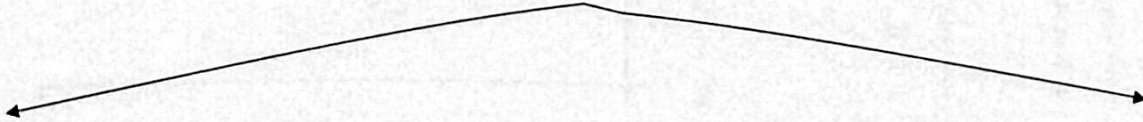


# ΓΡΑΜΜΙΚΑ ΜΟΝΤΕΛΑ

Μοντέλα για τη διερεύνηση και μελέτη της σχέσης μεταξύ μιας εξαρτημένης μεταβλητής και μιας ή περισσότερων ανεξάρτητων μεταβλητών



## ΜΟΝΤΕΛΑ ΠΑΛΙΝΔΡΟΜΗΣΗΣ

α) Η εξαρτημένη και οι ανεξάρτητες μεταβλητές είναι ποσοτικές.

β) Αναζητείται, κατασκευάζεται και ελέγχεται η ορθότητα της γραμμικής σχέσης που συνδέει την εξαρτημένη με τις ανεξάρτητες μεταβλητές. Η γραμμική σχέση που κατασκευάζεται χρησιμοποιείται για την πρόβλεψη της εξαρτημένης μεταβλητής για δεδομένες τιμές των ανεξάρτητων μεταβλητών.

## ΜΟΝΤΕΛΑ ΑΝΑΛΥΣΗΣ ΔΙΑΚΥΜΑΝΣΗΣ

α) Η εξαρτημένη μεταβλητή είναι ποσοτική και οι ανεξάρτητες μεταβλητές είναι, κατά κανόνα, ποιοτικές.

β) Δεν κατασκευάζεται κάποια γραμμική σχέση μεταξύ της εξαρτημένης και των ανεξάρτητων μεταβλητών αλλά το ενδιαφέρον εστιάζεται στην ανάλυση των επιδράσεων μιας ή περισσότερων ανεξάρτητων μεταβλητών στην εξαρτημένη μεταβλητή. Αναζητούνται οι τιμές των ανεξάρτητων μεταβλητών που ασκούν τη σημαντικότερη επίδραση στην εξαρτημένη μεταβλητή.

# Μοντέλα ανάλυσης διακύμανσης



## Γραμμικά μοντέλα

### Μοντέλα παλινδρόμησης

1) Ποσοτικές μεταβλητές  
 $Y$   $X_1, \dots, X_p$

2) Προσπαθούμε να κατασκευάσουμε τη γραμμική σχέση που συνδέει την  $Y$  με τις  $X_1, \dots, X_p$ , να ελέγχουμε την εχθρότητα της σχέσης και να τη χρησιμοποιήσουμε προκειμένου να κάνουμε προβλέψεις.

### Μοντέλα ανάλυσης διακύμανσης

1) Θεωρούμε μια εξαρτημένη ποσοτική μεταβλητή  $Y$  και μια ή περισσότερες κατά κανόνα ποιοτικές μεταβλητές

2) Στα μοντέλα ανάλυσης διακύμανσης

δεν ενδιαφερόμαστε για την κατασκευή κάποιας συγκεκριμένης σχέσης, αλλά ενδιαφερόμαστε στο να ευτοπίσουμε, ποια ή ποιές από τις κατηγορίες των ποιοτικών μεταβλητών ασκούν τη σημαντικότερη επίδραση στην  $Y$ .

Καποιες φορές ποσοτικά δεδομένα μετατρέπουμε σε ποιοτικά π.χ ύψος  
150 και κάτω → κοντός  
150-165 → μέτριος  
165-180 → ψηλός  
180-2 → πολύ ψηλός

## Ορολογία-Συμβολισμός

Παράγοντας → ταυτόσημη με την ανεξάρτητη μεταβλητή.

Κάθε τιμή του παράγοντα (που είναι μια κατηγορία, αφού ο παράγοντας είναι ποιοτική μεταβλητή) ονομάζεται επίπεδο του παράγοντα.

Προβλήματα στα οποία υπεισέρχεται ένας παράγοντας, ονομάζονται προβλήματα ανάλυσης διακύμανσης κατά ένα παράγοντα.

Πρόβλημα στα οποία υπεισέρχονται δύο παράγοντες, ονομάζονται προβλήματα ανάλυσης διακύμανσης κατά δύο παράγοντες

Πρόβλημα με περισσότερους από δύο παράγοντες, ονομάζονται προβλήματα πειραματικών σχεδιασμών

Δοκιμασία: Κάθε συνδυασμός επιπέδων των παραχόντων (δύο ή περισσότερων)

$Y \equiv$  επίδοση (επίπεδα)  
 $\uparrow$  (παράγοντας)  
 $X \equiv$  δόση της ντομάτας  
χαμηλό  
μέτριο  
υψηλό  
δεν έχει νόημα να συνδέω κάτι που μετρείται με κάτι που χαρακτηρίζεται



**ΠΑΡΑΔΕΙΓΜΑ:** Ενδιαφέρη επίδοση  $Y$  μαθητών (ποσοτική). Ενδιαφέρεται να διερευνήσει πως το  $Y$  επηρεάζεται από το επίπεδο μόρφωσης Πατέρα (Ε.Μ.Π). Το επίπεδο μόρφωσης πατέρα είναι ποιοτική μεταβλητή η οποία κατηγοριοποιείται ως εξής: 1) Απόφοιτος Υποχρεωτικής (ΑΥ) 2) Λυκείου (Λ) 3) Πανεπιστημίου (Π) 4) Μεταπτυχιακό (Μ) 5) Διδακτορικό (Δ)

Πρόβλημα ανάλυσης διακύμανσης κατά ένα παράγοντα.

Παράγοντας  $\equiv$  ΕΜΠ, Επίπεδα παράγοντα  
Αν εισάγουμε το επίπεδο μόρφωσης μητέρας  $\rightarrow$  κατά δύο παράγοντες



1<sup>ος</sup> παράγοντας: ΕΜΠ.

2<sup>ος</sup> παράγοντας: ΕΜΜ.



Δοκιμασία:  $(AY, AX)$ ,  $(AY, \Lambda)$  κ.ο.κ.

(Περιγράφει πρόβλημα μέση τιμή, μια διακύμανση, μια κατανομή)

(I) Ανάλυση διακύμανσης κατά ένα παράγοντα

Θεωρώ μια ποσοτική  $Y$  και ένα παράγοντα ο οποίος απαντάται σε  $I$  επίπεδα.

Μορφή δεδομένων

Δεδομένα για τη  $Y$  (τ.δ. μεγέθους  $J_i$ )  
 δεδομένα από τον πληθυσμό του 1<sup>ου</sup> επιπέδου  
 δεδομένα από τον πληθυσμό του 2<sup>ου</sup> επιπέδου

Παράγοντας	Δεδομένα για την $Y$	Σύνολα	Μέσος όρος	
1 επίπεδο	$Y_{11}, Y_{12}, \dots, Y_{1J_1}$	$Y_1$	$\bar{Y}_1$	$N_1 (J_1, \sigma_1^2)$
2 επίπεδο	$Y_{21}, Y_{22}, \dots, Y_{2J_2}$	$Y_2$	$\bar{Y}_2$	$N_2 (J_2, \sigma_2^2)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$I$ -επίπεδο	$Y_{I1}, Y_{I2}, \dots, Y_{IJ_I}$	$Y_I$	$\bar{Y}_I$	$N_I (J_I, \sigma_I^2)$

όπου  $Y_{i\cdot} = \sum_{j=1}^{J_i} Y_{ij}$   $\bar{Y}_i = \frac{1}{J_i} \sum_{j=1}^{J_i} Y_{ij} = \frac{1}{J_i} Y_{i\cdot}$   $\bar{Y}_{..} = \frac{1}{N} \sum_{i=1}^I \sum_{j=1}^{J_i} Y_{ij}$   $N = \sum_{i=1}^I J_i$   $N = J_1 + J_2 + \dots + J_I$

συνολικό δειγματικό μέσο

Μοντέλο ανάλυσης διακύμανσης κατά ένα παράγοντα

$$Y_{ij} = \mu_i + \varepsilon_{ij}, \quad i=1, \dots, I, \quad j=1, \dots, J_i$$

Εκτιμητές ελαχίστων τετραγώνων

$$S^2 = \sum_{i=1}^I \sum_{j=1}^{J_i} \varepsilon_{ij}^2 = \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \mu_i)^2$$

$$\frac{\partial S}{\partial \mu_i} = 0 \Rightarrow \dots \Rightarrow \text{ΕΕΤ} : \hat{\mu}_i = \bar{Y}_i, \quad i=1, \dots, I$$

(Υποθέσεις για σφάλματα):

- |   |               |  |
|---|---------------|--|
| 1) $E(\varepsilon_{ij}) = 0, \quad i=1, \dots, I, j=1, \dots, J_i$      | } αντανάκλαση | ① $E(Y_{ij}) = \mu_i, \quad i=1, \dots, I$                   |
| 2) $\text{Var}(\varepsilon_{ij}) = \sigma^2, \quad j=1, \dots, J_i$     |               | ② $\text{Var}(Y_{ij}) = \sigma^2 / J_i, \quad i=1, \dots, I$ |
| 3) $\text{Cov}(\varepsilon_{ij}, \varepsilon_{kl}) = 0, \quad k \neq l$ |               | ③ $\text{Cov}(Y_{ij}, Y_{kl}) = 0$                           |
| 4) $\varepsilon_{ij} \sim N(0, \sigma^2)$                               |               | ④ $Y_{ij} \sim N(\mu_i, \sigma^2)$                           |

ΙΔΙΟΤΗΤΕΣ:

- 1)  $E(\hat{\mu}_i) = \mu_i, \quad i=1, \dots, I$
- 2)  $\text{Var}(\hat{\mu}_i) = \frac{\sigma^2}{J_i}, \quad i=1, \dots, I$

ΑΠΟΔΕΙΞΗ:

$$1) E(\hat{\mu}_i) = E(\bar{Y}_i) = E\left(\frac{1}{J_i} \sum_{j=1}^{J_i} Y_{ij}\right) = \frac{1}{J_i} \sum_{j=1}^{J_i} E(Y_{ij}) = \frac{1}{J_i} \sum_{j=1}^{J_i} \mu_i = \frac{1}{J_i} \cdot J_i \cdot \mu_i = \mu_i \quad (E(\varepsilon_{ij})=0 \text{ από υπόθεση})$$

Ισοδύναμο μοντέλο (μαθηματικά  $\mu_i = \mu + a_i$ )

$$Y_{ij} = \mu + a_i + \varepsilon_{ij}, \quad i=1, \dots, I, \quad j=1, \dots, J_i$$

Εκπληκτές ελαχίστων τετραγώνων

$$S^2 = \sum_{i=1}^I \sum_{j=1}^{J_i} \varepsilon_{ij}^2 = \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \mu - a_i)^2$$



$$\left. \begin{array}{l} \frac{\partial S'}{\partial \mu} = 0 \\ \frac{\partial S'}{\partial a_i} = 0 \end{array} \right\} \Rightarrow \left. \begin{array}{l} \sum_{i=1}^I a_i J_i + \mu \sum_{i=1}^I J_i = \sum_{i=1}^I \sum_{j=1}^{J_i} Y_{ij} \\ a_{ij} J_i + \mu J_i = \sum_{j=1}^{J_i} Y_{ij}, \quad i=1, \dots, I \end{array} \right\} \begin{array}{l} \text{πρώτη εξίσωση} \\ \text{καυονικές εξισώσεις} \end{array}$$

Το σύστημα καυονικών εξισώσεων δεν οδηγεί σε καυονική λύση, αφού αν αθροίσω την 2<sup>η</sup> εξίσωση προκύπτει η 1<sup>η</sup> ως προς  $a_i$  (γραμ. εξαρτημένα)

Για να πετύχω μοναδική λύση θεωρώ πλευριική συνθήκη. Το ερώτημα είναι ποια θα είναι η πλευριική συνθήκη;

Έχουν προταθεί διάφορες. Από τις διάφορες καλύτερη φυσική ερμηνεία έχει εκείνη που οδηγεί σε αποδεικτούς διασθητικά επιτημητές

Π.χ. διασθητικά αποδεικτός επιτημητές για το  $\mu$  είναι  $\hat{\mu} = \bar{Y}_{..}$   $\rightarrow$  γενικό δείγμα μέσος (αφού το  $\mu$  σχετίζεται με όλες τις παρατηρήσεις δηλαδή διασθητικά του δείγμα)

Για να πετύχω ως επιτημητή του γενικού δείγματός μέσου αρκεί να υποθέσω ότι  $\sum_{i=1}^I a_i J_i = 0$  (από την πρώτη εξίσωση)

οπότε θεωρώ πλευριική συνθήκη:  $\sum_{i=1}^I a_i J_i = 0$ .

Υπό αυτή τη πλευριική συνθήκη οι ΕΕΤ είναι  $\hat{\mu} = \bar{Y}_{..}$   
 $\hat{a}_i = \bar{Y}_{i.} - \bar{Y}_{..}$

Πίνακας Ανάδια του μοντέλου ανάλυσης διακύμανσης κατά ένα παράγοντα

- ολική μεταβλητότητα = δειγματική διακύμανση χωρίς τα μερέδια δείγματος
- δειγματική διακύμανση =  $\frac{1}{N-1} \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{..})^2$   
 $\hookrightarrow$  χωρίς μερέδια δείγματος

ΘΕΩΡΗΜΑ:  $SS_{tot} = \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{..})^2 = \sum_{i=1}^I J_i (\bar{Y}_i - \bar{Y}_{..})^2 + \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_i)^2$

$\underbrace{\hspace{10em}}_{SS_{tr}} \qquad \qquad \qquad \underbrace{\hspace{10em}}_{SS_{res}}$

ΑΠΟΔΕΙΞΗ:  $\sum \sum (Y_{ij} - \bar{Y}_{..} + \bar{Y}_i - \bar{Y}_i)^2 = \dots$  ηρά Γεις

ΠΙΝΑΚΙΣ ΑΝΑΔΙΑ

πηγή μεταβλητότητας μοτέλο	$SS_{tr} = \sum_{i=1}^I J_i (\bar{Y}_i - \bar{Y}_{..})^2$	ΒΕ I-1	MS $MS_{tr} = \frac{SS_{tr}}{I-1}$	F-οηηιο $F = \frac{MS_{tr}}{MS_{res}}$
υπόηοηα	$SS_{res} = \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_i)^2$	N-I	$MS_{res} = \frac{SS_{res}}{N-I}$	$MS_{res}$
οηηη	$SS_{tot} = \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{..})^2$	N-1		