

Τ.Ε.Ι. ΜΕΣΣΟΓΕΩΝ
ΣΧΟΛΗ ΔΙΟΙΚΗΣΗΣ & ΟΙΚΟΝΟΜΙΑΣ
ΤΜΗΜΑ ΜΑΡΚΕΤΙΝΓΚ



ΕΙΣΗΓΗΤΗΣ: Ι. ΚΟΥΓΙΑΣ
ΣΠΟΥΔΑΣΤΕΣ: ΚΥΡΙΑΚΗ ΝΙΚΟΛΕΤΟΥ
ΒΑΣΙΛΕΙΟΣ ΝΑΝΟΣ

ΣΥΣΧΕΤΙΣΗ-ΠΑΛΙΝΔΡΟΜΗΣΗ
ΣΤΟΝ ΕΠΙΧΕΙΡΗΣΙΑΚΟ ΚΑΙ
ΟΙΚΟΝΟΜΙΚΟ ΤΟΜΕΑ

ΑΘΗΝΑ-ΙΟΥΝΙΟΣ 1998

ΕΛΛΗΝΙΚΗ ΔΗΜΟΚΡΑΤΙΑ
ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ & ΘΡΗΣΚΕΥΜΑΤΩΝ
ΙΝΣΤΙΤΟΥΤΟ ΤΕΧΝΟΛΟΓΙΑΣ ΥΠΟΛΟΓΙΣΤΩΝ ΚΑΙ ΕΚΔΟΣΕΩΝ ΔΙΟΙΚΗΤΙΚΗΣ ΥΠΗΡΕΣΙΑΣ (ΙΤΥΥ Δ.Ι.Κ.Ε.Υ.)
7967

**Τ.Ε.Ι ΜΕΣΟΛΟΓΓΙΟΥ
ΣΧΟΛΗ ΔΙΟΙΚΗΣΗΣ & ΟΙΚΟΝΟΜΙΑΣ
ΤΜΗΜΑ ΛΟΓΙΣΤΙΚΗΣ**

**ΕΙΣΗΓΗΤΗΣ ΞΚΟΥΓΙΑΣ
ΣΠΟΥΔΑΣΤΕΣ : ΚΥΡΙΑΚΗ ΝΙΚΟΛΕΤΟΥ
ΒΑΣΙΛΕΙΟΣ ΝΑΝΟΣ**

**ΣΥΣΧΕΤΙΣΗ-ΠΑΛΙΝΔΡΟΜΗΣΗ
ΣΤΟΝ ΕΠΙΧΕΙΡΗΣΙΑΚΟ ΚΑΙ
ΟΙΚΟΝΟΜΙΚΟ ΤΟΜΕΑ**

ΑΘΗΝΑ - ΙΟΥΝΙΟΣ 1998

ΚΕΦΑΛΑΙΟ ΠΡΩΤΟ

ΕΙΣΑΓΩΓΙΚΕΣ ΕΝΝΟΙΕΣ

Εισαγωγή

Για να πραγματοποιήσουν πωλήσεις οι επιχειρήσεις χρειάζονται περιουσιακά στοιχεία. Αν οι πωλήσεις πρόκειται να αυξηθούν, θα απαιτηθούν πρόσθετα περιουσιακά στοιχεία. Οι αναπτυσσόμενες επιχειρήσεις χρειάζονται νέες επενδύσεις-άμεσες επενδύσεις σε κυκλοφορούντα στοιχεία- και καθώς εξαντλείται η δυναμικότητα παραγωγής και νέες επενδύσεις σε πάγιες εγκαταστάσεις. Οι πρόσθετες επενδύσεις πρέπει να χρηματοδοτηθούν και η νέα χρηματοδότηση δημιουργεί δεσμεύσεις και υποχρεώσεις για την εξυπηρέτηση των κεφαλαίων που θα αποκτηθούν. Μια αναπτυσσόμενη κερδοφόρα επιχείρηση πιθανόν να χρειαστεί πρόσθετα χρήματα για επενδύσεις σε αποθέματα και πάγιες εγκαταστάσεις, καθώς και για την δημιουργία εισπρακτέων λογαριασμών. Μια τέτοια επιχείρηση είναι πιθανό να οδηγηθεί σε ταμιακό πρόβλημα, εκτός αν η απαιτούμενη ανάγκη για χρηματοδότηση επισημανθεί έγκαιρα.

Μέ την ανάπτυξη πολύπλοκων υποδειγμάτων για την επιχείρηση που κάνουν χρήση ηλεκτρονικών υπολογιστών, η πρόβλεψη των επενδυτικών αναγκών απαιτεί την εφαρμογή πολύ λεπτομερών διαδικασιών. Όλα πάντως τα υποδείγματα στηρίζονται στις βασικές μεθοδολογίες που θα παρουσιάσουμε εδώ, δηλαδή στην μέθοδο παλινδρόμηση-συσχέτιση.

1.1 Μοντέλο

Στα πλαίσια μιας επιστημονικής ανάλυσης συχνά μας ενδιαφέρει να προσδιορίσουμε τη σχέση που συνδέει μια μεταβλητή, έστω y , με μία ή περισσότερες άλλες, έστω τις x_1, x_2, \dots, x_k .

Στην περίπτωση αυτή λέμε ότι θέλουμε να προσδιορίσουμε ένα μοντέλο το οποίο να ερμηνεύει τη συμπεριφορά της y με βάση τη συμπεριφορά των x_j , $j=1, \dots, k$. Γενικά, το μοντέλο ή υπόδειγμα είναι ένας

μαθηματικός τύπος ο οποίος παρουσιάζει απλουστευτικά σχέσεις που στην πραγματικότητα μπορεί να είναι πολύπλοκες.

Όταν η ανάλυση είναι θεωρητική το μοντέλο μπορεί να είναι **προσδιοριστικό** (deterministic) δηλαδή να ορίζει πλήρως την y αν είναι γνωστή η x . Έτσι π.χ το μοντέλο κατανάλωσης

$$y=b_0+b_1x$$

προσδιορίζει μια τιμή για την κατανάλωση y ενός ατόμου όταν το εισόδημά του είναι x . Τα b_0 και b_1 είναι οι παράμετροι του μοντέλου και μπορεί να είναι οποιοδήποτε σταθερές.

1.2 Ανάλυση Παλινδρόμησης

Το καλύτερο που μπορούμε να κάνουμε για ένα στατιστικό μοντέλο είναι να το εκτιμήσουμε. Αυτό, γενικά, σημαίνει να παρατηρήσουμε τη συμπεριφορά της y και συγχρόνως, των μεταβλητών x_1, \dots, x_k και με τη στατιστική ανάλυση αυτών των παρατηρήσεων, να επιλέξουμε ένα μοντέλο το οποίο θα μπορούσε να τις είχε παράγει.

Όλες οι τεχνικές που χρησιμοποιούνται για το σκοπό αυτό αναφέρονται ως **Ανάλυση Παλινδρόμησης** (Regression Analysis) ή **Μεθοδολογία Προσαρμογής Καμπυλών** (Curve Fitting Methodology) ή **Εκτίμηση ενός Μοντέλου Παλινδρόμησης** (Estimation of a Regression Model).

Γενικά, στην Ανάλυση Παλινδρόμησης θεωρούμε ότι η τυχαία μεταβλητή y μπορεί να εκφραστεί ως άθροισμα δύο μερών:

μιας μέσης σχέσης ξ η οποία είναι συνάρτηση των μεταβλητών x_j , $j=1, \dots, k$ με παραμέτρους b_1, b_2, \dots, b_k και ενός τυχαίου όρου ε . Υποθέτουμε δηλαδή ότι ισχύει:

$$y=\xi+\varepsilon$$

και

$$\xi=\varepsilon(y)=f(x_1, x_2, \dots, x_k / b_1, b_2, \dots, b_k)$$

Η μεταβλητή y ονομάζεται **εξαρτημένη** (dependent) ή **ενδογενής** (endogenous) και οι x_j , $j=1, \dots, k$ ονομάζονται **ανεξάρτητες** (independent) ή **εξωγενείς** (exogenous) ή **ερμηνευτικές** (explanatory) ή **μεταβλητές ελέγχου** (control variables).

Οι ερμηνευτικές μεταβλητές θεωρούνται κατ'αρχήν ότι μπορούν να ελεγχθούν (δηλαδή ότι ο ερευνητής μπορεί να επιλέξει τα επίπεδα των x_j στα οποία θα παρατηρήσει την y) και μετριοούνται χωρίς σφάλμα. Όμως η ανάλυση παλινδρόμησης μπορεί να επεκταθεί και σε τυχαίες

ερμηνευτικές μεταβλητές, αρκεί να ισχύουν ορισμένες υποθέσεις, όπως ειδικότερα θα δούμε στη συνέχεια.

Η ξ ονομάζεται **συνάρτηση παλινδρόμησης** της y επί των χ_1, \dots, χ_k . Όταν οι ερμηνευτικές μεταβλητές είναι περισσότερες από μία ($k > 1$) τότε ονομάζεται **πολυμεταβλητή** ή **πολλαπλή** διαφορετικά ονομάζεται **απλή**.

1.3 Γραμμική Συνάρτηση Παλινδρόμησης

Εκτός από τις περιπτώσεις που υπάρχει μια καλά αναπτυγμένη σχετική θεωρία, ο μαθηματικός τύπος της συνάρτησης παλινδρόμησης είναι άγνωστος. Σχεδόν πάντα αρχίζουμε την ανάλυση υιοθετώντας την πιο απλή μορφή που είναι η γραμμική δηλαδή η εξής:

$$\xi = b_0 + b_1 \cdot \chi_1 + b_2 \cdot \chi_2 + \dots + b_k \cdot \chi_k$$

όπου b_0, b_1, \dots, b_k είναι άγνωστοι παράμετροι. Η μορφή αυτή είναι γενική και περιλαμβάνει οποιαδήποτε γραμμική ως προς τις παραμέτρους συνάρτηση. Αυτό σημαίνει οποιαδήποτε συνάρτηση που έχει το ακόλουθο χαρακτηριστικό: καμιά από τις ερμηνευτικές μεταβλητές δεν εξαρτάται από κάποια άγνωστη παράμετρο και καμιά άγνωστη παράμετρος δεν εξαρτάται από οποιαδήποτε άλλη παράμετρο. Έτσι π.χ οι συναρτήσεις

$$\xi = b_0 + b_1 \cdot \chi_1 + b_2^{b^3} \cdot \chi_2$$

$$\text{ή} \quad \xi = b_0 + b_1 \cdot \chi_1 + b_2 \cdot \chi_2^{b^3}$$

δεν είναι γραμμικές ενώ οι

$$\xi = b_0 + b_1 \cdot \chi_1 + b_2 \cdot \chi_1^2$$

$$\text{ή} \quad \xi = b_0 + b_1 \cdot \chi_1^{1/2} + b_2 \cdot \chi_1 \cdot \chi_2$$

είναι γραμμικές.

ΚΕΦΑΛΑΙΟ ΔΕΥΤΕΡΟ

ΑΠΛΗ ΓΡΑΜΜΙΚΗ ΠΑΛΙΝΔΡΟΜΗΣΗ

2.1 Το απλό γραμμικό μοντέλο παλινδρόμησης

Το απλό γραμμικό μοντέλο παλινδρόμησης είναι αυτό στο οποίο η μέση τιμή της εξαρτημένης μεταβλητής y είναι γραμμική συνάρτηση της ανεξάρτητης μεταβλητής x . Έτσι, αν y_1, \dots, y_n οι παρατηρήσεις της y στα επίπεδα, αντίστοιχα, x_1, \dots, x_n της x τότε ισχύει:

$$y_i = b_0 + b_1 \cdot x_i + \varepsilon_i, \quad i=1, \dots, n \quad (2.1.1)$$

όπου ε_i ο όρος σφάλματος. Για να εκτιμήσουμε τις παραμέτρους $b_0 + b_1$ με τη μέθοδο των ελαχίστων τετραγώνων κάνουμε τις ακόλουθες υποθέσεις:

Ι. Για την ερμηνευτική μεταβλητή

ι. Οι τιμές $x_i, i=1, \dots, n$ είναι σταθερές, όχι όλες ίσες μεταξύ τους που επιλέγονται από τον ερευνητή και μετριοούνται χωρίς σφάλμα.

Αυτό σημαίνει ότι η x μπορεί να προκαλεί μεταβολές στην y αλλά όχι το αντίστροφο. Καθώς η y είναι τυχαία μεταβλητή, αν προκαλούσε μεταβολές στην x , θα ήταν κι αυτή τυχαία. Η συνθήκη αυτή είναι πολύ περιοριστική, ιδιαίτερα για τις οικονομικές μεταβλητές οι οποίες είναι σχεδόν πάντα τυχαίες. Αποδεικνύεται όμως ότι η μέθοδος των ελαχίστων τετραγώνων ορίζει άριστους εκτιμητές και όταν ισχύει η λιγότερο περιοριστική συνθήκη :

ι'. Οι x_i είναι παρατηρήσεις της τυχαίας μεταβλητής x η οποία είναι ασυσχέτιστη με κάθε διαταρακτικό όρο $\varepsilon_i, i=1, \dots, n$.

Στην περίπτωση αυτή όλες οι εκτιμήσεις και οι αναγωγές γίνονται υπο συνθήκη των παρατηρήσεων x_i .

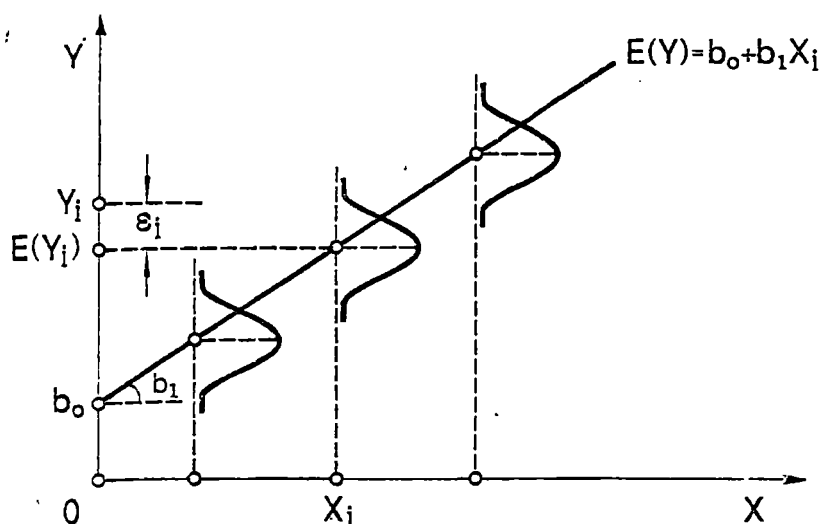
II. Για τον όρο σφάλματος

Ισχύει :

- i. $E(\varepsilon_i) = 0 \Leftrightarrow E(Y_i) = b_0 + b_1 \cdot X_i \quad i = 1, \dots, n$
- ii. $\text{Var}(\varepsilon_i) = \sigma^2 \Leftrightarrow \text{Var}(Y_i) = \sigma^2 \quad i = 1, \dots, n$
- iii. $\text{Cov}(\varepsilon_i - \varepsilon_j) = 0 \Leftrightarrow \text{Cov}(Y_i - Y_j) = 0 \quad \forall i \neq j$
- iv. $\varepsilon_i \sim N(0, \sigma^2) \Leftrightarrow Y_i \sim N(b_0 + b_1 X_i, \sigma^2) \quad i = 1, \dots, n$

Οι συνθήκες (i)-(iv) μαζί με την (2.1.1) ορίζουν ένα μοντέλο δειγματοληψίας στο οποίο οι παρατηρήσεις y_i είναι ανεξάρτητες, έχουν την ίδια κατανομή, ίδια διακύμανση και μέση τιμή που βρίσκεται επάνω σε μια ευθεία με σταθερό όρο b_0 και κλίση b_1 . Η ευθεία αυτή ονομάζεται **θεωρητική ή ευθεία παλινδρόμησης πληθυσμού**.

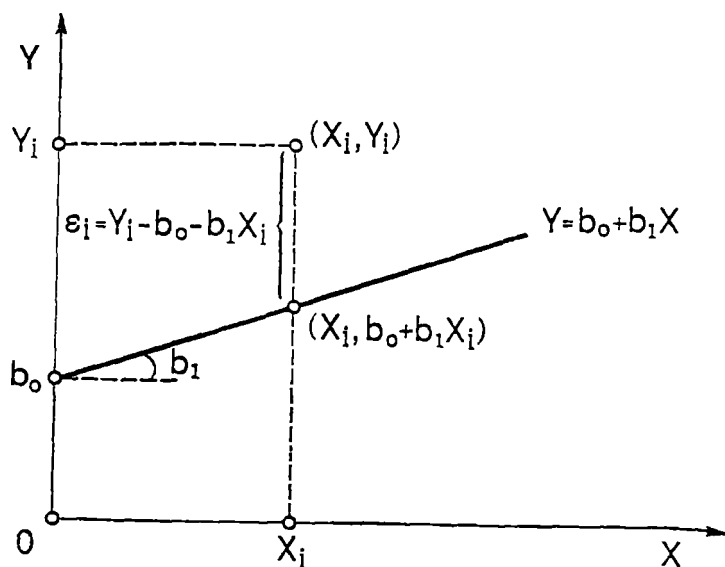
Η γραφική παράσταση ενός τέτοιου μοντέλου δίνεται στο ακόλουθο σχήμα :



2.2 Η μέθοδος των ελαχίστων τετραγώνων

Ας θεωρήσουμε το ακόλουθο σχήμα στο οποίο δίνεται η ευθεία L με εξίσωση $y=b_0+b_1 \cdot x$ και ένα σημείο (x_i, y_i) . Η απόσταση του σημείου αυτού από την L είναι η

$$\varepsilon_i = Y_i - (b_0 + b_1 \cdot X_i)$$



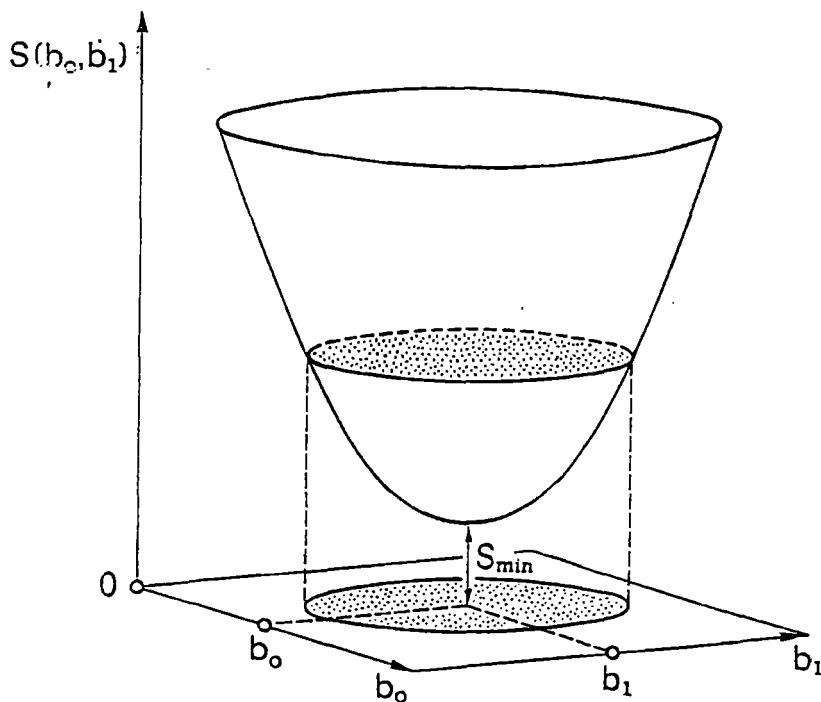
Παρατηρούμε ότι αν θεωρήσουμε μια άλλη ευθεία L' με εξίσωση $y=b'_0+b'_1 \cdot x$ τότε η απόσταση του σημείου (x_i, y_i) από την ευθεία L' θα είναι η

$$\varepsilon'_i = Y_i - (b'_0 + b'_1 \cdot X_i)$$

Από όλες τις δυνατές ευθείες στο επίπεδο xy η ευθεία ελαχίστων τετραγώνων είναι εκείνη που ελαχιστοποιεί το άθροισμα των τετραγώνων των ή των αποστάσεων $\varepsilon'_i = y_i - (b_0 + b_1 \cdot x_i)$ δηλαδή αυτή για την οποία η συνάρτηση

$$S = S(b_0, b_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - b_0 - b_1 \cdot x_i)^2 \quad (2.2.1)$$

έχει ελάχιστο. Η γραφική παράσταση μιας τέτοιας συνάρτησης είναι μια παραβολοειδής επιφάνεια σχήματος μπωλ όπως αυτή που δίνεται στο σχήμα.



Εστω \hat{b}_0, \hat{b}_1 οι τιμές των b_0, b_1 αντίστοιχα, οι οποίες όταν αντικατασταθούν στην εξίσωση (2.2.1) δίνουν την ελάχιστη δυνατή τιμή για την S . Για να προσδιορίσουμε τις \hat{b}_0 και \hat{b}_1 υπολογίζουμε τις μερικές παραγώγους της S

ως προς b_0 και b_1 και εξισώνουμε με μηδέν. Έχουμε (οι δείκτες στα Σ παραλείπονται) :

$$\frac{\partial S}{\partial b_0} = -2 \sum (y_i - b_0 - b_1 \cdot x_i)$$

$$\frac{\partial S}{\partial b_1} = -2 \sum x_i (y_i - b_0 - b_1 \cdot x_i)$$

Οι τιμές των \hat{b}_0, \hat{b}_1 προσδιορίζονται από τη λύση του συστήματος εξισώσεων:

$$\left. \begin{aligned} \sum (y_i - \hat{b}_0 - \hat{b}_1 \cdot x_i) &= 0 \\ \sum x_i (y_i - \hat{b}_0 - \hat{b}_1 \cdot x_i) &= 0 \end{aligned} \right\} \Leftrightarrow$$

$$\left. \begin{aligned} \sum y_i - n \hat{b}_0 - \hat{b}_1 \cdot \sum x_i &= 0 \\ \sum x_i \cdot y_i - \hat{b}_0 \sum x_i - \hat{b}_1 \sum x_i^2 &= 0 \end{aligned} \right\} \Leftrightarrow$$

$$\begin{aligned} n \hat{b}_0 + \hat{b}_1 \sum x_i &= \sum y_i & (2.2.2) \\ \hat{b}_0 \sum x_i + \hat{b}_1 \sum x_i^2 &= \sum x_i \cdot y_i \end{aligned}$$

Οι (2.2.2) ονομάζονται **κανονικές εξισώσεις** και είναι ένα γραμμικό σύστημα με αγνώστους τα b_0 και b_1 . Η λύση του συστήματος αυτού είναι η :

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \cdot \bar{x} \quad (2.2.3)$$

$$\hat{b}_1 = \frac{\sum \bar{x}_i \cdot \bar{y}_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \cdot \bar{x}^2} \quad (2.2.4)$$

Χρήσιμη είναι η ισοδύναμη έκφραση για το b_1

$$\hat{b}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad (2.2.4)'$$

Η ευθεία ελαχίστων τετραγώνων είναι η

$$\hat{y} = \hat{b}_0 + \hat{b}_1 \cdot x$$

και ονομάζεται δειγματική ή εκτιμηθείσα ευθεία παλινδρόμησης σε αντιδιαστολή με τη θεωρητική ή ευθεία παλινδρόμησης στον πληθυσμό.

2.3 Τα κατάλοιπα ή σφάλματα εκτίμησης

Αν στο εκτιμηθέν μοντέλο αντικαταστήσουμε όπου x τις τιμές x_i , $i=1, \dots, n$ των δεδομένων μας παίρνουμε τις εκτιμήσεις

$$\hat{y}_i = \hat{b}_0 + \hat{b}_1 \cdot x_i, \quad i=1, \dots, n$$

Οι διαφορές των εκτιμήσεων \hat{y}_i από τις παρατηρήσεις y_i συμβολίζονται με $\hat{\varepsilon}_i$ και ονομάζονται **σφάλματα εκτίμησης ή κατάλοιπα**. Δηλαδή έχουμε :

$$\begin{aligned} \hat{\varepsilon}_i &= Y_i - \hat{Y}_i \\ &= y_i - \hat{b}_0 - \hat{b}_1 \cdot X_i, \quad i=1, \dots, n \end{aligned}$$

Το κατάλοιπο $\hat{\varepsilon}_i$ είναι το μέρος της παρατήρησης Y_i το οποίο μένει ανερμήνευτο από το εκτιμηθέν μοντέλο και αποτελεί μια εκτίμηση του διαταρακτικού όρου ε_i . Επομένως, η συμπεριφορά των καταλοίπων $\hat{\varepsilon}_i$, $i=1, \dots, n$ μας δίνει χρήσιμες πληροφορίες για τη συμπεριφορά των διαταρακτικών όρων.

Παράδειγμα :

Οι εκτιμήσεις Y_i και τα κατάλοιπα $\hat{\varepsilon}_i$ του παραδείγματος της αξίας των πωλήσεων δίνονται στον παρακάτω πίνακα :

Y_i	X_i	$\hat{Y}_i =$	$\varepsilon_i = Y_i - \hat{Y}_i$
1100	170,0	1168,86	-68,86
1180	172,0	1407,22	-187,22
1124	169,5	1574,64	25,36
1350	190,0	1812,43	-52,43
1220	212,0	1835,13	-185,13
1590	228,0	1180,21	0,22
1340	202,0	1498,02	91,98
1570	235,0	1567,26	-17,26
1660	241,5	1694,95	-64,95
1550	240,2	1746,60	73,41
1820	283,3	1166,03	73,97
2040	300,4	1350,47	-10,47
1760	283,4	1811,86	8,14
1630	262,7	1916,28	3,72
1920	301,7	1282,37	67,63
1720	253,2	1537,75	32,25
1650	287,4	1908,90	131,10
1820	271,8	1641,04	78,96

Ιδιότητες των καταλοίπων

Τα κατάλοιπα $\hat{\varepsilon}_i$, $i=1, \dots, n$ της εκτίμησης ενός μοντέλου έχουν τις ακόλουθες ιδιότητες :

1^η ιδιότητα : Το άθροισμα τους ισούται με μηδέν, δηλαδή (ο συντελεστής Σ είναι για i από 1 ως n) :

$$\boxed{\Sigma \hat{\varepsilon}_i = 0}$$

Πράγματι έχουμε :

$$\hat{Y}_i = \hat{b}_0 + \hat{b}_1 X_i = \bar{Y} - \hat{b}_1 \cdot \bar{X} + \hat{b}_1 \cdot X_i = \bar{Y} + \hat{b}_1 (X_i - \bar{X})$$

και

$$\begin{aligned} \hat{\varepsilon}_i &= Y_i - \hat{Y}_i = Y_i - \bar{Y} - \hat{b}_1 (X_i - \bar{X}) \Leftrightarrow \\ \sum \hat{\varepsilon}_i &= \sum (Y_i - \bar{Y}) - \hat{b}_1 \sum (X_i - \bar{X}) = 0 \end{aligned}$$

2^η ιδιότητα : Το άθροισμα των γινομένων $X_i \cdot \hat{\varepsilon}_i$ ισούται με μηδέν δηλαδή

$$\boxed{\sum X_i \cdot \hat{\varepsilon}_i = 0}$$

Πράγματι έχουμε :

$$\begin{aligned} \sum X_i \cdot \hat{\varepsilon}_i &= \sum X_i (Y_i - \hat{b}_0 - \hat{b}_1 \cdot X_i) \\ &= \sum X_i Y_i - \hat{b}_0 \sum X_i - \hat{b}_1 \sum X_i^2 \\ &= 0 \text{ από τη δεύτερη κανονική εξίσωση.} \end{aligned}$$

Σημείωση : Η 1^η και η 2^η ιδιότητα χαρακτηρίζουν την ευθεία ελαχίστων τετραγώνων αρκεί μια ευθεία να έχει τις δύο αυτές ιδιότητες για να είναι ευθεία ελαχίστων τετραγώνων. Τα κατάλοιπα όμως έχουν και την ακόλουθη σημαντική ιδιότητα :

3^η ιδιότητα : Το άθροισμα των γινομένων $\hat{Y}_i \cdot \hat{\varepsilon}_i$ ισούται με μηδέν, δηλαδή

$$\boxed{\sum \hat{Y}_i \cdot \hat{\varepsilon}_i = 0}$$

Πράγματι έχουμε :

$$\begin{aligned} \hat{Y}_i \cdot \hat{\varepsilon}_i &= (\hat{b}_0 + \hat{b}_1 \cdot X_i) \cdot (Y_i - \hat{b}_0 - \hat{b}_1 \cdot X_i) \\ &= \hat{b}_0 (Y_i - \hat{b}_0 - \hat{b}_1 \cdot X_i) + \hat{b}_1 \cdot X_i (Y_i - \hat{b}_0 - \hat{b}_1 \cdot X_i) \end{aligned}$$

Αθροίζοντας για όλες τις τιμές του i παίρνουμε :

$$\begin{aligned} \sum \hat{Y}_i \cdot \hat{\varepsilon}_i &= \hat{b}_0 (\sum Y_i - n \hat{b}_0 - \hat{b}_1 \sum X_i) + \hat{b}_1 (\sum X_i Y_i - \hat{b}_0 \sum X_i - \hat{b}_1 \sum X_i^2) \\ &= 0 \end{aligned}$$

αφού από τις δύο κανονικές εξισώσεις οι παραστάσεις στις δύο παρενθέσεις μηδενίζονται

2.4 Μέση τιμή και διακύμανση των εκτιμητών ελαχίστων τετραγώνων.

Οι εκτιμήσεις των παραμέτρων b_0 και b_1 της θεωρητικής ευθείας παλινδρόμησης βασίζονται σε ένα τυχαίο δείγμα ή ζευγών παρατηρήσεων (x_i, Y_i) . Οι Y_i είναι τυχαίες μεταβλητές. Είναι προφανές επομένως, ότι από μια άλλη διαδικασία παρατήρησης στην οποία οι τιμές x_i παραμένουν ίδιες θα είχαμε, με πιθανότητα ίση με 1, διαφορετικές τιμές για τις Y_i και διαφορετικές εκτιμήσεις \hat{b}_0 και \hat{b}_1 . Επομένως οι \hat{b}_0 και \hat{b}_1 που ορίζονται από τους τύπους (2.2.3) και (2.2.4) ως συναρτήσεις των τυχαίων μεταβλητών Y_i , είναι τυχαίες μεταβλητές και θα αναφέρονται ως εκτιμητές των b_0 και b_1 . Αντίστοιχα, οι μεμονωμένες τιμές που υπολογίζουμε από ορισμένο δείγμα θα αναφέρονται ως εκτιμήσεις.

Αποδεικνύεται ότι οι εκτιμητές ελαχίστων τετραγώνων \hat{b}_0 και \hat{b}_1 είναι αμερόληπτοι εκτιμητές των b_0 και b_1 , δηλαδή ισχύει:

$$E(\hat{b}_0) = b_0 \quad \text{και} \quad E(\hat{b}_1) = b_1 \quad (2.4.1)$$

και οι διακυμάνσεις τους είναι, αντίστοιχα, ίσες με (σ^2) (ο συντελεστής σ^2 είναι για τιμές από 1 ως η):

$$\text{Var}(\hat{b}_0) = \sigma_{\hat{b}_0}^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum (x_i - \bar{X})^2} \right] \quad (2.4.2)$$

$$\text{Var}(\hat{b}_1) = \sigma_{\hat{b}_1}^2 = \frac{\sigma^2}{\sum (x_i - \bar{X})^2} \quad (2.4.3)$$

Για να χρησιμοποιήσουμε τους τύπους (2.4.2) και (2.4.3) θα πρέπει να εκτιμήσουμε την σ^2 . Αποδεικνύεται ότι, όταν ισχύουν οι υποθέσεις του γραμμικού μοντέλου ο τύπος:

$$S^2 = \frac{\sum \hat{\epsilon}_i^2}{n-2} = \frac{\sum (Y_i - \hat{b}_0 - \hat{b}_1 \cdot X_i)^2}{n-2} \quad (2.4.4)$$

ορίζει ένα αμερόληπτο εκτιμητή της κοινής διακύμανσης των διαταρακτικών όρων δηλαδή της σ^2 . Αν αναπτύξουμε τον αριθμητή της (2.4.4) προκύπτει ο ισοδύναμος τύπος:

$$S^2 = \frac{\sum Y_i^2 - \hat{b}_0 \sum Y_i - \hat{b}_1 \sum Y_i \cdot X_i}{n-2} \quad (2.4.5)$$

που είναι υπολογιστικά ευκολότερος αφού όλοι οι όροι του έχουν υπολογιστεί προηγουμένως. Το s^2 μέσω τετραγωνικό σφάλμα. Η τετραγωνική του ρίζα s ονομάζεται τυπικό σφάλμα της εκτίμησης και μετρά τη διασπορά των παρατηρήσεων Y_i γύρω από την ευθεία ελαχίστων τετραγώνων σε μονάδες της Y .

Αν στις (2.4.2) και (2.4.3) αντικαταστήσουμε την σ^2 με την εκτίμηση s^2 παίρνουμε τους αμερόληπτους εκτιμητές των $\text{Var}(\hat{b}_0)$ και $\text{Var}(\hat{b}_1)$ που συμβολίζονται με $s_{\hat{b}_0}^2$ και $s_{\hat{b}_1}^2$, αντίστοιχα. Οι τετραγωνικές ρίζες αυτών των εκτιμητών είναι τα τυπικά σφάλματα, αντίστοιχα, των εκτιμητών \hat{b}_0 και \hat{b}_1 .

Επομένως, το τυπικό σφάλμα του \hat{b}_0 είναι το εξής:

$$S_{\hat{b}_0} = s \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2}}$$

ενώ το τυπικό σφάλμα του \hat{b}_1 το εξής:

$$S_{\hat{b}_1} = s \sqrt{\frac{1}{\sum (X_i - \bar{X})^2}}$$

2.5 Πόσο καλοί είναι οι εκτιμητές ελαχίστων τετραγώνων

Όταν ισχύουν οι υποθέσεις του γραμμικού μοντέλου παλινδρόμησης οι εκτιμητές ελαχίστων τετραγώνων είναι:

Γραμμικοί: δηλαδή είναι γραμμικές συναρτήσεις των παρατηρήσεων Y_i . Πράγματι έχουμε:

$$\hat{b}_1 = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} (y_i - \bar{y})$$
$$\hat{b}_0 = \bar{y} - \hat{b}_1 \cdot \bar{x}$$

Επειδή τα X_i είναι σταθεροί αριθμοί οι όροι $\sum (x_i - \bar{x}) / \sum (x_i - \bar{x})^2$ και \bar{x} είναι επίσης σταθερές. Επομένως στους τύπους των \hat{b}_1 και \hat{b}_0 περιέχονται μόνο οι όροι οι οποίοι είναι γραμμικοί ως προς τις Y_i . Η σημασία της ιδιότητας αυτής οφείλεται στο ότι, συνήθως, η κατανομή μιας γραμμικής συνάρτησης τυχαίων μεταβλητών προσδιορίζεται εύκολα.

Αμερόληπτοι: δηλαδή $E(\hat{b}_0) = b_0$ και $E(\hat{b}_1) = b_1$.

Άριστοι: δηλαδή από όλους τους γραμμικούς αμερόληπτους εκτιμητές έχουν τη μικρότερη διακύμανση.

Έχει καθιερωθεί οι ιδιότητες αυτές να αναφέρονται συνοπτικά ως BLUE (=Best Linear Unbiased Estimators). Οι Άριστοι, Γραμμικοί, Αμερόληπτοι εκτιμητές δεν μας εξασφαλίζουν ότι οι εκτιμήσεις που θα πάρουμε από ένα μόνο δείγμα θα είναι κοντά στις πραγματικές τιμές των παραμέτρων, παρά μόνον πιθανοκρατικά. Έτσι, οι εκτιμήσεις που θα προκύψουν από τη μέθοδο των ελαχίστων τετραγώνων θα είναι, με μεγαλύτερη πιθανότητα, κοντά στις πραγματικές τιμές b_0 και b_1 από τις εκτιμήσεις που θα προκύψουν από οποιαδήποτε άλλη μέθοδο που χρησιμοποιεί γραμμικές συναρτήσεις των y_i .

Το αποτέλεσμα αυτό είναι γνωστό ως θεώρημα των Gauss-Markov. Σημειώνεται επίσης ότι επειδή υποθέσεις του γραμμικού μοντέλου είναι αποφασιστικές για την ισχύ του αποτελέσματος αυτού για αυτό ονομάζονται και συνθήκες Gauss-Markov.

2.6 Διαστήματα εμπιστοσύνης και έλεγχος των υποθέσεων.

Αποδεικνύεται ότι, όταν οι διαταρακτικοί όροι είναι ανεξάρτητες και κανονικές τυχαίες μεταβλητές με μέση τιμή μηδέν και διακύμανση σ^2 , δηλαδή συνοπτικά, όταν

$$\varepsilon_i \sim \text{IDN}(0, \sigma^2), \quad i = 1, \dots, n$$

τότε, οι τυχαίες μεταβλητές

$$\frac{\hat{b}_0 - b_0}{S \hat{b}_0} \quad \text{και} \quad \frac{\hat{b}_1 - b_1}{S \hat{b}_1}$$

ακολουθούν την κατανομή t-student με $\nu = n - 2$ βαθμούς ελευθερίας. Με βάση το αποτέλεσμα αυτό, μπορούμε να κατασκευάσουμε διαστήματα εμπιστοσύνης και να κάνουμε έλεγχο των υποθέσεων για τις παραμέτρους b_0 και b_1 της θεωρητικής ευθείας παλινδρόμησης. Ειδικότερα έχουμε :

Διαστήματα εμπιστοσύνης

Εστω $t_{n-2, \alpha/2}$ η τιμή της κατανομής t-student με $n-2$ βαθμούς ελευθερίας, για την οποία ισχύει $P(t_{n-2} > t_{n-2, \alpha/2}) = \alpha/2$ και επομένως $P(-t_{n-2, \alpha/2} \leq t_{n-2} \leq t_{n-2, \alpha/2}) = 1 - \alpha$

Τότε, το διάστημα το οποίο με επίπεδο εμπιστοσύνης $(1 - \alpha)$ θα περιέχει την παράμετρο b_j , $j=0, 1$ θα είναι το εξής :

$$\hat{b}_j - t_{n-2, \alpha/2} S \hat{b}_j \leq b_j \leq \hat{b}_j + t_{n-2, \alpha/2} S \hat{b}_j$$

Ειδικότερα για $j=0$ το διάστημα αυτό είναι το

$$\hat{b}_0 - t_{n-2, \alpha/2} S \sqrt{\frac{1}{n} + \frac{\bar{x}}{\sum (x_i - \bar{x})^2}} \leq b_0 \leq \hat{b}_0 + t_{n-2, \alpha/2} S \sqrt{\frac{1}{n} + \frac{\bar{x}}{\sum (x_i - \bar{x})^2}}$$

ενώ για $j=1$ είναι το

$$\hat{b}_1 - t_{n-2, \alpha/2} S \sqrt{\frac{1}{\sum (x_i - \bar{x})^2}} \leq b_1 \leq \hat{b}_1 + t_{n-2, \alpha/2} S \sqrt{\frac{1}{\sum (x_i - \bar{x})^2}}$$

Αυτό σημαίνει ότι η πιθανότητα για κάθε ένα από τα διαστήματα αυτά να περιέχει την αντίστοιχη παράμετρο ισούται με $1-\alpha$. Είναι προφανές ότι για μια παράμετρο b_j , όσο πιο μικρό είναι το αντίστοιχο διάστημα εμπιστοσύνης τόσο μεγαλύτερη η ακρίβεια της εκτίμησης. Το εύρος του διαστήματος εμπιστοσύνης είναι $2t_{n-2, \alpha/2} s_{b_j}$, επομένως όσο πιο μικρό το S_{b_j} τόσο μεγαλύτερη η αξιοπιστία του \hat{b}_j

Ελεγχος των υποθέσεων

Όταν εκτιμούμε ένα μοντέλο παλινδρόμησης μας ενδιαφέρει να ελέγξουμε την υπόθεση ότι η συμβολή της x στην ερμηνεία της εξαρτημένης δεν είναι στατιστικά σημαντική. Αυτό ισοδυναμεί με τον έλεγχο της μηδενικής υπόθεσης $H_0: b_1=0$, οπότε η παρατηρούμενη διαφορά του b_1 από το μηδέν μπορεί να αποδοθεί στις διακυμάνσεις της δειγματοληψίας. Ο δίπλευρος έλεγχος της υπόθεσης αυτής θα γίνει ως εξής:

Η μηδενική και η εναλλακτική υπόθεση

$$H_0: b_1 = 0$$

$$H_e: b_1 \neq 0$$

Το κριτήριο απόφασης

$$A_v \quad |t_{\eta}| = \frac{\hat{b}_1}{S_{\hat{b}_1}} > t_{n-2, \alpha/2}$$

Θα απορίψουμε την H_0 στο επίπεδο σημαντικότητας α . Αξίζει να σημειωθεί ότι αν $\alpha=0,05$ και $n-2>12$, η τιμή $t_{n-2, \alpha/2}$ είναι περίπου ίση με 2. Γι' αυτό συχνά χρησιμοποιείται ο ακόλουθος πρακτικός κανόνας: η τιμή του \hat{b}_1 θεωρείται στατιστικά σημαντική, όταν είναι τουλάχιστον διπλάσια

από το s_{b_1} . Πάντως ο κανόνας αυτός, όπως και οποιοσδήποτε άλλος πρακτικός κανόνας θα πρέπει να χρησιμοποιείται με προσοχή.

Είναι δυνατόν, να έχουμε α priori πληροφορίες για το πρόσημο του b_1 οπότε εξειδικεύεται μονόπλευρος έλεγχος ο οποίος είναι πιο ευαίσθητος από το δίπλευρο. Τότε, το στατιστικό του ελέγχου συγκρίνεται με την τιμή

$t_{n-2, \alpha}$ ή την $-t_{n-2, \alpha}$ ανάλογα με το αν ο έλεγχος είναι δεξιόπλευρος ή αριστερόπλευρος.

Εξάλλου, είναι δυνατό να έχουμε α priori πληροφορίες για την τιμή του b_1 . Τότε ο έλεγχος θα γίνει ως εξής:

Η μηδενική και η εναλλακτική υπόθεση

$$H_0: b_1 = b_1^0$$

$$H_1: \text{i) } b_1 \neq b_1^0$$

$$\text{ii) } b_1 > b_1^0$$

$$\text{iii) } b_1 < b_1^0$$

Το κριτήριο απόφασης

Υπολογίζουμε
$$t_n = \frac{\hat{b}_1 - b_1^0}{s_{\hat{b}_1}}$$

οπότε απορρίπτουμε με την H_0 στο επίπεδο σημαντικότητας α , για τα τρία είδη ελέγχου, αντίστοιχα αν:

$$\text{i) } |t_n| > t_{n-2, \alpha/2}$$

$$\text{ii) } t_n > t_{n-2, \alpha}$$

$$\text{iii) } t_n < -t_{n-2, \alpha}$$

Για το σταθερό όρο, συνήθως, γίνεται ο δίπλευρος έλεγχος για τη μηδενική τιμή του. Η διαδικασία είναι ανάλογη μ'αυτήν που περιγράψαμε για τον b_1 όπως και στο παράδειγμα που ακολουθεί.

Σημείωση 1: Όταν οι x και y δεν συνδέονται γραμμικά, τότε $b_1=0$. Και στην περίπτωση αυτή όμως, λόγω διακυμάνσεων της δειγματοληψίας μπορεί να υπολογίσουμε μια εκτίμηση \hat{b}_1 διαφορετική από το μηδέν. Τότε, με πιθανότητα $(1-\alpha)$ θα οδηγηθούμε στην αποδοχή της $H_0: b_1=0$. Θα πρέπει να τονιστεί ότι η αποδοχή της H_0 δεν σημαίνει αναγκαστικά

ότι οι x ή y είναι ανεξάρτητες και άρα ότι θα πρέπει να υιοθετήσουμε το βασικό μοντέλο

$$y_i = b_0 + \varepsilon_i \quad i=1, \dots, n$$

Είναι δυνατόν οι x και y να συνδέονται με μια σχέση η οποία γίνεται γραμμική με τον κατάλληλο μετασχηματισμό της x ή (και) της y .

Σημείωση 2: Είναι προφανές ότι από τους δύο ελέγχους ο πιο σημαντικός είναι ο έλεγχος για την παράμετρο b_1 . Έτσι, όταν δεν μπορούμε να απορίψουμε την $H_0: b_1 = 0$ θα κάνουμε νέα εξειδίκευση του μοντέλου ενώ συνήθως αφήνουμε το σταθερό όρο ακόμη και όταν η υπόθεση $H_0: b_0 = 0$ δεν μπορεί να απορριφθεί.

Παράδειγμα

Στο παράδειγμα της αξίας των πωλήσεων θα κάνουμε το δίπλευρο έλεγχο των υποθέσεων για μηδενικής τιμές των παραμέτρων b_0 και b_1 ως εξής :

Για το σταθερό όρο

$$H_0: b_0 = 0$$

$$H_1: b_0 \neq 0$$

$$\text{Επειδή: } |t_n| = \frac{|\hat{b}_0|}{S_{\hat{b}_0}} = \frac{204,09}{118,14} = 1,73 < t_{16(0,025)} = 2,12$$

δεν μπορούμε να απορίψουμε την H_0

Για την κλίση της ευθείας

$$H_0: b_1 = 0$$

$$H_1: b_1 \neq 0$$

$$\text{Επειδή: } |t_n| = \frac{\hat{b}_1}{S_{\hat{b}_1}} = \frac{5,675}{0,486} = 11,68 > 2,12$$

απορρίπτουμε την H_0 στο $\alpha=0,05$. Στην πραγματικότητα η H_0 απορρίπτεται και για μικρότερες τιμές του α .

2.7 Η ερμηνευτική ικανότητα του μοντέλου

Στο τμήμα αυτό θα παρουσιάσουμε ένα μέτρο της ικανότητας του μοντέλου να ερμηνεύσει τη συνολική μεταβλητότητα των παρατηρήσεων y_i γύρω από τη μέση τιμή τους. Κάθε παρατήρηση y_i μπορεί να γραφεί ως εξής:

$$y_i = \hat{y}_i + \hat{\varepsilon}_i \quad i=1, \dots, n \quad (2.7.1)$$

όπου $\hat{y}_i = \hat{b}_0 + \hat{b}_1 \cdot x_i$ η εκτίμηση που ορίζεται από την ευθεία ελαχίστων τετραγώνων στη τιμή $x=x_i$ και $\hat{\varepsilon}_i = y_i - \hat{y}_i$, τα κατάλοιπα της εκτίμησης. Η (2.7.1) μπορεί να γραφεί ισοδύναμη ως εξής:

$$y_i = \hat{y}_i + (y_i - \hat{y}_i)$$

$$\Leftrightarrow y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i) \quad i = 1, \dots, n \quad (2.7.2)$$

Αν υψώσουμε τις δύο πλευρές στο τετράγωνο και αθροίσουμε για όλες τις τιμές του i τότε, επειδή $\sum (\hat{y}_i - \bar{y}) \cdot (y_i - \hat{y}_i) = 0$ παίρνουμε :

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 \quad (2.7.3)$$

Αν θέσουμε :

$$\sum (y_i - \bar{y})^2 = \text{SST} = \text{Ολικό Αθροισμα Τετραγώνων}$$

$$\sum (\hat{y}_i - \bar{y})^2 = \text{SSR} = \text{Αθροισμα Τετραγώνων Παλινδρόμησης}$$

$$\sum (y_i - \hat{y}_i)^2 = \text{SSE} = \text{Αθροισμα Τετραγώνων των Σφαλμάτων}$$

τότε η (2.7.3.) γράφεται ισοδύναμα ως εξής :

$$\text{SST} = \text{SSR} + \text{SSE}$$

Σημειώνεται ότι αν για την Y εκτιμηθεί το βασικό μοντέλο τότε εκτιμούμε τη σταθερή μέση τιμή με την \bar{Y} και τα κατάλοιπα είναι οι διαφορές $Y_i - \bar{Y}$, $i=1, \dots, n$. Επομένως ο όρος $SST = \sum (Y_i - \bar{Y})^2$ είναι ένα συνολικό μέτρο σφάλματος όταν για την εκτίμηση των τιμών Y_i , $i=1, \dots, n$. Χρησιμοποιείται ο αριθμητικός τους μέσος. Ο όρος $SSE = \sum (Y_i - \hat{Y}_i)^2 = \sum (Y_i - \hat{b}_0 - \hat{b}_1 \cdot x_i)^2$ είναι το αντίστοιχο μέτρο όταν χρησιμοποιείται το απλό γραμμικό μοντέλο παλινδρόμησης. Η διαφορά $SST - SSE$ είναι η μείωση του σφάλματος εκτίμησης που επιτυγχάνεται με το μοντέλο παλινδρόμησης. Η ποσοστιαία μείωση του σφάλματος δηλαδή η

$$PRE = \frac{SST - SSE}{SST}$$

ονομάζεται **συντελεστής προσδιορισμού** και συμβολίζεται με R^2 . Δηλαδή έχουμε:

$$R^2 = PRE = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

Σημειώνεται ότι ο όρος SST είναι ένα μέτρο της συνολικής μεταβλητότητας των παρατηρήσεων Y_i γύρω από τη μέση τιμή τους ενώ ο όρος SSE ένα μέτρο της συνολικής μεταβλητότητας των Y_i γύρω από την ευθεία ελαχίστων τετραγώνων. Επομένως, ο συντελεστής R^2 ερμηνεύεται, ισοδύναμα, ως το ποσοστό της μεταβλητότητας των Y_i γύρω από τη μέση τιμή τους το οποίο ερμηνεύεται από το από το εκτιμηθέν μοντέλο παλινδρόμησης.

Από τον ορισμό αυτόν προκύπτει ότι ισχύει.

$$0 \leq R^2 \leq 1$$

Όσο ο όρος SSR είναι μεγαλύτερος από τον SSE , ή, ισοδύναμα, όσο πιο κοντά στη μονάδα είναι ο συντελεστής R^2 , τόσο μεγαλύτερη είναι η ερμηνευτική ικανότητα του μοντέλου παλινδρόμησης.

Παράδειγμα

Για την αξία των πωλήσεων εκτιμήσαμε το μοντέλο.

$$\hat{Y} = 205 + 5,66$$

Εχουμε επίσης υπολογίσει τα ακόλουθα.

$$\begin{aligned} SSE &= \sum \hat{\varepsilon}_i^2 = 130520 \\ SST &= \sum (Y_i - \bar{Y})^2 = \sum Y_i^2 - n\bar{Y}^2 \\ &= 451102000 - 18(1561,1)^2 \\ &= 1243602,22 \end{aligned}$$

οπότε υπολογίζουμε το συντελεστή προσδιορισμού ως εξής :

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{130520}{1243602,22} = 0,895$$

Δηλαδή το απλό μοντέλο παλινδρόμησης στο οποίο παίρνεται υπόψη η X στην ερμηνεία της συμπεριφοράς των παρατηρήσεων Y_i επιτυγχάνεται, σε σύγκριση με το βασικό μοντέλο μείωση του σφάλματος εκτίμησης ίση με 89,5%.

2.8 Προβλέψεις

Το μοντέλο που εκτιμήσαμε μπορεί να χρησιμοποιηθεί για προβλέψεις Διακρίνουμε δυο περιπτώσεις :

i) **Πρόβλεψη της μέσης τιμής της εξαρτημένης όταν $x = x_{n+1}$**

Αν στην ευθεία ελαχίστων τετραγώνων αντικαταστήσουμε την τιμή $X = X_{n+1}$ θα πάρουμε την

$$\hat{Y}_{n+1} = \hat{b}_0 + \hat{b}_1 X_{n+1}$$

που είναι μια εκτίμηση της $b_0 + b_1 X_{n+1}$ και επομένως είναι μια σημειακή εκτίμηση της μέσης τιμής της εξαρτημένης για $X = X_{n+1}$. Λόγω του σφάλματος δειγματοληψίας η εκτιμούμενη μέση τιμή θα διαφέρει από την πραγματική, με πιθανότητα ίση με 1. Σε

επαναληπτικές όμως προβλέψεις, η μέση τιμή αυτών των διαφορών θα ισούτε με μηδέν διότι

$$E(\hat{Y}_{n+1}) = E(\hat{b}_0) + E(\hat{b}_1) X_{n+1} \\ = b_0 + b_1 X_{n+1}.$$

Όταν κάνουμε μια πρόβλεψη, θα εκτιμήσουμε τη διακύμανση της εκτιμούμενης μέσης τιμής με τον τύπο

$$S^2 \left(\frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right) \quad (2.8.1)$$

όπου

$$S^2 = \frac{\sum \hat{\epsilon}_i^2}{n-2}$$

Αποδεικνύεται ότι, κάτω από την υπόθεση της κανονικότητας των διαταρακτικών όρων, το διάστημα το οποίο, με πιθανότητα $1 - \alpha$, θα περιέχει την

$E(Y_{n+1})$ είναι το

$$\hat{Y}_{n+1} \pm t_{n-2, \alpha/2} S \sqrt{\frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum (X_i - \bar{X})^2}} \quad (2.8.2)$$

όπου $t_{n-2, \alpha/2}$ η τιμή της κατανομής t -student με $v = n - 2$ βαθμού ελευθερίας για την οποία ισχύει $P(t_{n-2} > t_{n-2, \alpha/2})$

Το διάστημα αυτό ερμηνεύεται ως εξής: Εστω ότι έχουμε επαναληπτικές λήψεις δειγμάτων ή ζευγών (X_i, Y_i) στις οποίες οι τιμές των X_i επιλέγονται ίδιες όπως στο δείγμα μας και από κάθε δείγμα εκτιμούμε την ευθεία παλινδρόμησης. Τότε το $(1 - \alpha)\%$ των διαστημάτων τα οποία θα εκτιμήσουμε για $X = X_{n+1}$ θα περιέχουν την $E(Y_{n+1})$. Όταν υπολογίζουμε ένα μόνο διάστημα τότε η πιθανότητα να περιέχει την $E(Y_{n+1})$ ισούτε με $(1 - \alpha)$.

ii) Πρόβλεψη μιας μεμονομένης τιμής της Y όταν $X = X_{n+1}$

Στην περίπτωση αυτή θα πρέπει να προβλέψουμε και μια τιμή για τον E_{n+1} . Επειδή $E(\varepsilon + 1) = 0$, προβλέπουμε την τιμή μηδέν, οπότε προκύπτει και στην περίπτωση αυτή η ίδια σημειακή εκτίμηση. Η αβεβαιότητα όμως στην πρόβλεψη είναι μεγαλύτερη απ' ό,τι στην προηγούμενη περίπτωση λόγω της επίδρασης και του διαταρακτικού όρου E . Έτσι θα εκτιμήσουμε τη διακύμανση της προβλεπόμενης τιμής για την Y_{n+1} με το άθροισμα της διακύμανσης της μέσης τιμής, συν τη διακύμανση του διαταρακτικού όρου E δηλαδή θα είναι ίση με

$$S^2 + S^2 \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum (X_i - \bar{X})^2} = S^2 \left[1 + \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right] \quad (2.8.3)$$

Όταν η κατανομή των διαταρακτικών όρων είναι κανονική, τότε το διάστημα στο οποίο θα βρίσκεται η y_{n+1} με πιθανότητα $(1-\alpha)$, είναι το εξής:

$$\hat{Y}_{n+1} \pm t_{n-2, \alpha/2} S \sqrt{1 + \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum (X_i - \bar{X})^2}} \quad (2.8.4)$$

Η ερμηνεία του διαστήματος αυτού είναι ανάλογη με την ερμηνεία του διαστήματος (2.8.2). Από τις (2.8.2) και (2.8.4) προκύπτουν μερικές γενικές παρατηρήσεις για το εύρος των αντιστοιχών διαστημάτων εμπιστοσύνης:

- Όσο πιο μεγάλο είναι το μέγεθος του δείγματος n , τόσο πιο μικρό το εύρος του αντιστοιχού διαστήματος εμπιστοσύνης, επομένως τόσο μικρότερη η αβεβαιότητα για την αντίστοιχη εκτίμηση, *ceteris paribus*.
- Όσο μεγαλύτερη η ποσότητα $(x_{n+1} - \bar{x})^2$ δηλαδή όσο μεγαλύτερη η απόσταση του x_{n+1} από το δειγματικό μέσο της ερμηνευτικής μεταβλητής τόσο μεγαλύτερο το εύρος του αντιστοιχού διαστήματος, *ceteris paribus*.

- Όσο μεγαλύτερο το πληροφοριακό περιεχόμενο του δείγματος για την εκτίμηση του μοντέλου, δηλαδή ο όρος $\sum_{i=1}^n (x_i - \bar{x})^2$ τόσο πιο στενά τα όρια του αντιστοίχου διαστήματος εμπιστοσύνης, *ceteris paribus*.
- Τέλος, όσο πιο μεγάλη η τιμή του μέσου τετραγωνικού σφάλματος της εκτίμησης s^2 δηλαδή όσο πιο μεγάλη είναι η μεταβλητότητα των παρατηρήσεων γύρω από την εκτιμώμενη ευθεία παλινδρόμησης, τόσο πιο ευρύ είναι το αντίστοιχο διάστημα εμπιστοσύνης.

Παράδειγμα

Εστω ότι θέλουμε να χρησιμοποιήσουμε το μοντέλο που εκτιμήσαμε προηγουμένως για να προβλέψουμε τη μηνιαία αξία των πωλήσεων της επιχείρησης όταν αυτή δαπανά μηνιαίως για διαφήμιση 250 χιλ. δρχ. Η σημειακή εκτίμηση των πωλήσεων για $x_{n+1}=250$ ισούται με :

$$\begin{aligned}\hat{y}_{n+1} &= 204,9 + 5,675 x_{n+1} \\ &= 204,9 + (5,675)250 = 1623,65 \text{ χιλ. δρχ.}\end{aligned}$$

Για να υπολογίσουμε διαστήματα εμπιστοσύνης για την αξία των πωλήσεων και τη μέση αξία των πωλήσεων όταν $y_{n+1}=250$ χρειαζόμαστε τις τυπικές αποκλίσεις των αντιστοίχων εκτιμητών.

Προηγουμένως υπολογίσαμε :

$$\bar{X} = 239,12 \quad \sum (x_i - \bar{x})^2 = 34541 \quad \text{και} \quad S = 90,32$$

οπότε έχουμε

$$S \sqrt{\frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum (x_i - \bar{x})^2}} = 90,32 \sqrt{\frac{1}{18} + \frac{(250 - 239,12)^2}{34541}} = 21,94$$

και

$$S \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum (x_i - \bar{x})^2}} = 90,32 + 21,94 = 112,26$$

Για επίπεδο εμπιστοσύνης $(1 - \alpha) = 0,95$ βρίσκουμε από τον πίνακα της κατανομής t- student την τιμή $t_{n-2, \alpha/2} = t_{16(0,025)} = 2,12$ οπότε υπολογίζουμε τα ακόλουθα διαστήματα εμπιστοσύνης.

Για τη μέση αξία των πωλήσεων όταν $X_{n+1} = 250$

$$1623,65 - (2,12)(21,94) \leq E(Y_{n+1}) \leq 1623,65 + (2,12)(21,94)$$

$$\hat{\eta} \quad 1577,14 \leq E(Y_{n+1}) \leq 1670,16$$

Το εύρος του διαστήματος αυτού ισούτε με 93,02

Για την αξία των πωλήσεων όταν $X_{n+1} = 250$

$$1623,65 - (2,12)(112,26) \leq Y_{n+1} \leq 1623,65 + (2,12)(112,26)$$

$$1385,66 \leq Y_{n+1} \leq 1861,64$$

Το εύρος του διαστήματος αυτού ισούτε με 476 περίπου

Παρατηρούμε ότι όσο απομακρυνόμαστε από τη μέση τιμή των δεδομένων μας τόσο αυξάνει το εύρος των αντίστοιχων διαστημάτων δηλαδή η αβεβαιότητα της αντίστοιχης εκτίμησης και επομένως θα πρέπει να αποφεύγουμε τις προβλέψεις για τιμές της X , οι οποίες είναι μακριά από τα όρια των δεδομένων μας.

2.9 Ανάλυση διακυμάνσεων

Σε κάθε άθροισμα τετραγώνων αντιστοιχεί ένας αριθμός που ονομάζεται βαθμοί ελευθερίας και δείχνει πόσα ανεξάρτητα μέρη πληροφορίας χρησιμοποιούμε από τις η ανεξάρτητες παρατηρήσεις Y_1, Y_2, \dots, Y_n για να το υπολογίσουμε. Έτσι π.χ για το SST χρησιμοποιούμε τις $(n-1)$ ανεξάρτητες διαφορές $(Y_i - \bar{Y})$ οπότε στο SST αντιστοιχούν $n-1$ βαθμοί ελευθερίας. Αντίστοιχα, επειδή

$$\begin{aligned} SSR &= \sum (\hat{Y}_i - \bar{Y})^2 = \sum (\hat{b}_0 + \hat{b}_1 x_i - \hat{b}_0 - \hat{b}_1 \bar{x})^2 \\ &= \sum [\hat{b}_1 (x_i - \bar{x})]^2 \\ &= \hat{b}_1^2 \sum (x_i - \bar{x})^2 \end{aligned}$$

δηλαδή μπορούμε να υπολογίσουμε τα SSR από μια μονάχα συνάρτηση των Y_i έχουμε 1 βαθμό ελευθερίας. Επομένως, με απλή αφαίρεση, στα SSE αντιστοιχούν $(n-1)-1 = n-2$ βαθμοί ελευθερίας. Αυτό αντικατοπτρίζει το γεγονός ότι τα κατάλοιπα ε1 προκύπτουν από ένα εκτιμώμενο μοντέλο με 2 παραμέτρους.

Όλα τα προγράμματα υπολογιστή για τη γραμμική παλινδρόμηση δίνουν μαζί με τους υπόλοιπους υπολογισμούς έναν πίνακα ανάλυσης διακυμάνσεως (Analysis of variance of Anova+able) του οποίου η τυπική μορφή στο απλό μοντέλο είναι η εξής :

SOURCE OF VARIATION	SUM OF SQUARES	D.F	MEAN OF SQUARES
REGRESSION	$\sum (\hat{y}_i - \bar{y})^2$	1	$SSR / 1$
ERROR	$\sum (\hat{y}_i - y_i)^2$	$n-2$	$SSE / (n-2)$
TOTAL	$\sum (y_i - \bar{y})^2$	$n-1$	$SST / (n-1)$

Ετσι στο παράδειγμα της αξίας των πωλήσεων ο πίνακας ANOVA θα δοθεί ως εξής :

SOURCE OF VARIATION	SUM OF SQUARES	D.F	MEAN SQUARE
REGRESSION	1112458,06	1	1112458,06
ERROR	130519,7	16	8157,48
TOTAL	1242977,76	17	73116,34

2.10 Το μοντέλο χωρίς σταθερό όρο.

Όταν θέλουμε η συνάρτηση παλινδρόμησης να περνά από την αρχή των αξόνων τότε το μοντέλο εξειδικεύεται χωρίς σταθερό όρο. Σ' ένα τέτοιο μοντέλο μπορεί να καταλήξουμε :

i. **Υποχρεωτικά** : Όταν τα δεδομένα υπαγορεύουν κάποιον μετασχηματισμό που απαλοίζει το σταθερό όρο. Ετσι π.χ όταν η διακύμανση των διαταρακτικών όρων φαίνεται να αυξάνει με τις τιμές x_i τότε στο μοντέλο

$$Y_i = b_0 + b_1 x_i + \varepsilon_i$$

θα δοκιμάσουμε τον ακόλουθο μετασχηματισμό :

$$\frac{Y_i}{\sqrt{x_i}} = \frac{b_0}{\sqrt{x_i}} + b_1 \sqrt{x_i} + v_i$$

όπου $v_i = \varepsilon_i / \sqrt{x_i}$. Το μετασχηματισμένο αυτό μοντέλο δεν έχει σταθερό όρο αλλά δυο ερμηνευτικές μεταβλητές, τις $X_1 = 1/\sqrt{x}$ και $X_2 = \sqrt{x}$

ii. **Μετά από ανάλυση** : Η σχετική θεωρία μπορεί να υπαγορεύει ότι στη μηδενική τιμή της ερμηνευτικής μεταβλητής αντιστοιχεί μηδενική τιμή της εξαρτημένης. Ετσι π.χ ας θεωρήσουμε την Κενσιανή συνάρτηση κατανάλωσης :

$$E(Y) = b_0 + b_1 x$$

όπου Y_i κατανάλωση όταν το εισόδημα ισούτε με x . Η Οριακή Ροπή προς κατανάλωση (OPK) ισούτε με b_1 ενώ η Μέση Ροπή προς Κατανάλωση (MPK) ισούτε με :

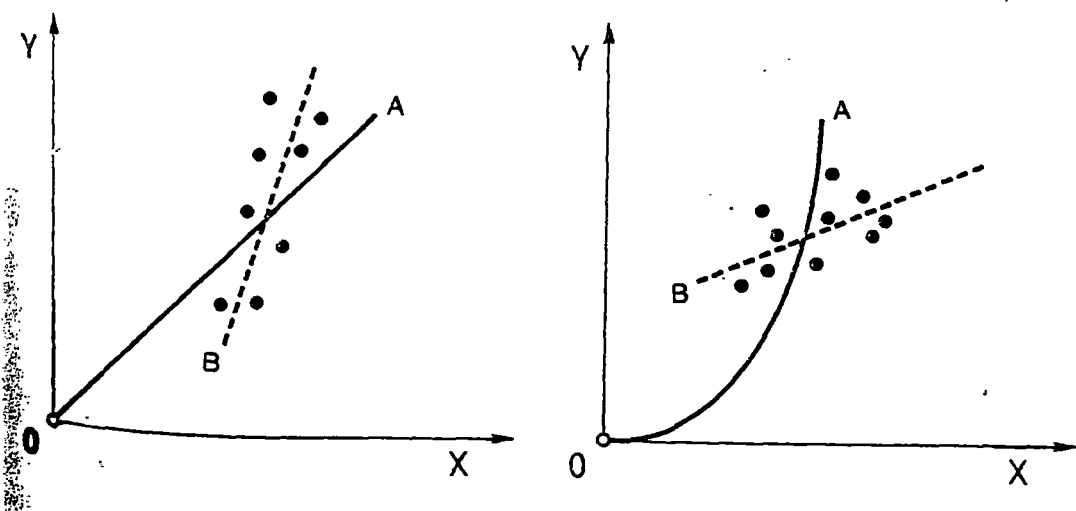
$$\frac{E(Y)}{X} = \frac{b_0}{X} + b_1$$

Εστω ότι έχουμε λόγους να πιστεύουμε ότι $OPK=MPK$ δηλαδή ότι

$$b_1 = \frac{b_0}{X} + b_1$$

Αυτό ισχύει μόνο όταν $b_0=0$. Συμπεραίνουμε ότι αν έχουμε λόγο να πιστεύουμε ότι $OPK=MPK$ τότε θα θεωρήσουμε το Κενσιανικό μοντέλο κατανάλωσης χωρίς σταθερό όρο.

Η σχετική θεωρία πάντως μπορεί να μη ληφθεί υπόψη αν η συμπεριφορά των δεδομένων είναι διαφορετική. Στο σχήμα που ακολουθεί βλέπουμε δυο περιπτώσεις στις οποίες παρόλο που το θεωρητικό μοντέλο παλινδρόμησης (γραμμή A) περνά από την αρχή των αξόνων δηλαδή το σημείο $(X=0, Y=0)$ εν τούτοις για τα δεδομένα που διαθέτουμε, η ευθεία με μη μηδενικό σταθερό όρο (γραμμή B) έχει καλύτερη προσαρμογή.



Τέλος, η επιλογή ανάμεσα σε ένα μοντέλο με και ένα χωρίς σταθερό όρο μπορεί να γίνει με στατιστικά κριτήρια όπως είναι ο t- έλεγχος της υπόθεσης $b_0=0$.

Θα πρέπει πάντως να τονιστεί ότι ακόμη και αν οι σχετικοί στατιστικοί έλεγχοι απορίψουν το σταθερό όρο δεν κρίνεται πάντα σκόπιμη η αφαίρεση του από το μοντέλο. Επισημαίνεται η σχετική παρατήρηση του A. Spanos (1986) : << ο σταθερός όρος θα πρέπει πάντα να περιλαμβάνεται σε ένα μοντέλο παλινδρόμησης για μεταβλητές σε αναλογική κλίμακα.>>

Στη συνέχεια της ανάλυσης επισημαίνονται τα κυριότερα σημεία στα οποία το μοντέλο χωρίς σταθερό όρο διαφέρει από το αντίστοιχο μοντέλο με το σταθερό όρο. Η αναφορά γίνεται στο απλό μοντέλο αλλά η γενίκευση στο πολύ μεταβλητό είναι άμεση.

Η μέθοδος των ελαχίστων τετραγώνων στο μοντέλο χωρίς σταθερό όρο.

Η μέθοδος των ελαχίστων τετραγώνων στο μοντέλο

$$Y_i = b_1 x_i + \varepsilon_i \quad i=1, \dots, n$$

ελαχιστοποιεί τη συνάρτηση

$$S = S(b_1) = \sum_{i=1}^n (Y_i - b_1 x_i)^2 = \sum_{i=1}^n \varepsilon_i^2$$

Επειδή

$$\frac{dS}{db_1} = -2 (\sum Y_i - b_1 \sum x_i) x_i$$

η εκτίμηση των ελαχίστων τετραγώνων του b_1 είναι αυτή που ικανοποιεί την εξίσωση

$$\begin{aligned} -2 (\sum Y_i - \hat{b}_1 \sum x_i) x_i &= 0 \\ \Leftrightarrow \sum Y_i x_i - \hat{b}_1 \sum x_i^2 &= 0 \\ \Leftrightarrow \hat{b}_1 &= \frac{\sum x_i Y_i}{\sum x_i^2} \end{aligned} \quad (2.10.1)$$

Παράδειγμα

Μετρήσαμε την ηλικία X σε λεπτά, ενός ψεκασμού αεροσόλ και την παρατηρούμενη διασπορά του Y σ' αυτόν τον χρόνο. Η διασπορά μετρήθηκε ως το αντίστροφο του αριθμού των σωματιδίων ανά μονάδα όγκου. Σε 9 παρατηρήσεις (X_i, Y_i) πήραμε τα ακόλουθα αποτελέσματα :

X_i 8 22 35 40 57 73 78 87 98
 Y_i 6,16 9,88 14,35 24,06 30,34 32,17 42,18 43,23 48,76

(Τα δεδομένα δίνονται από τους Box - Hunter - Hunter, 1978)

θα εκτιμήσουμε ένα απλό μοντέλο παλινδρόμησης χωρίς σταθερό όρο. Υπολογίζουμε :

$$\sum X_i Y_i = 17638,61 \quad \sum x_i^2 = 35208 \quad \sum Y_i^2 = 8901,31$$

οπότε

$$\hat{b}_1 = \frac{\sum x_i Y_i}{\sum x_i^2} = \frac{17638,61}{35208} = 0,501$$

και το μοντέλο που εκτιμήσαμε είναι το

$$\hat{Y} = 0,501 X$$

Δηλαδή, εκτιμούμε ότι κάθε επιπλέον λεπτό που περνά, αυξάνει τη διασπορά του ψεκασμού κατά μισό περίπου.

Τα κατάλοιπα της εκτίμησης

Στο μοντέλο χωρίς σταθερό όρο τα κατάλοιπα της εκτίμησης $\hat{e}_i = Y_i - \hat{Y}_i = Y_i - \hat{b}_1 x_i$, $i = 1, \dots, n$ ικανοποιούν τις σχέσεις :

$$\begin{aligned} \sum \hat{Y}_i e_i &= 0 \\ \sum Y_i \hat{e}_i &= 0 \end{aligned}$$

αλλά γενικά

$$\sum \hat{\varepsilon}_i \neq 0$$

Παράδειγμα

Υπολογίζουμε τα κατάλοιπα της εκτίμησης στο μοντέλο του προηγούμενου παραδείγματος :

X_i	Y_i	$\hat{Y}_i=0,501X_i$	$\varepsilon_i=Y_i-\hat{Y}_i$
8	6,16	4,0079	2,1521
22	9,88	11,0216	-1,1416
35	14,35	17,5344	-3,1844
40	24,06	20,0393	4,0207
57	30,34	28,5660	1,7840
73	32,17	36,5718	-4,4018
78	42,18	39,0767	3,1033
87	43,23	43,5855	-0,3555
98	48,76	49,0963	-0,3363
Άθροισμα 498	251,13		1,6405

Ανάλυση Διακυμάνσεως

Γενικά

$$\sum (Y_i - \hat{Y}_i)^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum \hat{\varepsilon}_i^2$$

ή $SST = SSR + SSE$

ενώ ισχύει

$$\sum Y_i^2 = \sum \hat{Y}_i^2 + \sum \hat{\varepsilon}_i^2$$

$$(n) \cdot (\bar{y})^2 + (n-k)$$

και οι αριθμοί στις παρενθέσεις είναι οι βαθμοί ελευθερίας του κάθε όρου. Για το απλό μοντέλο $k=1$. Από τη σχέση αυτή προκύπτει ότι γενικά

$$\frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = 1 - \frac{\sum \hat{\varepsilon}_i^2}{\sum (Y_i - \bar{Y})^2}$$

Ο συντελεστής προσδιορισμού

Μπορούμε να πάρουμε για το συντελεστή προσδιορισμού R^2 μια τιμή που να ικανοποιεί τη σχέση $0 \leq R^2 \leq 1$ μόνο από το τύπο

$$R^2 = 1 - \frac{\sum \hat{\varepsilon}_i^2}{\sum Y_i^2} \quad (2.10.2)$$

Στα περισσότερα προγράμματα H/Y, πάντως (μεταξύ αυτών και το RATS) ο συντελεστής R^2 υπολογίζεται σε όλες τις περιπτώσεις από τον τύπο

$$R^2 = 1 - \frac{\sum \hat{\epsilon}_i^2}{\sum (y_i - \bar{y})^2} \quad (2.10.3)$$

Ετσι, στα μοντέλα χωρίς σταθερό όρο, μπορεί να πάρουμε αρνητική τιμή για τον R^2 . Ο αρνητικός συντελεστής προσδιορισμού είναι φυσικά δείκτης κακής προσαρμογής του μοντέλου. Θα πρέπει να σημειωθεί ότι αυτό ισχύει για οποιονδήποτε περιορισμό στις παραμέτρους του μοντέλου, ένας λανθασμένος περιορισμός που επιβάλλεται στις παραμέτρους μπορεί να έχει ως αποτέλεσμα το SSE να γίνει μεγαλύτερο από το SST.

Η διακύμανση του \hat{b}_1

Αποδεικνύεται ότι η διακύμανση του b_1 δίνεται από τον τύπο :

$$\text{Var}(\hat{b}_1) = \frac{\sigma^2}{\sum x_i^2} < \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

Επομένως στο μοντέλο χωρίς σταθερό όρο, η εκτίμηση του \hat{b}_1 γίνεται, γενικά, με μεγαλύτερη ακρίβεια απ' ό,τι στο μοντέλο με σταθερό όρο. Ο αμερόληπτος εκτιμητής της σ^2 για το απλό μοντέλο χωρίς σταθερό όρο υπολογίζεται από τον τύπο

$$S^2 = \frac{\sum \hat{\varepsilon}_i^2}{n-1}$$

Αντίστοιχα, ο αμερόληπτος εκτιμητής της $\text{Var}(\hat{b}_1)$ δίνεται από τον τύπο

$$S_{\hat{b}_1}^2 = \frac{S^2}{\sum x_i^2}$$

και ισχύει :

$$\frac{\hat{b}_1 - b_1}{S_{\hat{b}_1}} \sim t_{n-1}$$

Ο έλεγχος των υποθέσεων και το διάστημα εμπιστοσύνης για το b_1 γίνεται όπως και στο μοντέλο με σταθερό όρο.

Παράδειγμα

Εκτιμούμε την διακύμανση του \hat{b}_1 στο παράδειγμα της διασποράς του αεροσόλ ως εξής :

$$S_{\hat{b}_1}^2 = \frac{S^2}{\sum x_i^2} = \frac{8,08}{35208} = 0,000229$$

και το τυπικό σφάλμα της εκτίμησης ως εξής :

$$S_{\hat{b}_1} = \sqrt{0,000229} = 0,015$$

Ο έλεγχος για μηδενική τιμή του b_1 θα γίνει ως εξής :

$$H_0 : b_1 = 0$$

$$H_1 : b_1 \neq 0$$

Επειδή

$$t_p = \frac{\hat{b}_1}{s_{\hat{b}_1}} = \frac{0,501}{0,015} = 33,4 > t_{\frac{\alpha}{2}(n-2)} = 2,306$$

απορρίπτουμε την H_0 . Το διάστημα το οποίο με πιθανότητα 95% θα περιέχει την τιμή του b_1 είναι το εξής :

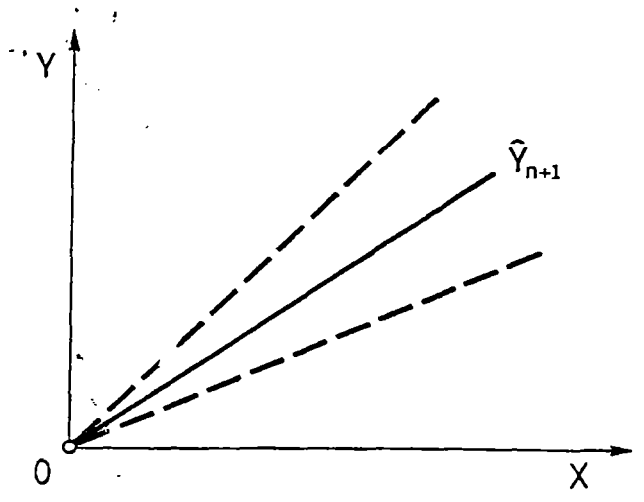
$$\begin{aligned} & \hat{b}_1 \pm (0,01515)(2,306) \\ & \text{ή } 0,501 \pm (0,01515)(2,306) \\ & \text{ή } 0,466 \leq b_1 \leq 0,536 \end{aligned}$$

Προβλέψεις και αντίστοιχες λωρίδες εμπιστοσύνης

Στην τιμή X_{n+1} προβλέπουμε την μέση τιμή $E(Y_{n+1})$ της εξαρτημένης με την $\hat{Y}_{n+1} = \hat{b}_1 X_{n+1}$. Το διάστημα στο οποίο θα βρίσκεται η $E(Y_{n+1})$ με επίπεδο εμπιστοσύνης $(1 - \alpha)$ είναι το :

$$\hat{Y}_{n+1} \pm t_{n-1, \alpha/2} S \sqrt{\frac{X_{n+1}^2}{\sum X_i^2}}$$

Οι λωρίδες εμπιστοσύνης τείνουν στο μηδέν όσο πλησιάζουμε στην αρχή των αξόνων όπως φαίνεται και στο επόμενο σχήμα. Αυτό είναι σύμφωνο με την υπόθεση ότι με πιθανότητα ίση με 1 η $E(Y_i)$ ισούτε με μηδέν όταν $X_i=0$



Λωρίδες εμπιστοσύνης του απλού μοντέλου χωρίς σταθερό όρο.

ΚΕΦΑΛΑΙΟ ΤΡΙΤΟ

Συσχέτιση

3.1 Εισαγωγή

Στην ανάλυση συσχέτισης μας ενδιαφέρει να δούμε αν δυο τυχαίες μεταβλητές X και Y συνδέονται με μια σχέση και πόσο στενή είναι αυτή η σχέση. Η ανάλυση συσχέτισης εφαρμόζεται είτε αυτοτελώς, είτε ως ένα προκαταρκτικό στάδιο πριν από την ανάλυση παλινδρόμησης. Έτσι, όταν έχουμε πολλές υποψήφιες ερμηνευτικές μεταβλητές η ανάλυση συσχέτισης μπορεί, καταρχήν, να δείξει ποιές απ' αυτές συνδέονται περισσότερο με την εξαρτημένη μεταβλητή και θα πρέπει να περιληφθούν στο μοντέλο.

Οι συντελεστές συσχέτισης συνδέονται στενά με τους συντελεστές παλινδρόμησης όμως, το αναλυτικό πλαίσιο, συνεπώς και το πεδίο εφαρμογών των δυο μεθόδων, είναι διαφορετικό. Στην ανάλυση παλινδρόμησης έχουμε μια τυχαία μεταβλητή Y της οποίας τη συμπεριφορά θέλουμε να ερμηνεύσουμε ή να προβλέψουμε με τη συμπεριφορά της μεταβλητής X .

Η X είναι ελεγχόμενη δηλαδή οι τιμές της έχουν επιλεγεί με βάση ένα σχέδιο πειράματος. Η X μπορεί να είναι και τυχαία μεταβλητή αλλά και στην περίπτωση αυτή αντιμετωπίζεται ως ελεγχόμενη με την έννοια ότι δεν γίνεται καμιά αναφορά στην πιθανοκρατική συμπεριφορά της- οι παρατηρήσεις X_i θεωρούνται ως οι τιμές της X στις οποίες επιλέξαμε να παρατηρήσουμε την Y και τα αποτελέσματα ερμηνεύονται υπό συνθήκη των X_i .

Στην ανάλυση συσχέτισης οι δυο μεταβλητές X και Y είναι τυχαίες και αντιμετωπίζονται συμμετρικά. Εξάλλου, όταν οι X και Y είναι μετρήσιμες μόνο στην ταξική κλίμακα μέτρησης, συνήθως περιοριζόμαστε στον υπολογισμό του αντίστοιχου συντελεστή συσχέτισης. Έτσι δικαιολογείται και η μάλλον υπερβολική αλλά όχι αβάσιμη δήλωση του Tukey (1954) ότι οι συντελεστές συσχέτισης είναι χρήσιμοι σε δυο και μόνον περιπτώσεις : όταν είναι συντελεστές παλινδρόμησης ή όταν η μέτρηση μιας ή περισσοτέρων μεταβλητών δεν μπορεί να γίνει παρά μόνον σε ταξική κλίμακα. Σημειώνεται πάντως ότι τα τελευταία χρόνια έχουν αναπτυχθεί τεχνικές της ανάλυσης παλινδρόμησης που είναι κατάλληλες και για τις μεταβλητές αυτές και είναι γνωστές με τον όρο **στιβαρή παλινδρόμηση**.

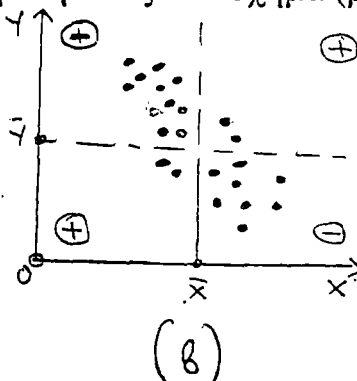
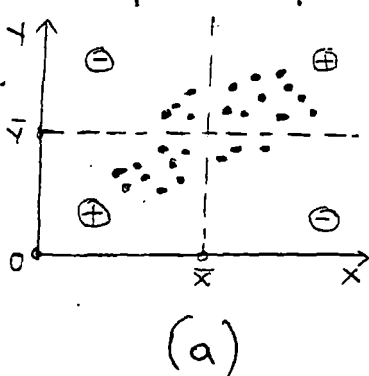
3.2 Η συνδιακύμανση

Εστω ότι σε κάθε στοιχείο το οποίο παίρνεται από τον πληθυσμό, με απλή τυχαία δειγματοληψία παρατηρείται η μεταβλητή X και συγχρόνως η μεταβλητή Y και διαθέτουμε τα n ζεύγη παρατηρήσεων $(X_1, Y_1), \dots, (X_n, Y_n)$. Κάθε ζεύγος (X_i, Y_i) , $i=1, \dots, n$ αποτελεί μια διμεταβλητή παρατήρηση και παριστάνεται με ένα σημείο στο επίπεδο XY . Η θέση των παρατηρήσεων X_i στον άξονα των X δίνεται από τον αριθμητικό τους μέσο $\bar{X} = \sum X_i / n$ ενώ η διασπορά τους από τη διακύμανση τους $S^2_x = \sum (X_i - \bar{X})^2 / (n - 1)$. Αντίστοιχα, η θέση και διασπορά των Y_i δίνεται από τα $\bar{Y} = \sum Y_i / n$ και $S^2_y = \sum (Y_i - \bar{Y})^2 / (n - 1)$. Οι παράμετροι \bar{X} , \bar{Y} , S^2_x , S^2_y , δεν μας δίνουν καμιά πληροφορία για την κοινή κατανομή των X και Y δηλαδή για τον τρόπο που τα σημεία (X_i, Y_i) κατανέμονται στο επίπεδο XY . Για το σκοπό αυτό υπολογίζουμε τη συνδιακύμανση η οποία ορίζεται ως εξής :

$$S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

και μετρά την γραμμική συμμεταβολή δηλαδή την τάση των σημείων (X_i, Y_i) να συγκεντρώνονται κατά μήκος μιας ευθείας. Ανάλογα με το αν η συμμεταβολή των παρατηρήσεων X και Y είναι ευθεία ή αντίστροφη, το πρόσημο της συνδιακύμανσης είναι θετικό ή αρνητικό. Ειδικότερα παρατηρούμε τα εξής :

Αν στο ζεύγος (X_i, Y_i) η παρατήρηση X_i είναι μεγαλύτερη (μικρότερη) από τον αριθμητικό μέσο \bar{X} και η Y_i μεγαλύτερη (μικρότερη) από τον \bar{Y} τότε το γινόμενο $(X_i - \bar{X})(Y_i - \bar{Y})$ είναι θετικό. Αν η παρατήρηση της μιας μεταβλητής είναι μεγαλύτερη από το μέσο της, ενώ της άλλης είναι μικρότερη τότε το γινόμενο $(X_i - \bar{X})(Y_i - \bar{Y})$ είναι αρνητικό. Όταν οι X και Y τείνουν να συμμεταβάλλονται ευθέως τότε τα θετικά γινόμενα $(X_i - \bar{X})(Y_i - \bar{Y})$ είναι περισσότερα από τα αρνητικά και η συνδιακύμανση είναι θετική, όπως στο σχήμα (α) που ακολουθεί. Αντίστοιχα, όταν οι X και Y τείνουν να συμμεταβάλλονται αντιστρόφως τότε τα αρνητικά γινόμενα $(X_i - \bar{X})(Y_i - \bar{Y})$ είναι περισσότερα από τα θετικά και η συνδιακύμανση αρνητική όπως στο σχήμα (β).



Η μηδενική συνδιακύμανση δείχνει αν οι X και Y δεν συµµεταβάλλονται ούτε ευθέως ούτε αντιστρόφως. Είναι προφανές ότι όταν $X_i = Y_i$, V_i , τότε :

$$S_{xy} = \frac{1}{1-n} \sum_{i=1}^n (X_i - \bar{X})^2 = S_x^2$$

δηλαδή, η συνδιακύμανση των παρατηρήσεων X_i είναι απλώς η διακύμανση τους.

3.3 Ο συντελεστής συσχέτισης

Η συνδιακύμανση εκφράζεται σε (μονάδες X) (μονάδες Y). Διαιρώντας την με το γινόμενο $S_x S_y$ προκύπτει ένα μέτρο γραμμικής συµµεταβολής που είναι απαλλαγμένο από μονάδες μέτρησης και ονομάζεται **συντελεστής συσχέτισης του Pearson** ή του γινομένου των ροπών. Ο συντελεστής αυτός συµβολίζεται με r_{xy} οπότε έχουμε :

$$r_{xy} = \frac{S_{xy}}{S_x \cdot S_y}$$

ή

$$r_{xy} = \frac{\sum (X_i - \bar{X}) \cdot (Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \sqrt{\sum (Y_i - \bar{Y})^2}} = \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sqrt{\sum X_i^2 - n \bar{X}^2} \sqrt{\sum Y_i^2 - n \bar{Y}^2}}$$

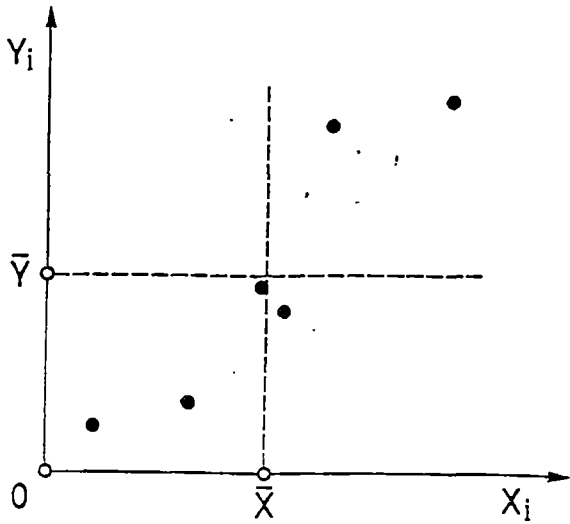
και ο συντελεστής Σ είναι για τιμές του i από 1 ως n .

Ο συντελεστής r_{xy} είναι καθαρός αριθμός επιτρέποντας έτσι τις συγκρίσεις, ενώ παράλληλα έχει ενδιαφέρουσες στατιστικές ιδιότητες. Γι' αυτό χρησιµοποιείται πολύ πιο συχνά από την συνδιακύμανση. Τις βασικές ιδιότητες του r_{xy} θα δούµε στο επόμενο τμήμα. Προηγουµένως θα δούµε την πληροφορία που δίνουν οι συντελεστές S_{xy} και r_{xy} για τη συµµεταβολή η διµεταβλητών παρατηρήσεων με το ακόλουθο.

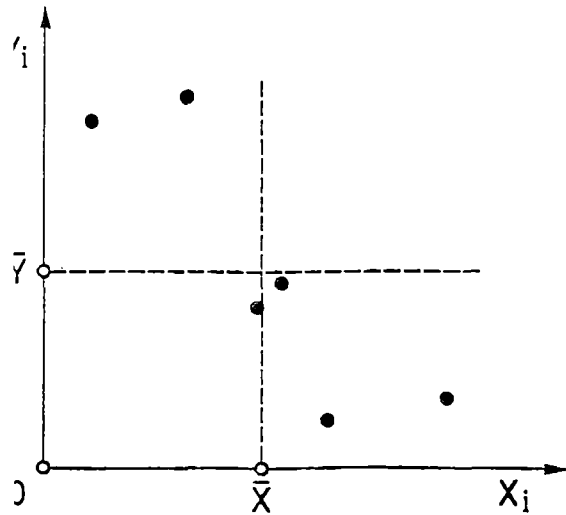
Παράδειγμα

Δίνεται η γραφική παράσταση τριών διµεταβλητών δειγµάτων με $n=6$. Οι παρατηρήσεις (X_i , Y_i) σηµειώνονται στον άξονα X και Y αντίστοιχα,

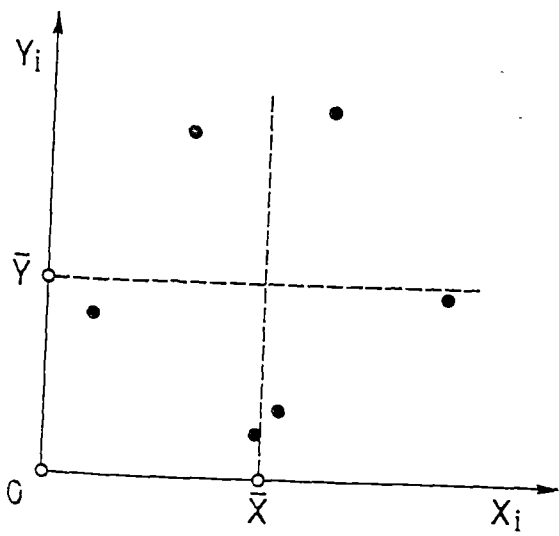
με τις τελείες και είναι ίδιες και στα τρία δείγματα. Η, διαφορετικά, οι περιθώριες κατανομές των X_i και Y_i είναι ίδιες και στα τρία δείγματα. Η κοινή κατανομή όμως των 6 σημείων (X_i, Y_i) στο επίπεδο XY είναι εντελώς διαφορετική όπως φαίνεται στα αντίστοιχα διαγράμματα διασποράς. Ομοίως, διαφέρουν οι συντελεστές συνδιακύμανσης και συσχέτισης.



X_i 2 6 9 10 12 17
 Y_i 2 3 8 7 15 16
 $X = 9,3$ $S_y = 10,65$ $S_{xy} = 28,34$
 $Y = 8,5$ $S_y = 5,89$ $r_{xy} = 0,45$



X_i 2 6 9 10 12 17
 Y_i 15 16 7 8 2 3
 $X = 9,3$ $S_x = 10,65$ $S_{xy} = 28,34$
 $Y = 8,5$ $S_y = 5,89$ $r_{xy} = 0,78$



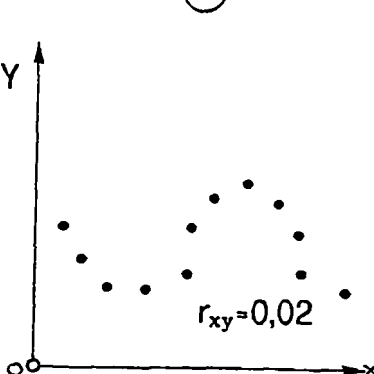
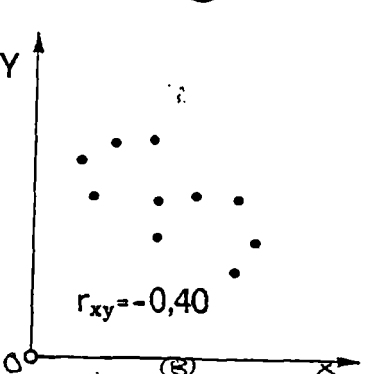
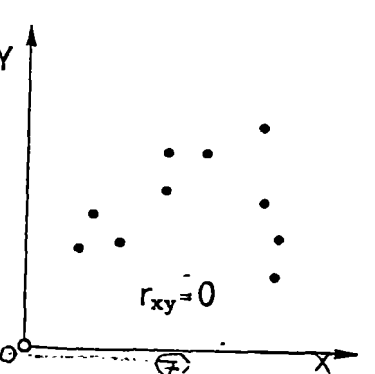
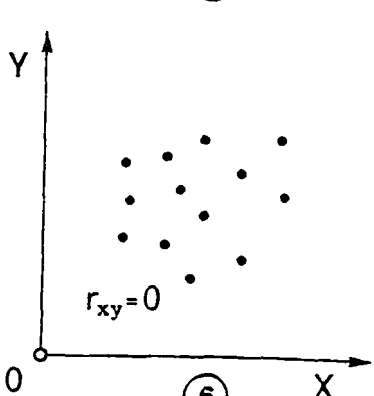
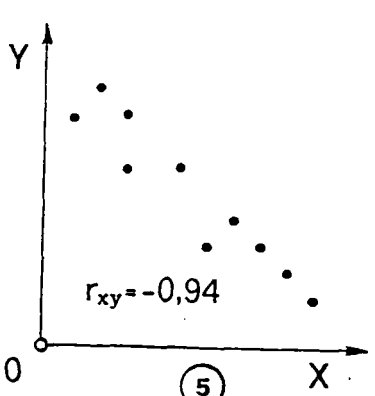
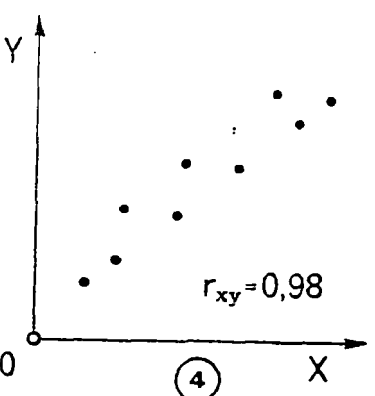
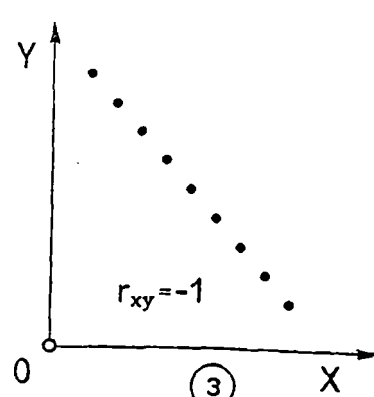
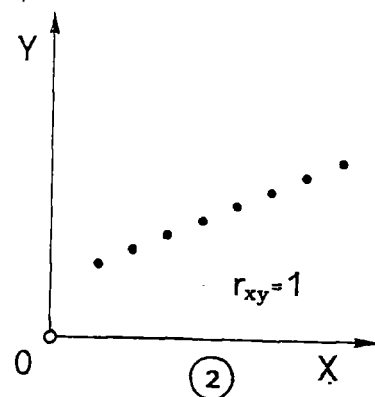
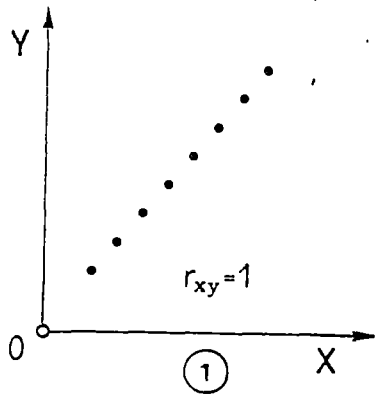
X_i 2 6 9 10 12 17
 Y_i 7 15 2 3 16 8
 $\bar{X} = 9,3$ $S_x = 10,65$ $S_{xy} = 21,81$
 $\bar{Y} = 8,5$ $S_y = 5,89$ $r_{xy} = -0,35$

3.4 Ιδιότητες του συντελεστή συσχέτισης

Αποδεικνύεται ότι ο συντελεστής r_{xy} έχει τις ακόλουθες ιδιότητες

- Ιδιότητα 1η : Παίρνει τιμές στο κλειστό διάστημα $[-1,1]$

Οι ακραίες τιμές -1 και 1 αντιστοιχούν στην περίπτωση που όλα τα σημεία (X_i, Y_i) , $i = 1, \dots, n$ βρίσκονται επάνω σε μια ευθεία με αρνητική ή θετική κλίση αντίστοιχα. Όταν οι X_i και Y_i είναι παρατηρήσεις τυχαίων μεταβλητών, η πιθανότητα να υπολογίσουμε συντελεστή συσχέτισης ίσο με 1 ή -1 είναι μηδενική (σημειώνεται πάντως το προφανές $r_{xy} = 1$). Οι ακραίες τιμές όμως χρησιμεύουν για να ερμηνεύσουμε τις διάμεσες. Όσο πιο κοντά σε μια απ' αυτές βρίσκεται ο r_{xy} τόσο πιο έντονη η γραμμική συμμεταβολή των παρατηρήσεων X και Y . Μια τιμή για τον r_{xy} ίση ή πολύ κοντά στο μηδέν, δηλώνει απουσία οποιασδήποτε σχέσης. Θα πρέπει να τονιστεί ότι απρόβλεπτη είναι η επίδραση στην τιμή του r_{xy} μιας μη γραμμικής όπως και μιας ή περισσότερων ακραίων τιμών.



- Ιδιότητα 2^η : Η τιμή του δεν μεταβάλλεται με τον γραμμικό μετασχηματισμό των X_i, Y_i .

Ετσι, αν $Z_i = a + b X_i, W_i = c + d Y_i, i = 1, \dots, n$ και τα b, d έχουν το ίδιο πρόσημο, τότε

$$\Gamma_{xy} = \Gamma_{zw}$$

Αντίστοιχα, αν τα b, d έχουν διαφορετικό πρόσημο, τότε

$$\Gamma_{xy} = - \Gamma_{zw}$$

Η ιδιότητα αυτή είναι ιδιαίτερα χρήσιμη διότι μας επιτρέπει να κάνουμε μετασχηματισμούς κλίμακας και αρχής των δεδομένων μας, έτσι ώστε να διευκολύνουμε τους υπολογισμούς. Ορισμένα μάλιστα υποπρογράμματα μετασχηματίζουν αυτομάτως τα δεδομένα στο διάστημα $(0,1)$ πριν υπολογίσουν τον

Παράδειγμα

Δίνεται ο ετήσιος αριθμός των σκαφών X_i και ο αριθμός των πληρωμάτων Y_i της Ελλάδος κατά τα έτη 1968-1981. Οι παρατηρήσεις είναι εκφρασμένες σε δεκάδες.

Έτη	1968	1969	1970	1971	1972	1973	1974	1975	1977	1978
X_i	9,5	9,9	7,0	5,7	6,0	6,2	6,2	6,0	8,3	8
Y_i	69,7	84,3	49,0	30,7	33,5	34,9	35,0	29,2	37,9	36,5

1979	1980	1981
8,1	7,1	8,2
37,7	25,5	23,3

Πηγή : ΕΣΥΕ, Στατιστική Επετηρίς της Ελλάδος 1984.

Υπολογίζουμε :

$$\sum X_i = 102,2$$

$$\sum X_i^2 = 770,72$$

$$\bar{X} = 7,3$$

$$\sum Y_i = 564,5$$

$$\sum Y_i^2 = 26520,35$$

$$\bar{Y} = 40,3$$

$$\sum Y_i X_i = 4333,24$$

ΟΠΟΤΕ :

$$r_{xy} = \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sqrt{\sum X_i^2 - n \bar{X}^2} \sqrt{\sum Y_i^2 - n \bar{Y}^2}}$$

$$= \frac{4333,24 - 14(7,3) \cdot (40,3)}{\sqrt{770,72 - 14(7,3)^2} \sqrt{26520,35 - 14(40,3)^2}} = 0,70$$

Αν οι παρατηρήσεις δεν εκφράζονται σε δεκάδες, αλλά σε μονάδες, ο συντελεστής συσχέτισης παραμένει ο ίδιος.

3.5 Ερμηνεία του συντελεστή συσχέτισης

Επαναλαμβάνεται ότι ο συντελεστής συσχέτισης δεν είναι παρά μόνο ένα μέτρο γραμμικής συμμεταβολής δυο τυχαίων μεταβλητών. Επειδή συχνά παρερμηνεύεται από τους μη ειδικούς κρίθηκε σκόπιμο να γίνουν οι ακόλουθες διευκρινίσεις:

Συσχέτιση και αιτιότητα

Η συσχέτιση δεν σημαίνει αιτιότητα. Στη διεθνή βιβλιογραφία αναφέρεται χαρακτηριστικά το εξής παράδειγμα : ο υψηλός συντελεστής συσχέτισης ανάμεσα στις αφίξεις πελαργών σε μια περιοχή και στον αριθμό των γεννήσεων, δεν μας επιτρέπει να συμπεραίνουμε ούτε ότι οι πελαργοί φέρνουν τα μωρά (ούτε πολύ περισσότερο για τον έλεγχο των γεννήσεων να προτείνουμε τη θανάτωση τους) ούτε ότι τα μωρά προσελκύουν τους πελαργούς.

Γενικά, υψηλός συντελεστής συσχέτισης μπορεί να παρατηρηθεί όταν η X προκαλεί μεταβολές στην Y ή η Y προκαλεί μεταβολές στην X ή όταν μια Τρίτη μεταβλητή η οποία ονομάζεται **κεκαλυμμένη** προκαλεί μεταβολές στην ίδια ή αντίθετη κατεύθυνση και στις δυο μεταβλητές ή τέλος, μπορεί να οφείλεται στην τύχη. Στις δυο τελευταίες περιπτώσεις η συσχέτιση ονομάζεται **νόθα ή χωρίς νόημα**.

Στις περισσότερες περιπτώσεις αρκεί ο κοινός νους για να συμπεράνουμε αν μια παρατηρούμενη συσχέτιση μπορεί να ερμηνευτεί ή είναι χωρίς νόημα.

Συσχέτιση και ανεξαρτησία

Ο συντελεστής συσχέτισης μετρά την ένταση γραμμικής συμμεταβολής. Έτσι δυο μεταβλητές X, Y μπορεί να έχουν συντελεστή συσχέτισης ίσο με μηδέν και οι μεταβλητές να μην είναι ανεξάρτητες αλλά να συνδέονται με μια σχέση η οποία δεν είναι γραμμική. Μηδενική συσχέτιση δεν συνεπάγεται ανεξαρτησία εκτός από την ακόλουθη περίπτωση : Όταν η κοινή κατανομή των X και Y είναι η κανονική όπως θα δούμε στο τμήμα 8 αυτού του κεφαλαίου.

Συσχέτιση και το εύρος των δεδομένων

Το μικρό εύρος των παρατηρήσεων της μιας ή και των δυο μεταβλητών έχει συνήθως ως αποτέλεσμα να υπολογίσουμε μια μικρή τιμή για τον r_{xy} . Αν υποπτευόμαστε ότι η μικρή τιμή του r_{xy} οφείλεται σ' αυτόν το λόγο και το επιτρέπουν οι συνθήκες δειγματοληψίας θα πρέπει να διευρύνουμε το δείγμα.

Έτσι π.χ αν μας ενδιαφέρει η συσχέτιση του ύψους X ενός παιδιού με το Νο Y του παπουτσιού του είναι προτιμότερο να μην πάρουμε δείγμα παιδιών μόνον από την ίδια τάξη γιατί τότε το εύρος των δυο μεταβλητών θα είναι πιθανότατα μικρό και συνεπώς η τιμή του r_{xy} μικρή ακόμη και αν η τιμή του ρ_{xy} είναι μεγάλη.

Συσχέτιση και ακραίες τιμές

Η ύπαρξη μιας ή περισσότερων ακραίων τιμών μπορεί να έχει σημαντική επίδραση στην τιμή του r_{xy} , η έκταση και η κατεύθυνση της οποίας μπορεί να είναι απρόβλεπτη. Το διάγραμμα διασποράς των τιμών (X_i, Y_i) , $i = 1, \dots, n$ είναι διαφωτιστικό στην περίπτωση αυτή.

3.6. Ο συντελεστής συσχέτισης ομαδοποιημένων παρατηρήσεων.

Η διμεταβλητή κατανομή συχνοτήτων.

Συχνά η διμεταβλητές παρατηρήσεις για τις οποίες θέλουμε να υπολογίσουμε τον συντελεστή συσχέτισης, είναι ομαδοποιημένες σε μια διμεταβλητή κατανομή συχνοτήτων όπως στο ακόλουθο,

Παράδειγμα

Τυχαίο δείγμα 100 γυναικών που γέννησαν μη πρόωρα βρέφη ομαδοποιήθηκε στην ακόλουθη κατανομή ανάλογα με το ύψος X της μητέρας και το ύψος Y του παιδιού, κατά τη γέννηση.

Υψος παιδιού Y	34 - 38	38 - 42	42 - 46	46 - 50	50 - 54	54 - 58	58 - 62	Άθροισμα
Y_i	36	40	44	48	52	56		
Υψος μητέρας X								
X_i								
154 - 158	4	2	1					7
156								
158 - 162	2	5	2					9

160								
162 - 166	1	4	6	5				16
164								
166 - 170		6	10	7	1			24
168								
170 - 174	1		14	12				27
172								
174 - 178		2	4	2	1	1		10
176								
178 - 182				2	4	1		7
180								
Άθροισμα	8	19	37	28	6	2		100

Υπολογισμός του συντελεστή συσχέτισης

Για να υπολογίσουμε το συντελεστή συσχέτισης ή ομαδοποιημένων διμεταβλητών παρατηρήσεων εργαζόμαστε ως εξής :

- Από τα περιθώρια κατανομή των X υπολογίζουμε τα \bar{X} , S_x^2 .
- Από τα περιθώρια κατανομή των Y υπολογίζουμε τα \bar{Y} , S_y^2 .
- Αν F_{ij} η συχνότητα στη θέση ή "κελί" ij , δηλαδή στη θέση που είναι κοινή στην i -γραμμή και στην j -στήλη του διμεταβλητού πίνακα συχνοτήτων και X_i, Y_j οι κεντρικές τιμές των αντίστοιχων F_{ij}, X_i, Y_j και αθροίζουμε για όλα τα i και j .
- Υπολογίζουμε την συνδιακύμανση

$$S_{xy} = \frac{1}{n} \sum_i \sum_j (x_i - \bar{x}) \cdot (y_j - \bar{y}) f_{ij} \quad (3.6.1.)$$

$$S_{xy} = \frac{1}{n} \sum_i \sum_j x_i y_j f_{ij} - \bar{x} \bar{y}$$

και το συντελεστή συσχέτισης

$$r_{xy} = \frac{S_{xy}}{S_x \cdot S_y} \quad (3.6.2.)$$

Παράδειγμα

Υπολογίζουμε τον συντελεστή συσχέτισης του ύψους X της μητέρας και του ύψους Y του παιδιού στις 100 διμεταβλητές παρατηρήσεις του προηγούμενου παραδείγματος, σταδιακά, ως εξής:

- Από την περιθώρια κατανομή των X υπολογίζουμε την μέση τιμή και την διακύμανση τους. Ειδικότερα έχουμε:

Τάξεις	X_i	F_i	X_i'	$X_i F_i$	$X_i^2 F_i$
154 - 158	156	7	-3	-21	63
158 - 162	160	9	-2	-18	36
162 - 166	164	16	-1	-16	16
166 - 170	168	24	0	0	0
170 - 174	172	27	1	27	27
174 - 178	176	10	2	20	40
178 - 182	180	7	3	21	63
Άθροισμα		100		13	245

Η βοηθητική μέση τιμή

$$\bar{X}' = \frac{1}{n} \sum X_i' F_i = \frac{13}{100} = 0,13$$

και η βοηθητική διακύμανση

$$S'^2_x = \frac{1}{n} \sum X_i'^2 F_i - \bar{X}'^2 = \frac{1}{100} 245 - (0,13)^2 = 2,433$$

οπότε

$$\bar{X} = 168 + (0,13) 4 = 168,52$$

$$S^2_x = 4^2 (2,433) = 38,93 \text{ και } S_x = 6,24$$

- Ομοίως, από την περιθώρια κατανομή των Y υπολογίζουμε

$$\bar{Y} = 44,44$$

$$S^2_y = 19,166 \text{ και } S_y = 4,38$$

- Για να υπολογίσουμε την συνδιακύμανση υπολογίζουμε αρχικά τους όρους:

$$\begin{aligned} \sum_i \sum_j X_i Y_j F_{ij} &= (156)(36)4 + (156) \cdot (40) \cdot 2 + (156) \cdot 44 \\ &+ (160)(36)2 + (160) \cdot (40) 5 + (160) \cdot (44) \cdot 2 \\ &+ (164)(36) + (164) \cdot (40) \cdot 4 + (164)(44)6 + (164) \cdot \\ &\cdot (48) \cdot 5. \end{aligned}$$

$$\begin{aligned}
& + (168)(40)6 + (168)(44)10 + (168)(48)7 + (168)(52) \\
& + (172)(36) + (172)(44) \cdot 14 + (172)(48) \cdot 12 \\
& + (176)(40)2 + (176)(44)4 + (176)(48)2 + (176)52 + \\
& \cdot (176)56 \\
& + (180)(48)2 + (180)(52) \cdot 4 + (180) \cdot 56 = 750608
\end{aligned}$$

και

$$\bar{X} \bar{Y} = (168, 52)(44, 44) = 7489,03$$

οπότε

$$S_{xy} = \frac{750608}{100} - 7489,03 = 17,05$$

και

$$r_{xy} = \frac{S_{xy}}{S_x \cdot S_y} = \frac{17,05}{(6,24)(4,38)} = 0,624$$

3.7 Ο Θεωρητικός συντελεστής συσχέτισης

Ο συντελεστής συσχέτισης δύο τυχαίων μεταβλητών X και Y με μέση τιμή μ_x και μ_y , αντίστοιχα, συμβολίζεται με ρ_{xy} και ορίζεται ως εξής:

$$\rho_{xy} = \frac{Cov(X, Y)}{\sigma_x \sigma_y} = \frac{E(X - \mu_x) \cdot (Y - \mu_y)}{\sqrt{E(X - \mu_x)^2} \sqrt{E(Y - \mu_y)^2}} \quad (3.7.1.)$$

Όταν οι X και Y είναι διακριτές, τότε η συνδιακύμανση τους δηλαδή ο αριθμός της (3.7.1.) υπολογίζεται ως εξής:

$$Cov(X, Y) = \sum_i \sum_j (X_i - \mu_x) \cdot (Y_j - \mu_y) \cdot P(X_i, Y_j) \quad (3.7.2.)$$

ενώ οι τυπικές αποκλίσεις ως εξής:

$$\begin{aligned}
\sigma_x &= \sqrt{\sum_i (X_i - \mu_x)^2 P(X_i)} \\
\sigma_y &= \sqrt{\sum_j (Y_j - \mu_y)^2 P(Y_j)} \quad (3.7.3.)
\end{aligned}$$

Συγκρίνοντας τις (3.7.1.), (3.7.2.), (3.7.3.) με τους τύπους που εφαρμόσαμε για να υπολογίσουμε τον συντελεστή συσχέτισης ομαδοποιημένων παρατηρήσεων παρατηρούμε ότι αυτοί διαφέρουν μόνο στο ότι οι σχετικές συχνότητες F_{ij}/n έχουν αντικατασταθεί από τις πιθανότητες $P(X_i, Y_j)$ και οι \bar{X} , \bar{Y} από τις μ_x , μ_y .

Πράγματι, ο \hat{r}_{xy} είναι το δειγματικό ανάλογο του ρ_{xy} , όπως τα \bar{x} , \bar{y} , s^2x , s^2y είναι τα δειγματικά ανάλογα των παραμέτρων μ_x , μ_y , σ^2x , σ^2y αντίστοιχα.

Για να υπολογίσουμε τον \hat{r}_{xy} πρέπει να γνωρίζουμε την κοινή κατανομή πιθανοτήτων των X και Y κάτι που στην πράξη δεν συμβαίνει σχεδόν ποτέ. Συνήθως έχουμε ένα δείγμα η παρατηρήσεων των τυχαίων μεταβλητών X και Y για τις οποίες υπολογίζουμε τον συντελεστή συσχέτισης r_{xy} που στην περίπτωση αυτή ονομάζεται **δειγματικός συντελεστής συσχέτισης** ενώ ο ρ_{xy} ονομάζεται **θεωρητικός ή συντελεστής συσχέτισης πληθυσμού**.

Ο r_{xy} είναι εκτιμητής του ρ_{xy} του οποίου η αξιοπιστία εξαρτάται από το μέγεθος του δείγματος αλλά και την κατανομή πληθυσμού δηλαδή την κοινή κατανομή των X και Y . Όταν η κατανομή αυτή είναι η διμεταβλητή κανονική μπορούμε να προσδιορίσουμε την αξιοπιστία του r_{xy} , να υπολογίσουμε διαστήματα εμπιστοσύνης και να κάνουμε τους σχετικούς ελέγχους.

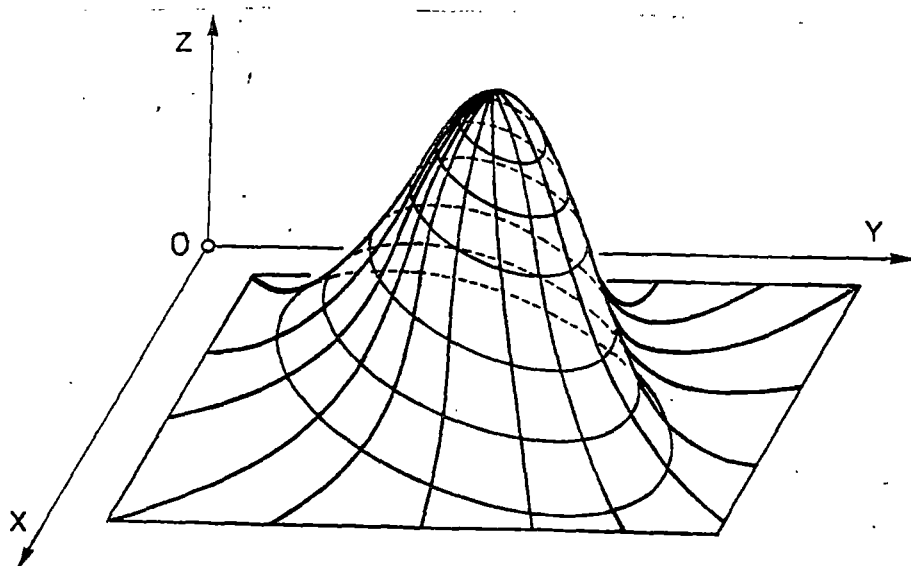
3.8. Η διμεταβλητή κανονική κατανομή.

Ένα μοντέλο για την κοινή κατανομή δύο συνεχών τυχαίων μεταβλητών X και Y είναι η **κοινή κανονική κατανομή**.

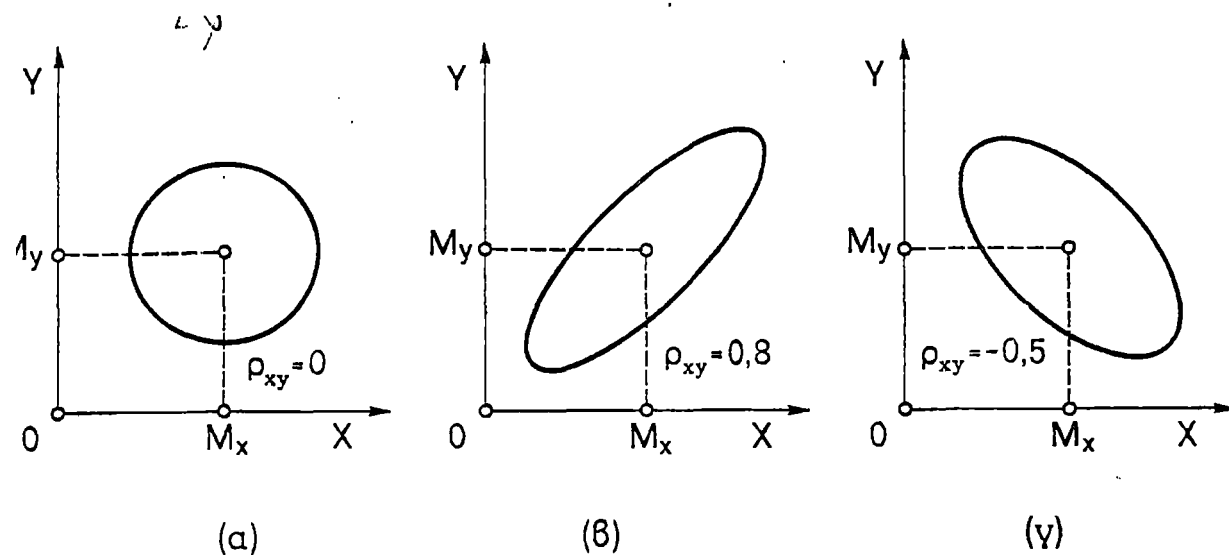
Οι παράμετροι της διμεταβλητής κανονικής κατανομής είναι οι μέσες τιμές μ_x, μ_y ,

οι διακυμάνσεις σ^2x , σ^2y των x και y αντίστοιχα και ο συντελεστής συσχέτισης ρ_{xy} .

Η γραφική παράσταση της συνάρτησης πυκνότητας πιθανότητας είναι η επιφάνεια ενός λόφου που υψώνεται πάνω στο επίπεδο xy όπως στο ακόλουθο σχήμα και έχει τα εξής χαρακτηριστικά :



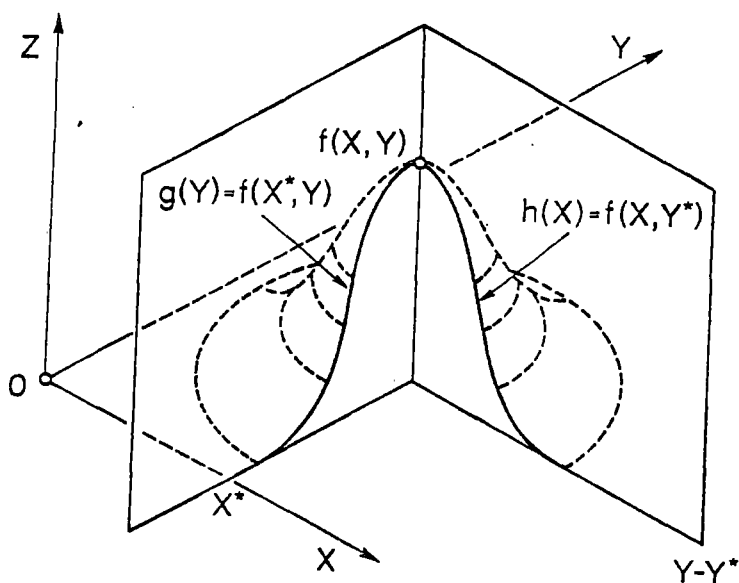
- Ο όγκος του χώρου που περικλείεται από την επιφάνεια αυτή και το επίπεδο XY ισούται με 1.
- Η κορυφή του "λόφου" βρίσκεται στο σημείο (μ_x, μ_y)
- Για δεδομένες τις μέσες τιμές μ_x, μ_y και τις διακυμάνσεις σ^2_x, σ^2_y , το σχήμα του "λόφου" εξαρτάται από τον συντελεστή συσχέτισης ρ_{xy} . Όταν $\rho_{xy}=0$, ο λόφος είναι συμμετρικά κυκλικός και η τομή του με το επίπεδο XY είναι ένας κύκλος (εφ. σχ. α). Όσο αυξάνει η απόλυτη τιμή του ρ_{xy} , ο λόφος πιέζεται από την μία πλευρά έτσι ώστε η τομή του με το επίπεδο XY είναι μια έλλειψη με κέντρο στο σημείο (μ_x, μ_y) (σχ. β, γ). Αν $\rho_{xy}>0$ η μεγάλη ακτίνα της έλλειψης έχει θετική κλίση διαφορετικά, έχει αρνητική. Όταν $\rho_{xy}=1$ ο λόφος δεν έχει καθόλου θετική ή αρνητική κλίση ανάλογα με το αν το πρόσημο του ρ_{xy} είναι (+) ή (-).



• Μπορούμε να κόψουμε το λόφο που σχηματίζει η διμεταβλητή κανονική συνάρτηση πυκνότητας πιθανότητας με ένα επίπεδο παράλληλο στο XY . Η τομή του βουνού με το επίπεδο αυτό ονομάζεται **ισοϋψής καμπύλη πιθανότητας** και είναι επίσης έλλειψη όπως και η τομή με το XY . Ο όγκος που περικλείεται από την επιφάνεια που ορίζει μια ισοϋψής καμπύλη και το υπερκείμενο μέρος του λόφου αντιστοιχεί σε ορισμένη πιθανότητα. Έτσι, αρκεί να γνωρίζουμε τις ισοϋψής καμπύλες

που αντιστοιχούν σε διαφορετικές πιθανότητες για να γνωρίζουμε την διμεταβλητή κατανομή πιθανοτήτων . Στο (σχ. δ) δίνονται οι ισοϋψής καμπύλες που αντιστοιχούν σε πιθανότητες 50% , 75% και 99% για διαφορετικές τιμές των σ^2_x , σ^2_y και ρ_{xy} ενώ $\mu_x=\mu_y=0$.

• Η τομή της επιφάνειας που ορίζει η διμεταβλητή κανονική συνάρτηση πιθανότητας με ένα επίπεδο κάθετο στο XY είναι μία κανονική καμπύλη . Η διαφορετικά , η δεσμευμένη κατανομή της Y (αντίστοιχα , της X) για οποιαδήποτε τιμή της X (αντίστοιχα , της Y) είναι κανονική (σχ. ε)

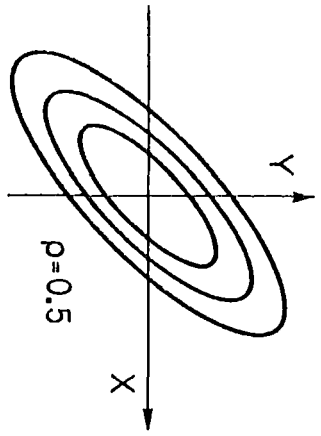


(ε)

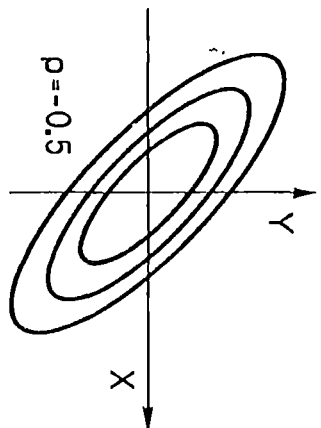
$$\sigma_X^2 = \sigma_Y^2$$

$$\sigma_X^2 > \sigma_Y^2$$

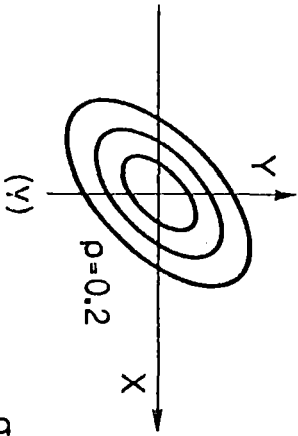
$$\sigma_X^2 < \sigma_Y^2$$



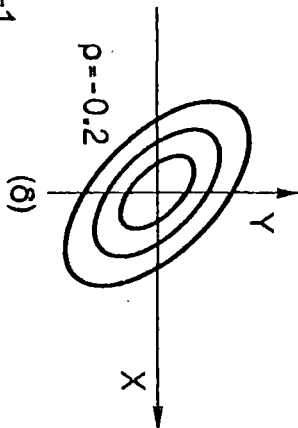
(a)



(b)



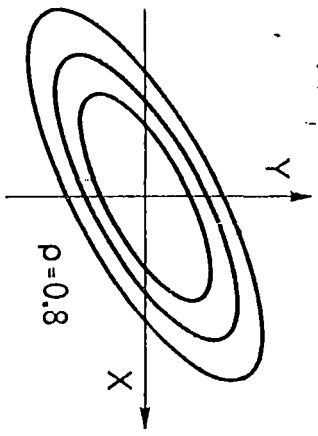
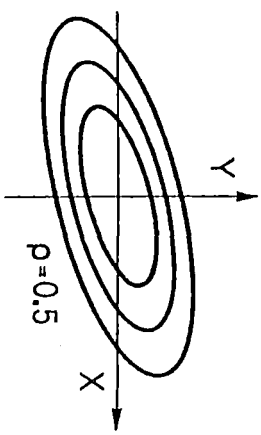
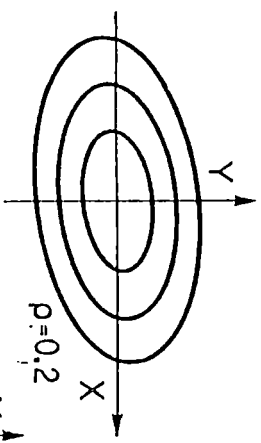
(γ)



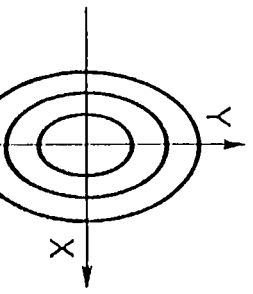
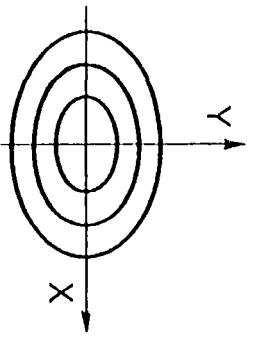
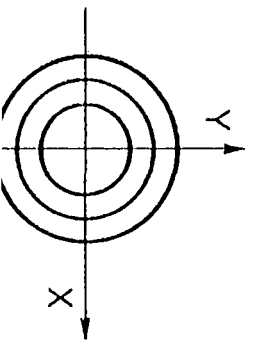
(δ)

$$\sigma_X^2 = \sigma_Y^2 = 1$$

(8)



$\sigma_X^2 = 4$ $\sigma_Y^2 = 1$



Βασικές Ιδιότητες της Διμεταβλητής Κανονικής Κατανομής

Η διμεταβλητή κανονική κατανομή έχει ορισμένες ιδιότητες οι οποίες θα φανούν χρήσιμες στην συνέχεια της ανάλυσης . Έτσι, όταν η κοινή κατανομή των X, Y είναι κανονική, ισχύουν τα εξής :

1. Οι περιθώριες κατανομές των X και Y είναι επίσης κανονικές δηλαδή $X \sim N(\mu_x, \sigma^2_x)$ και $Y \sim N(\mu_y, \sigma^2_y)$. Το αντίθετο όμως δεν ισχύει . Πολλές διμεταβλητές μη κανονικές κατανομές έχουν κανονικές περιθώριες .
2. Κάθε γραμμικός συνδυασμός των X, Y κατανέμεται επίσης κανονικά . Απ'όλες τις διμεταβλητές κατανομές μόνο η κανονική έχει αυτή την ιδιότητα .
3. Αν $\rho_{xy} = 0 \Leftrightarrow$ οι X, Y είναι ανεξάρτητες . Δηλαδή στην διμεταβλητή κανονική κατανομή και μόνον α'αυτήν οι έννοιες " ανεξάρτατες " και " ασυσχέτιστες " είναι ταυτόσημες .
4. Η συνάρτηση παλινδρόμησης της Y επί της X όπως και η συνάρτηση παλινδρόμησης της X επί της Y είναι γραμμικές .
5. Η δεσμευμένη κατανομή της Y για κάθε X_i έχει σταθερή διακύμανση και ίση με $\sigma^2_y (1 - \rho^2_{xy})$.

Πως θα ξέρουμε αν η κοινή κατανομή δύο συνεχών τυχαίων μεταβλητών είναι κανονική .

Ο έλεγχος προσαρμογής της διμεταβλητής κανονικής κατανομής σε ένα δείγμα η ζευγών παρατηρήσεων είναι πολύπλοκος και δεν μπορεί να γίνει χωρίς ηλεκτρονικό υπολογιστή .

Αρχικά γίνεται έλεγχος προσαρμογής της κανονικής κατανομής σε κάθε μια από τις περιθώριες κατανομές των X και Y . Αν αυτές δεν μπορούν να θεωρηθούν κανονικές τότε μπορούμε να απορρίψουμε την υπόθεση της κανονικότητας για την κοινή κατανομή πληθυσμού . Αν όμως οι περιθώριες κατανομές μπορούν να θεωρηθούν κανονικές τότε , σύμφωνα με την ιδιότητα 1 της διμεταβλητής κανονικότητας είναι πολύ πιθανόν αλλά όχι βέβαιο ότι η κοινή κατανομή είναι κανονική . Η διαδικασία αυτή αναφέρεται πληροφοριακά παρόλο που υπάρχει κίνδυνος να χρησιμοποιηθεί λανθασμένα .

3.9 Έλεγχος των υποθέσεων και διαστήματα εμπιστοσύνης .

Εστω $(X_1, Y_1), \dots, (X_n, Y_n)$ ένα τυχαίο δείγμα η παρατηρήσεων από μια διμεταβλητή κατανομή πιθανοτήτων και r_{xy} ο δειγματικός συντελεστής συσχέτισης .

Έχει αποδειχθεί ότι , ακόμη και όταν η κατανομή πληθυσμού είναι η διμεταβλητή κανονική , η κατανομή δειγματοληψίας του r_{xy} δεν είναι η κανονική. Εξαρτάται από τον συντελεστή συσχέτισης ρ_{xy} της κατανομής πληθυσμού και το μέγεθος του δείγματος n , προσεγγίζει Δε την κανονική μόνον για πάρα πολύ μεγάλες τιμές του n - μερικές εκατοντάδες . Αν όμως κάνουμε στον r_{xy} τον ακόλουθο μετασχηματισμό :

$$W = \frac{1}{2} \log_e \left(\frac{1+r_{xy}}{1-r_{xy}} \right)$$

τότε η κατανομή της W προσεγγίζει κανονική με μέση τιμή του W - μετασχηματισμού του ρ_{xy} και τυπική απόκλιση ίση με $1/\sqrt{n-3}$. Η προσέγγιση της κατανομής της W στην κανονική είναι ικανοποιητική ακόμη και για μικρά , σχετικά , δείγματα ($n \geq 15$) .

Επομένως , ο έλεγχος των υποθέσεων και ο υπολογισμός του διαστήματος εμπιστοσύνης για τον συντελεστή συσχέτισης πληθυσμού θα γίνει , αντίστοιχα , ως εξής :

Έλεγχος των υποθέσεων

α . Υποθέσεις που γίνονται δεκτές .

Η κατανομή πληθυσμού είναι η διμεταβλητή κανονική

β . Η μηδενική και η εναλλακτική υπόθεση .

$$H_0 : \rho_{xy} = \rho_0$$

$$H_1 : \rho_{xy} \neq \rho_0$$

γ . Το κριτήριο αποφάσεως .

Έστω W_1 και W_0 ο W - μετασχηματισμός του r_{xy} και του ρ_0 αντίστοιχα . Τους μετασχηματισμούς αυτούς βρίσκουμε απ'ευθείας στον πίνακα 7 στο τέλος της εργασίας .

$$\text{Τότε , } \alpha_n : |z_n| = \frac{|W_1 - W_0|}{1/\sqrt{n-3}} > z_{1-\alpha/2}$$

απορίπτουμε την H_0 στο επίπεδο σημαντικότητας α .

δ . Η απόφαση

Παράδειγμα

Επιφανής ανθρωπολόγος δήλωσε ότι ο συντελεστής ανάμεσα στο μήκος του ποδιού X μιας ενήλικης γυναίκας και το μήκος του χεριού της Y , είναι ίσος με $0,80$. Να ελεγχθεί η υπόθεση αυτή αν στις παρατηρήσεις τυχαίου δείγματος μεγέθους $n=307$ ενηλίκων γυναικών υπολογίσαμε συντελεστή συσχέτισης $r_{xy} = 0,72$ και μπορούμε να υποθέσουμε ότι η κοινή κατανομή των X και Y είναι η κανονική. Ο έλεγχος να γίνει στο επίπεδο σημαντικότητας $\alpha=0,05$.

α. Υποθέσεις που γίνονται δεκτές

Το δείγμα των 307 γυναικών είναι τυχαίο, η κοινή κατανομή των X και Y είναι η κανονική.

β. Η μηδενική και η εναλλακτική υπόθεση

$$H_0: \rho_{xy} = 0,80$$

$$H_1: \rho_{xy} \neq 0,80$$

γ. Το κριτήριο αποφάσεως και η απόφαση

Από τον πίνακα 7, στο τέλος της εργασίας, βρίσκουμε ότι η τιμή του συντελεστή συσχέτισης $0,80$ μετασχηματίζεται στην W - τιμή $W_0 = 1,0986$, ενώ η τιμή $0,72$ μετασχηματίζεται στην $W_1 = 0,9076$.

$$\text{Υπολογίζουμε: } |z_n| = \frac{|W_1 - W_0|}{1/\sqrt{n-3}} = \frac{|0,9076 - 1,0986|}{1/\sqrt{307-3}} = 3,33$$

Επειδή $|z_n| > z_{1-\alpha/2} = z_{0,975} = 1,96$ απορρίπτουμε την H_0 στο $\alpha=0,05$

Το διάστημα εμπιστοσύνης

Το διάστημα το οποίο με πιθανότητα $(1-\alpha)$ θα περιέχει την ρ_{xy} είναι το $\rho_l \leq \rho_{xy} \leq \rho_u$.

Το ρ_l είναι η τιμή του συντελεστή συσχέτισης που στον πίνακα στο τέλος της εργασίας, αντιστοιχεί στην $W_1 - z_{1-\alpha/2} \frac{\alpha}{2 \sqrt{1/n-3}}$ όπου W_1 είναι ο W -μετασχηματισμός του r_{xy} και $z_{1-\alpha/2}$ η τιμή της τυπικής κανονικής

κατανομής που αντιστοιχεί σε αθροιστική συχνότητα $1-\alpha/2$. Ομοίως το r_{xy} είναι η τιμή του r που αντιστοιχεί στην $W_{1-\alpha/2} + Z_{1-\alpha/2} \sqrt{1/n-3}$.

Ο Έλεγχος της υπόθεσης για μηδενική συσχέτιση

Όταν η κατανομή πληθυσμού είναι η διμεταβλητή κανονική και ισχύει $\rho_{xy}=0$ τότε, η τυχαία μεταβλητή

$$t = \frac{r_{xy}}{\sqrt{\frac{1-r_{xy}^2}{n-2}}}$$

ακολουθεί την κατανομή t- student με $v=n-2$ βαθμούς ελευθερίας. Επομένως ο έλεγχος για μηδενική συσχέτιση δύο τυχαίων μεταβλητών X και Y θα γίνει σταδιακά ως εξής :

α. Υποθέσεις που γίνονται δεκτές :

Η κοινή κατανομή πληθυσμού είναι η κανονική .

β. Η μηδενική και η εναλλακτική υπόθεση :

$$H_0: \rho_{xy} = 0$$

$$H_1: \rho_{xy} \neq 0$$

γ. Το κριτήριο απόφασης :

Αν r_{xy} η τιμή του δειγματικού συντελεστή συσχέτισης, n το μέγεθος του δείγματος και ισχύει :

$$|t_{n-2}| = \frac{|r_{xy}|}{\sqrt{\frac{1-r_{xy}^2}{n-2}}} > t_{n-2, \alpha/2}$$

τότε απορρίπτεται η H_0 στο επίπεδο σημαντικότητας α .

δ. Η απόφαση

Παράδειγμα

Σε τυχαίο δείγμα $n=4$ νεογέννητων μετρήσαμε το βάρος τους X και το ύψος τους Y και υπολογίσαμε δειγματικό συντελεστή συσχέτισης ίσο με $r_{xy}=0,76$. Γνωστού όντως ότι η κοινή κατανομή των μεταβλητών ανθρώπινο βάρος και ύψος, σε οποιαδήποτε ομάδα μπορεί να υποτεθεί κανονική, να ελεγχθεί η υπόθεση ότι οι δύο μεταβλητές είναι, στα νεογέννητα, ανεξάρτητες.

Αρκεί να ελέγξουμε αν οι γραμμικώς ανεξάρτητες δηλαδή αν $\rho_{xy}=0$. Ο ζητούμενος έλεγχος θα γίνει σταδιακά ως εξής:

α. Υποθέσεις που γίνονται δεκτές:

Η κοινή κατανομή των X και Y είναι η κανονική.

β. Η μηδενική και η εναλλακτική υπόθεση:

$$H_0: \rho_{xy} = 0$$

$$H_1: \rho_{xy} \neq 0$$

γ. Το κριτήριο απόφασης και η απόφαση:

Επειδή $t_{n-2, 9/2} = t_{12(0,025)} = 2,179$ και ισχύει:

$$|t_n| = \frac{|r_{xy}|}{\sqrt{\frac{1-r_{xy}^2}{n-2}}} = \frac{0,76}{\sqrt{\frac{1-(0,76)^2}{14-2}}} = 4,05 > 2,179$$

απορρίπτουμε την H_0 .

Σημείωση. Στην κανονική διμεταβλητή κατανομή οι έννοιες ασυσχέτιστες και ανεξάρτητες μεταβλητές είναι συνώνυμες.

Επομένως στην περίπτωση αυτή και μόνο ο προηγούμενος έλεγχος είναι και έλεγχος ανεξαρτησίας. Θα πρέπει πάντως να σημειωθεί ότι ο έλεγχος για μηδενική τιμή του ρ_{xy} μπορεί να εφαρμοστεί ακόμη και όταν η κοινή κατανομή των X, Y δεν μπορεί να υποτεθεί κανονική, αρκεί να ισχύει $n \geq 30$.

3.10 Παλινδρόμηση και Συσχέτιση

Ο συντελεστής συσχέτισης r_{xy} που υπολογίζεται για τα ζεύγη παρατηρήσεων $(X_1, Y_1), \dots, (X_n, Y_n)$ μετρά την τάση τους να συγκεντρώνονται γύρω από μια ευθεία. Ο r_{xy} είναι συμμετρικός ως προς τις X και Y με την έννοια ότι είναι αδιάφορο αν η ευθεία ορίζεται ως $L = \{(X, Y) : Y = b_0 + b_1 X\}$ ή ως $L' = \{(Y, X) : X = a_0 + a_1 Y\}$.

Στην εκτίμηση ενός μοντέλου παλινδρόμησης οι X και Y δεν είναι συμμετρικές. Έτσι, αν αλλάξουμε τη θέση των X και Y δεν θα πάρουμε τα ίδια αποτελέσματα. Είναι όμως ενδιαφέρον να δούμε πως ο συντελεστής συσχέτισης r_{xy} συνδέεται με τις παραμέτρους και, γενικότερα, τα στατιστικά ενός γραμμικού μοντέλου παλινδρόμησης που εκτιμήσαμε με βάση τις ίδιες παρατηρήσεις. Στη συνέχεια, η αναφορά θα γίνεται στο μοντέλο $\hat{Y} = \hat{b}_0 + \hat{b}_1 X$ λόγω συμμετρίας του r_{xy} όμως τα αποτελέσματα ισχύουν και για το μοντέλο $\hat{X} = \hat{a}_0 + \hat{a}_1 Y$.

Ο r_{xy} και ο συντελεστής παλινδρόμησης \hat{b}_1

Από τους τύπους

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

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

προκύπτει ότι

$$\hat{b}_1 = r_{xy} \frac{\sqrt{\sum (y_i - \bar{y})^2}}{\sqrt{\sum (x_i - \bar{x})^2}} \quad (3.10.1)$$

και επειδή

$$\sum (x_i - \bar{x})^2 = (n-1) s_x^2 \quad \text{και} \quad \sum (y_i - \bar{y})^2 = (n-1) s_y^2$$

η (3.10.1) γράφεται, ισοδύναμα ως εξής :

$$\hat{b}_1 = r_{xy} \frac{s_y}{s_x} \iff r_{xy} = \hat{b}_1 \frac{s_x}{s_y} \quad (3.10.2)$$

Από τη σχέση αυτή προκύπτει ότι ο r_{xy} έχει το ίδιο πρόσημο με το συντελεστή παλινδρόμησης \hat{b}_1 . Όταν η συμμεταβολή των X και Y είναι προς την ίδια κατεύθυνση, τότε και η κλίση \hat{b}_1 της ευθείας παλινδρόμησης είναι θετική. Διαφορετικά, είναι αρνητική.

Όταν θέλουμε να ελέγξουμε την υπόθεση ότι οι X και Y δεν συνδέονται γραμμικά, αρκεί να κάνουμε έλεγχο για μηδενική τιμή του b_1 είτε έλεγχο για μηδενική τιμή του r_{xy} . Οι δύο έλεγχοι είναι ισοδύναμοι.

Παλινδρόμηση προς το μέσο

Στην κεντραρισμένη μορφή του απλού μοντέλου παλινδρόμησης δηλαδή την

$$\hat{y} - \bar{y} = b_1(x - \bar{x})$$

αν αντικαταστήσουμε την $\hat{b}_1 = r_{xy} \frac{s_y}{s_x}$ παίρνουμε: $\hat{y} - \bar{y} = r_{xy} \frac{s_y}{s_x} (x - \bar{x})$

Διαιρώντας και τα δύο σκέλη με s_y προκύπτει η

$$\frac{\hat{y} - \bar{y}}{s_y} = r_{xy} \frac{x - \bar{x}}{s_x} \quad (3.10.3)$$

Για να εκτιμήσουμε το απλό μοντέλο παλινδρόμησης μπορούμε να εκφράσουμε τις παρατηρήσεις X_i , όπως και τις Y_i , σε τυπικές αποκλίσεις από τον αριθμητικό τους μέσο δηλαδή αντίστοιχα, ως:

$$z_{x_i} = \frac{X_i - \bar{x}}{s_x} \quad i = 1, \dots, n \quad (3.10.4)$$

και

$$z_{y_i} = \frac{Y_i - \bar{y}}{s_y}$$

Είναι εύκολο να δείχτεί ότι, τότε θα πάρουμε το ακόλουθο μοντέλο:

$$\hat{z}_y = r_{xy} z_x \quad (3.10.5)$$

δηλαδή το (3.10.3). Ο μετασχηματισμός (3.10.4) αναφέρεται ως τυποποίηση των παρατηρήσεων ή μετασχηματισμός τους σε Z -τιμές (z -scores). Είναι προφανές ότι οι Z_{x_i} , $i=1, \dots, n$, έχουν μέση τιμή μηδέν και διακύμανση 1 όπως οι Z_{y_i} , $i=1, \dots, n$.

Η εξίσωση (3.10.3) ή η ισοδύναμη της (3.10.5) περιγράφει ένα φαινόμενο το οποίο είναι γνωστό ως **παλινδρόμηση ή επαναφορά προς το μέσο** και μας λέει ότι σε μια τιμή X η οποία απέχει π.χ μία τυπική απόκλιση από το μέσο x εκτιμούμε ότι αντιστοιχεί μια τιμή Y η οποία θα απέχει από τον y όχι μια τυπική αλλά ένα ποσοστό της ίσο με r_{xy} . Η ονομασία οφείλεται στον F. Galton (1822-1911) ο οποίος μελετώντας την ανθρώπινη κληρονομικότητα παρατήρησε ότι οι πατεράδες οι οποίοι είναι πολύ ψηλότεροι από το μέσο όρο, τείνουν να έχουν γιους οι οποίοι είναι λιγότερο ψηλότεροι από τον αντίστοιχο μέσο όρο. Έτσι αν X = το ύψος του πατέρα και Y = το ύψος του γιου και $r_{xy}=0,6$ τότε οι πατεράδες που έχουν ύψος 2 τυπικές αποκλίσεις πάνω από το μέσο όρο έχουν γιους

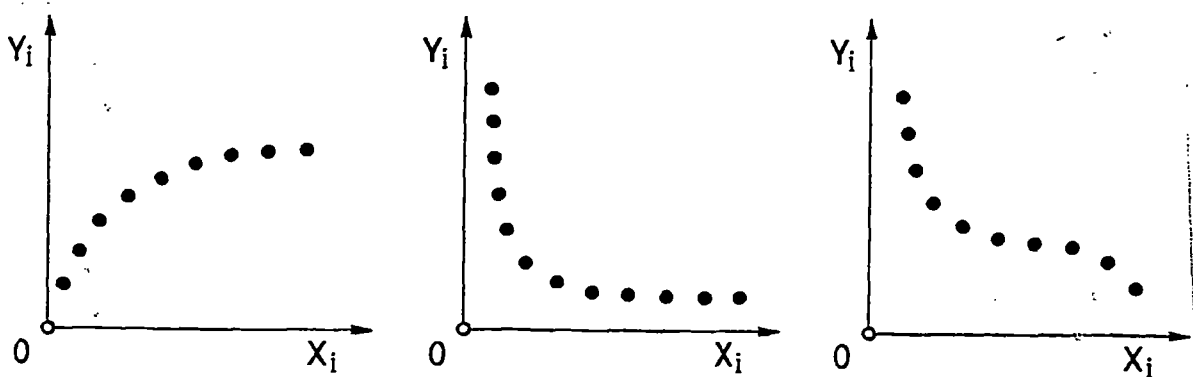
που κατά μέσο όρο έχουν ύψος $(0,6)^2=1,2$ τυπικές αποκλίσεις πάνω από το μέσο όρο .

Στην διεθνή βιβλιογραφία αναφέρονται και άλλα παραδείγματα παλινδρόμησης προς το μέσο . Έτσι π.χ σε δύο διαδοχικές μετρήσεις των ιδίων στοιχείων, εξαιρετικά μεγάλα σφάλματα μέτρησης την πρώτη φορά τείνουν να ακολουθούνται από λιγότερο υψηλά σφάλματα την επόμενη . Εξαιρετικά υψηλές επιδόσεις των αθλητών μία φορά τείνουν να ακολουθούνται από λιγότερο υψηλές επιδόσεις την επόμενη λόγω παλινδρόμησης προς το μέσο . Ομοίως εξαιρετικά χαμηλές επιδόσεις την μία φορά τείνουν να ακολουθούνται από λιγότερο χαμηλές επιδόσεις την άλλη .

Το φαινόμενο της παλινδρόμησης προς το μέσο θα πρέπει να παίρνεται υπόψη στα προβλήματα επαναληπτικών ελέγχων . Έτσι π.χ έστω ότι σε μία ομάδα παιδιών προσχολικής ηλικίας δίνεται ένα IQ test στα 4 χρόνια τους και ένα δεύτερο στα 5 . Τα αποτελέσματα των δύο τεστ φυσικά συνεχίζονται . Αν με βάση τα αποτελέσματα του πρώτου τεστ τα παιδιά με χαμηλό IQ επιλέγηκαν για συμπληρωματική εκπαίδευση η τάση για μεγαλύτερες αποδόσεις στο Β τεστ μπορεί λανθασμένα να αποδοθεί στην συμπληρωματική εκπαίδευση . Στην περίπτωση αυτή η σύγκριση πρέπει να γίνει συγκρίσιμη ομάδα ελέγχου (control group)

3.11 Ο συντελεστής συσχέτισης κατά τάξεις

Είναι δυνατόν οι παρατηρήσεις (X_i, Y_i) , $i = 1, \dots, n$ να μην προέρχονται από μία κανονική διμεταβλητή κατανομή και να συνδέονται ξεκάθαρα με μία σχέση που όμως δεν είναι γραμμική , όπως στα ακόλουθα σχήματα :



Σχέσεις όπως αυτές ονομάζονται μονοτικές . Γενικά ,

Δύο μεταβλητές X και Y συνδέονται με μία αυστηρά αύξουσα (φθίνουσα) μονοτική σχέση αν και μόνο αν , αυξάνουν (μειώνονται) οι τιμές της μίας με κάθε αύξηση των τιμών της άλλης . Η σχέση δεν είναι αυστηρά μονοτική αν στην αύξηση της μίας η άλλη μπορεί και να μην μεταβάλλεται .

Αν οι (X_i, Y_i) συνδέονται με μια μονοτική σχέση και στις τιμές X_i αντίστοιχα τις τάξεις W_i και ομοίως στις Y_i τις τάξεις V_i , $i=1, \dots, n$ τότε τα σημεία (W_i, V_i) θα βρίσκονται πάνω σε μια ευθεία . Επομένως , αν υπολογίσουμε τον συντελεστή συσχέτισης r_{wv} αυτός θα είναι ίσος με 1 ή -1 .

Ο συντελεστής συσχέτισης που υπολογίζεται για τις τάξεις η ζευγών παρατηρήσεων ονομάζεται **συντελεστής συσχέτισης κατά τάξεις ή συντελεστής συσχέτισης του Spearman** και συμβολίζεται με r_s . Δηλαδή έχουμε (ο τελεστής S είναι για όλες τις τιμές $i = 1, \dots, n$) :

$$r_s = - \frac{\sum (W_i - \bar{W})(V_i - \bar{V})}{\sqrt{\sum (W_i - \bar{W})^2} \sqrt{\sum (V_i - \bar{V})^2}} \quad (3.11.1)$$

Επειδή $1+2+\dots+n=n(n+1)/2$ και

$$\bar{V}^2 = \bar{W}^2 = \bar{V} * \bar{W} = \frac{n(n+1)^2}{4} \equiv C$$

μπορούμε να υπολογίσουμε εύκολα τον r_s υπολογίζοντας το C και στη συνέχεια αντικαθιστώντας στην (3.11.1) η οποία έτσι γίνεται :

$$r_s = \frac{\sum W_i V_i - C}{\sqrt{\sum W_i^2 - C} \sqrt{\sum V_i^2 - C}} \quad (3.11.2)$$

Αποδεικνύεται , ότι όταν δεν υπάρχουν επαναλαμβανόμενες παρατηρήσεις στα δεδομένα , οι τύποι (3.11.1) και (3.11.2) είναι ισοδύναμοι με τον υπολογιστικά ευκολότερο τύπο : $r_s = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$, όπου $d_i = w_i - v_i$, $i=1, \dots, n$

Χρησιμότητα του r_s

Ο συντελεστής συσχέτισης κατά τάξεις είναι ιδιαίτερα χρήσιμος στις κοινωνικές επιστήμες όπου οι παρατηρήσεις είναι πολλές φορές στην τακτική κλίμακα μέτρησης . Είναι συχνές οι περιπτώσεις όπου δύο άτομα αξιολογούν τα ίδια αντικείμενα είτε βαθμολογώντας τα σε ορισμένη κλίμακα είτε συνηθέστερα, κατατάσσοντας τα σε αύξουσα σειρά .

Στην πρώτη περίπτωση, συνήθως οι δυνατές τιμές της κλίμακας είναι λίγες ώστε να μην μπορούμε σε καμία περίπτωση να θεωρήσουμε

ότι οι μεταβλητές είναι συχνές και άρα να υποθέσουμε κανονική κατανομή πληθυσμού. Γι'αυτό θα αντιστοιχήσουμε στις παρατηρήσεις τις τάξεις και θα υπολογίσουμε τον συντελεστή r_s . Στη δεύτερη περίπτωση, οι παρατηρήσεις είναι τάξεις και ο μόνος συντελεστής που μπορεί να υπολογιστεί είναι ο r_s .

Όταν ο συντελεστής r_s ανάμεσα σε δύο βαθμολογίες ή σε δύο κατατάξεις είναι υψηλός, τότε μπορούμε να δεχτούμε ότι τα δύο άτομα που βαθμολογούν ή κατατάσσουν χρησιμοποιούν τα ίδια κριτήρια αξιολόγησης.

Σημειώνεται ότι όταν η κλίμακα μέτρησης μας επιτρέπει να υπολογίσουμε και τους δύο συντελεστές συσχέτισης τότε, συνήθως στους συνήθως στους δύο συντελεστές οφείλονται συνήθως σε ακραίες τιμές οι οποίες επηρεάζουν πολύ τον r_{xy} και λιγότερο τον r_s .

Παράδειγμα

Ζητήσαμε από δύο προπονητές να δουν το φιλμ του τελευταίου αγώνα μπάσκετ μεταξύ δύο ομάδων και να βαθμολογήσουν την απόδοση των 10 παικτών σε μία κλίμακα από το 1 ως το 20. Πήραμε τα ακόλουθα αποτελέσματα :

Προπονητής A X_i 4 8 20 12 14 15 18 20 10
 Προπονητής B Y_i 2 8 4 20 8 16 10 19 18 15

Για να υπολογίσουμε τον συντελεστή συσχέτισης κατά τάξεις αντιστοιχούμε τις τάξεις W_i στις παρατηρήσεις X_i και ομοίως τις τάξεις V_i στις παρατηρήσεις Y_i .

Σημειώνεται ότι όταν στις παρατηρήσεις υπάρχουν επαναλαμβανόμενες τιμές αντιστοιχούμε τον αριθμητικό μέσο των τάξεων τις οποίες θα έπρεπε να καταλάβουν.

W_i 1 2 3 9,5 5 6 7 8 9,5 4
 Y_i 1 3,5 2 10 3,5 7 5 9 8 6

Στη συνέχεια υπολογίζουμε :

$$C = n(n+1)^2/4 = 10 \cdot 11^2/4 = 302,5$$

$$\sum W_i^2 = \sum V_i^2 = 383,5$$

$$\sum V_i W_i = 375,5$$

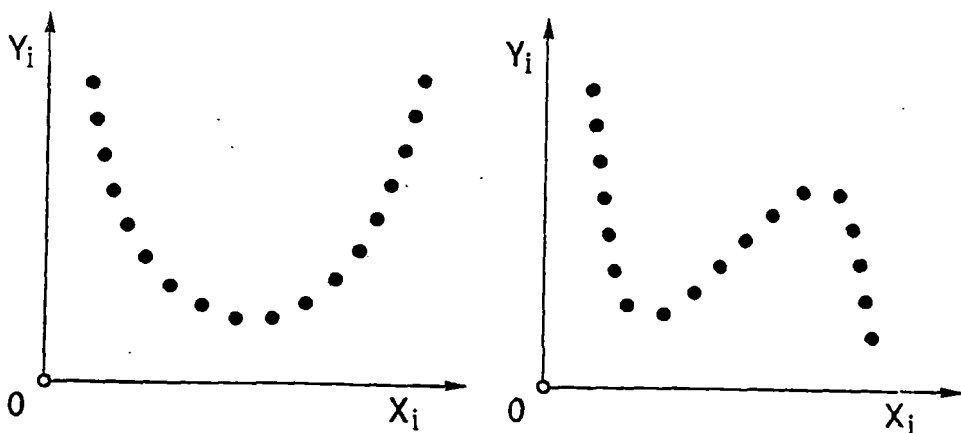
οπότε

$$r_s = \frac{\sum W_i V_i - C}{\sqrt{\sum W_i^2 - C} \sqrt{\sum V_i^2 - C}} = \frac{375,5 - 302,5}{383,5 - 302,5} = \frac{72,5}{81} = 0,895$$

που δείχνει ότι οι δύο προπονητές τείνουν να συμφωνούν στην αξιολόγηση των παικτών .

Έλεγχος των υποθέσεων.

Έστω ότι σε ορισμένο τυχαίο δείγμα η ζευγών παρατηρήσεων δύο μεταβλητών X και Y υπολογίσαμε το συντελεστή κατά τάξεις r_s και μας ενδιαφέρει ο έλεγχος σημαντικότητας του r_s . Η μηδενική υπόθεση, συνήθως εξειδικεύεται ως υπόθεση ανεξαρτησίας των X και Y παρόλο που ο σχετικός έλεγχος δεν ανιχνεύει μη μονοτονικές σχέσεις . Έτσι, αν οι παρατηρήσεις συνδέονται με μια μη μονοτονική σχέση όπως στα ακόλουθα σχήματα, θα υπολογίσουμε μια τιμή για τον r_s , η οποία δεν είναι στατιστικά σημαντική . Το διάγραμμα διασποράς μπορεί να είναι διαφωτιστικό και στην περίπτωση αυτή .



Αποδεικνύεται ότι όταν οι X και Y είναι ανεξάρτητες , τότε η κατανομή του στατικού

$$t = r_s \sqrt{\frac{n-2}{1-r_s^2}}$$

προσεγγίζει την κατανομή t - student με $v=n-2$ βαθμός ελευθερίας . Η προσέγγιση είναι πολύ καλή για $n \geq 18$. Στον πίνακα 6 στο τέλος της εργασίας δίνονται ποσοστιαία σημεία από την ακριβή κατανομή του r_s

για μικρές τιμές του n . Επομένως, ο έλεγχος σημαντικότητας του r_s θα γίνει σταδιακά ως εξής:

- Η μηδενική και η εναλλακτική υπόθεση.

H_0 : Οι X και Y είναι ανεξάρτητες

H_e : i) Οι X και Y συνδέονται με μια μονοτονική σχέση.

ii) Οι X και Y συνδέονται με μια αύξουσα μονοτονική σχέση.

iii) Οι X και Y συνδέονται με μια φθίνουσα μονοτονική σχέση.

- Το κριτήριο απόφασης.

Για $n \geq 18$ Υπολογίζουμε

$$t_n = r_s \sqrt{\frac{n-2}{1-r_s^2}}$$

οπότε, απορρίπτουμε την H_0 , στο επίπεδο α , για τα τρία είδη ελέγχου, αντίστοιχα, αν

i) $|t_n| > t_{n-2, 9/2}$

ii) $t_n > t_{n-2, \alpha}$

iii) $t_n < -t_{n-2, \alpha}$

3.12 Οι συντελεστές αυτοσυσχέτισης

Οι παρατηρήσεις μιας τυχαίας μεταβλητής που συλλέγονται σε διαδοχικά ισαπέχοντα σημεία του χρόνου ή του χώρου αποτελούν μια **χρονική σειρά**. Στην ανάλυση χρονικών σειρών μας ενδιαφέρει να ελέγξουμε αν υπάρχει αλληλεξάρτηση μεταξύ διαδοχικών παρατηρήσεων οπότε παραβιάζεται μια από τις υποθέσεις του τυχαίου δείγματος όπως ορίστηκε στο πρώτο κεφάλαιο. Για το σκοπό αυτό υπολογίζονται οι συντελεστές συσχέτισης μεταξύ διαδοχικών παρατηρήσεων που ονομάζονται συντελεστές αυτοσυσχέτισης ή **σειριακής συσχέτισης**.

Ειδικότερα, έστω X_1, X_2, \dots, X_n μια χρονική σειρά. Ο συντελεστής αυτοσυσχέτισης πρώτης τάξης ή συντελεστής σειριακής συσχέτισης σε υστέρηση 1, συμβολίζεται με r_1 και είναι ο συντελεστής συσχέτισης που υπολογίζεται στα $(n-1)$ ζεύγη παρατηρήσεων $(X_2, X_1), (X_3, X_2), \dots$

(X_n, X_{n-1}). Συνήθως χρησιμοποιείται ο ακόλουθος , ελαφρώς τροποποιημένος τύπος .

$$r_1 = \frac{\sum_{i=2}^n (x_i - \bar{x})(x_{i-1} - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (3.12.1)$$

Γενικά , ο συντελεστής αυτοσυσχέτισης κ τάξης συμβολίζεται με r_k και υπολογίζεται ως εξής :

$$r_k = \frac{\sum_{i=k+1}^n (x_i - \bar{x})(x_{i-k} - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (3.12.2)$$

Παράδειγμα

Θα υπολογίσουμε τους συντελεστές αυτοσυσχέτισης πρώτης και δεύτερης τάξης για τις καταθέσεις όψεως X_i (σε δ. . δρχ) της Ελλάδας των ετών 1977-86 .

Έτος	X_i	X_{i-1}	X_{i-2}
1977	32	-	-
1978	41	32	-
1979	51	41	32
1980	57	51	41
1981	73	57	51
1982	88	73	57
1983	105	88	73
1984	140	105	88
1985	168	140	105
1986	195	168	140

Πηγή : Μηνιαίος Στατιστικό Δελτίο της Τράπεζας της Ελλάδας ,
Οκτώβριος 1987.

Υπολογίζουμε : $\sum_{i=1}^{10} X_i = 950$ $\bar{X} = 95$

$$\sum_{i=2}^{10} (X_i - \bar{X})(X_{i-1} - \bar{X}) = 19405$$

$$\sum_{i=3}^{10} (X_i - \bar{X})(X_{i-2} - \bar{X}) = 10753$$

$$\sum_{i=1}^{10} (X_i - \bar{X})^2 = 28252$$

οπότε :

$$r_1 = \frac{\sum_{i=2}^{10} (X_i - \bar{X})(X_{i-1} - \bar{X})}{\sum_{i=1}^{10} (X_i - \bar{X})^2} = \frac{19405}{28252} = 0,687$$

$$r_2 = \frac{\sum_{i=3}^{10} (X_i - \bar{X})(X_{i-2} - \bar{X})}{\sum_{i=1}^{10} (X_i - \bar{X})^2} = \frac{10753}{28252} = 0,38$$

● Χρήσιμη είναι η ακόλουθη παρατήρηση : Αν σε δύο χρονικές σειρές υπολογίζουμε υψηλούς συντελεστές αυτοσυσχέτισης τότε είναι σχεδόν βέβαιο ότι και ο συντελεστής της μεταξύ τους συσχέτισης θα είναι υψηλός .

Έλεγχος σημαντικότητας

Ο συντελεστής αυτοσυσχέτισης r_k είναι δειγματικό ανάλογο του θεωρητικού συντελεστή αυτοσυσχέτισης ρ_k τάξης k ο οποίος συμβολίζεται με ρ_k και ορίζεται για την στοχαστική διαδικασία παραγωγής των παρατηρήσεων X_1, \dots, X_n .

Ο r_k είναι ένας εκτιμητής του ρ_k και , αποδεικνύεται ότι , όταν οι διαδοχικές παρατηρήσεις της τυχαίας μεταβλητής X είναι ασυσχέτιστες, δηλαδή ισχύει :

$$\rho_k = 0 \text{ για κάθε } k > 0$$

τότε η διακύμανση του r_k ισούται με $1/n$. Αν , επιπλέον , το n είναι μεγάλο τότε η κατανομή του r_k προσεγγίζει την κανονική . Ισοδύναμα , η κατανομή της $r_k \sqrt{n}$ προσεγγίζει την $N(0,1)$. Έτσι να τιμή του (r_k) μεγαλύτερη από το $2/\sqrt{n}$.. θα πρέπει να θεωρείται σημαντικά διαφορετική από το μηδέν .

Η γραφική παράσταση των r_k , $k=1,2,\dots$ στο διάστημα $\pm 2/\sqrt{n}$

ονομάζεται **συνάρτηση αυτοσυσχέτισης** και μας επιτρέπει να δούμε αμέσως για ποιές τιμές του k ο σχετικός συντελεστής είναι έξω από το διάστημα αυτό και άρα είναι στατιστικά σημαντικός . Η τυπική μορφή μιας τέτοιας συνάρτησης είναι όπως στο ακόλουθο σχήμα :

r_i			0		
1	*	o		o	-.405280
2		* o		o	-.257709
3		* o		o	-.229934
4		o *		o	-.055679
5		o *		o	-.070653
6		o *		o	.008952
7		o *		o	.050745
8		o *		o	-.004799
9		o *		o	-.075927
10		o *		o	.011173
11		o		o *	.319452
12	*	o		o	-.372311

ΚΕΦΑΛΑΙΟ ΤΕΤΑΡΤΟ

Το πολλάπλο μοντέλο παλινδρόμησης

4.1 Εισαγωγή

Το απλό μοντέλο του προηγούμενου κεφαλαίου γενικεύεται στο **πολλάπλο ή πολυμεταβλητό γραμμικό μοντέλο παλινδρόμησης**, στο οποίο η μέση τιμή της εξαρτημένης μεταβλητής Y είναι γραμμική συνάρτηση των ερμηνευτικών μεταβλητών X_1, \dots, X_k . Έτσι, σε κάθε δείγμα παρατηρήσεων της Y αντιστοιχεί το $n \times k$ σύστημα εξισώσεων :

$$Y_i = b_0 + b_1 X_{1i} + \dots + b_k X_{ki} + \varepsilon_i \quad i=1, \dots, n \quad (4.1.1)$$

όπου X_{ij} , $i=1, \dots, k$, $j=1, \dots, n$ είναι η i -τιμή της ερμηνευτικής μεταβλητής X_j και ε_i ο όρος σφάλματος που αντιστοιχεί στην παρατήρηση Y_i . Μπορούμε να γράψουμε την (4.1.1) ισοδύναμα, ως εξής :

$$Y_i = b_0 X_{0i} + b_1 X_{1i} + \dots + b_k X_{ki} + \varepsilon_i \quad i=1, \dots, n \quad (4.1.2)$$

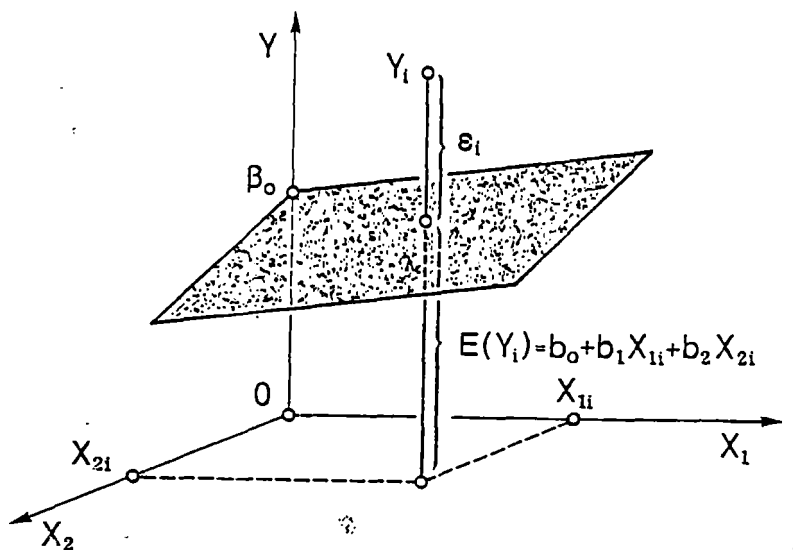
όπου $X_{0i} = 1$, $i=1, \dots, n$

Σε ένα πρακτικό πρόβλημα, υποθέτοντας ότι ισχύει το μοντέλο αυτό υποθέτουμε ότι σε επαναληπτικές μετρήσεις στις οποίες οι τιμές των ερμηνευτικών μεταβλητών παραμένουν σταθερές και ίσες με X_{1i}, \dots, X_{ki} αντιστοιχεί η τυχαία μεταβλητή Y_i . Η στοχασμική συμπεριφορά της Y_i προσδιορίζεται από την αντίστοιχη συμπεριφορά του όρου σφάλματος ε_i ο οποίος είναι μια τυχαία μεταβλητή.

Υποθέτοντας ότι $E(\varepsilon_i) = 0$ έχουμε :

$$E(Y_i) = b_0 + b_1 X_{1i} + \dots + b_k X_{ki}, \quad i=1, \dots, n$$

που είναι η συνάρτηση παλινδρόμησης της Y επί των X_1, X_2, \dots, X_k και ορίζει ένα επίπεδο στο χώρο των $k+1$ διαστάσεων. Στο επόμενο σχήμα δίνεται το γράφημα μιας γραμμικής συνάρτησης παλινδρόμησης με δυο ερμηνευτικές μεταβλητές.



Στην ανάλυση παλινδρόμησης η εξαρτημένη ονομάζεται και **μεταβλητή ανταπόκρισης**. Επειδή η συνάρτηση πολλαπλής παλινδρόμησης ορίζει ένα επίπεδο ή επιφάνεια, το γράφημα της ονομάζεται και **επιφάνεια ανταπόκρισης**.

Το διάνυσμα των τιμών (X_{1i}, \dots, X_{ki}) στις οποίες παρατηρείται η Y_i ονομάζεται το i - **σημείο σχεδιασμού**. Μια τιμή Y_i μαζί με το αντίστοιχο σημείο σχεδιασμού σημειώνονται ως $(X_{1i}, \dots, X_{ki}, Y_i)$ και αναφέρεται ως το i - **σημείο των δεδομένων ή περίπτωση**. Η εκτίμηση των αγνώστων παραμέτρων b_0, b_1, \dots, b_k του μοντέλου θα γίνει από ένα σύνολο n ($n > k$) περιπτώσεων που θα αναφέρεται ως **δείγμα δεδομένων**.

4.2 Η εκτίμηση των παραμέτρων με τη μέθοδο των ελαχίστων τετραγώνων.

Έστω ότι παρατηρήσαμε την εξαρτημένη μεταβλητή Y σε n σημεία σχεδιασμού (X_{1i}, \dots, X_{ki}) όχι όλα ίδια μεταξύ τους και διαθέτουμε τις n περιπτώσεις $(X_{1i}, \dots, X_{ki}, Y_i)$ $i=1, \dots, n$ για τις οποίες υποθέτουμε ότι ισχύει:

$$Y_i = b_0 + b_1 X_{1i} + \dots + b_k X_{ki} \quad i=1, \dots, n$$

Όπως και στο μοντέλο, με την μέθοδο των ελαχίστων τετραγώνων ζητούμε να βρούμε τις τιμές b_0, b_1, \dots, b_k των παραμέτρων b_0, b_1, \dots, b_k αντίστοιχα, οι οποίες για τις n περιπτώσεις των δεδομένων μας ελαχιστοποιούν τη συνάρτηση

$$S = S(b_0, b_1, \dots, b_k) = \sum_{i=1}^n \hat{\epsilon}_i^2 \quad (4.2.1)$$

$$= \sum_{i=1}^n (Y_i - b_0 - b_1 X_{1i} - \dots - b_k X_{ki})^2$$

Οι τιμές $\hat{b}_0, \hat{b}_1, \dots, \hat{b}_k$ που ελαχιστοποιούν την $S(b_0, b_1, \dots, b_k)$ ικανοποιούν το σύστημα εξισώσεων που προκύπτει όταν θέτουμε τις $k+1$ μερικές παραγωγούς $\frac{\partial S}{\partial b_j}$, $j=0, \dots, k$ ίσες με μηδέν.

Έχουμε δηλαδή το σύστημα :

$$\begin{aligned} \frac{\partial S}{\partial b_0} &= -2 \sum [Y_i - (\hat{b}_0 + \hat{b}_1 X_{1i} + \hat{b}_2 X_{2i} + \dots + \hat{b}_k X_{ki})] = 0 \\ \frac{\partial S}{\partial b_1} &= -2 \sum [X_{1i} [Y_i - (\hat{b}_0 + \hat{b}_1 X_{1i} + \hat{b}_2 X_{2i} + \dots + \hat{b}_k X_{ki})]] = 0 \\ &\vdots \\ \frac{\partial S}{\partial b_k} &= -2 \sum X_{ki} [Y_i - (\hat{b}_0 + \hat{b}_1 X_{1i} + \hat{b}_2 X_{2i} + \dots + \hat{b}_k X_{ki})] = 0 \end{aligned} \quad (4.2.2)$$

και ο τελεστής \sum είναι για $i=1, \dots, n$. Αν πολλαπλασιάσουμε με $-1/2$ κάθε εξίσωση του (4.2.2) το σύστημα αυτό γράφεται ισοδύναμα ως :

$$\begin{aligned} \hat{b}_0 + \hat{b}_1 \sum X_{1i} + \hat{b}_2 \sum X_{2i} + \dots + \hat{b}_k \sum X_{ki} &= \sum Y_i \\ \hat{b}_0 \sum X_{1i} + \hat{b}_1 \sum X_{1i}^2 + \hat{b}_2 \sum X_{1i} X_{2i} + \dots + \hat{b}_k \sum X_{1i} X_{ki} &= \sum X_{1i} Y_i \\ &\vdots \\ \hat{b}_0 \sum X_{ki} + \hat{b}_1 \sum X_{ki} X_{1i} + \hat{b}_2 \sum X_{ki} X_{2i} + \dots + \hat{b}_k \sum X_{ki}^2 &= \sum X_{ki} Y_i \end{aligned} \quad (4.2.3)$$

Το (4.2.3) είναι η γενική μορφή των κανονικών εξισώσεων για $k \geq 1$. Για $k=1$ είδαμε στο δεύτερο κεφάλαιο ότι είναι οι εξής (οι δέκτες παραλείπονται):

$$\begin{aligned} n \hat{b}_0 + \hat{b}_1 \sum X_1 &= \sum Y \\ \hat{b}_0 \sum X_1 + \hat{b}_1 \sum X_1^2 &= \sum X_1 Y \end{aligned}$$

ενώ για $k=2$ οι κανονικές εξισώσεις είναι οι εξής :

$$\begin{aligned} n \hat{b}_0 + \hat{b}_1 \sum X_1 + \hat{b}_2 \sum X_2 &= \sum Y \\ \sum X_1 \hat{b}_0 + \hat{b}_1 \sum X_1^2 + \hat{b}_2 \sum X_1 X_2 &= \sum X_1 Y \\ \sum X_2 \hat{b}_0 + \hat{b}_1 \sum X_2 X_1 + \hat{b}_2 \sum X_2^2 &= \sum X_2 Y \end{aligned}$$

Έτσι,

$$X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{k1} & X_{k2} & \dots & X_{kn} \end{bmatrix} \quad \hat{b} = \begin{bmatrix} \hat{b}_0 \\ \hat{b}_1 \\ \vdots \\ \hat{b}_k \end{bmatrix}$$

γράφουμε το σύστημα των κανονικών εξισώσεων συμπαγώς ως εξής :

$$\mathbf{X}^T \mathbf{Y} = (\mathbf{X}^T \mathbf{X}) \hat{\mathbf{b}}$$

(kxn)(nx1) (kxn)(nxk)(kx1)

και τη λύση του ως εξής :

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (4.2.4)$$

όπου $(\mathbf{X}^T \mathbf{X})^{-1}$ η αντίστροφη του $\mathbf{X}^T \mathbf{X}$. Να σημειωθεί ότι η λύση του συστήματος των κανονικών εξισώσεων , δηλαδή το διάνυσμα που ορίζεται από την (4.2.4) υπάρχει αν και μόνον αν η $(\mathbf{X}^T \mathbf{X})$ είναι μη ιδιάζουσα δηλαδή έχει αντίστροφη .Αυτό εξασφαλίζεται όταν οι $k+1$ στήλες της X είναι γραμμικώς ανεξάρτητες .

Η λύση αυτή οδηγεί σε μια προφανή διαδικασία υπολογισμού του διανύσματος $\hat{\mathbf{b}}$. Υπολογίζουμε αρχικά τη $\mathbf{X}^T \mathbf{X}$, την αντίστροφη της $(\mathbf{X}^T \mathbf{X})^{-1}$ και το διάνυσμα $\mathbf{X}^T \mathbf{Y}$. Στη συνέχεια υπολογίζουμε το $\hat{\mathbf{b}}$ πολλαπλασιάζοντας την $(\mathbf{X}^T \mathbf{X})^{-1}$ με το διάνυσμα $\mathbf{X}^T \mathbf{Y}$.

Όταν τα n και κυρίως το k είναι μικρά ($n \leq 15$, $k \leq 4$) αυτή η διαδικασία είναι εύκολη , και εκτός από το μέρος της αντίστροφης της $\mathbf{X}^T \mathbf{X}$, κατανοητή για το μέσο αναγνώστη . Η διαδικασία αυτή πάντως έχει σοβαρά υπολογιστικά μειονεκτήματα όπως αριθμητική αστάθεια και αριθμητικά σφάλματα . Τα γνωστά στατιστικά προγράμματα H/Y χρησιμοποιούν άλλες μεθόδους υπολογισμού του $\hat{\mathbf{b}}$. Η $(\mathbf{X}^T \mathbf{X})^{-1}$ πάντως είναι πολύ σημαντική στην ανάλυση παλινδρόμησης και υπεισέρχεται σε τύπους που ορίζουν σημαντικές παραμέτρους του εκτιμηθέντος μοντέλου.

4.3 Οι συντελεστές μερικής παλινδρόμησης

Στην εξίσωση παλινδρόμησης

$$Y_i = b_0 + b_1 X_{1i} + \dots + b_k X_{ki} + \varepsilon_i$$

κάθε Παράμετρος b_j , $j=1, \dots, k$ ονομάζεται συντελεστής μερικής ή καθαρής παλινδρόμησης της X_j , $j=1, \dots, k$ και μετρά την αναμενόμενη ή μέση μεταβολή της εξαρτημένης μεταβλητής όταν η τιμή της ερμηνευτικής μεταβλητής μεταβληθεί κατά μια μονάδα και οι τιμές των υπολοίπων παραμείνουν σταθερές . Έτσι π.χ αν η τιμή της X_k αυξηθεί κατά μία μονάδα , δηλαδή γίνει $X'_{ki} = X_{ki} + 1$ τότε ισχύει :

$$\begin{aligned}
 E(Y_i') &= b_0 + b_1 X_{1i} + \dots + b_k (X_{ki} + 1) \\
 &= b_0 + b_1 X_{1i} + \dots + b_k X_{ki} + b_k \\
 &= E(Y_i) + b_k
 \end{aligned}$$

Η ερμηνεία των εκτιμήσεων b_j , $j=1, \dots, k$ είναι ανάλογη.

Παράδειγμα

Για την πίεση Y ενός ενήλικα εκτιμήθηκε το ακόλουθο μοντέλο παλινδρόμησης με 20 παρατηρήσεις που έγιναν σε ισάριθμους ενήλικες τους οποίους επιλέξαμε με βάση την ηλικία τους X_1 , σε χρόνια και το βάρος τους X_2 , σε κιλά, από τα ενήλικα άτομα ίσου περίπου ύψους :

$$\hat{Y} = -65,10 + 0,425 X_1 + 1,077 X_2$$

Οι συντελεστές $\hat{b}_1 = 0,425$ και $\hat{b}_2 = 1,077$ ερμηνεύονται ως εξής : Στα άτομα ίδιου βάρους, κάθε επιπλέον χρόνο ηλικίας εκτιμούμε ότι αυξάνει την πίεση του αίματος κατά 0,425 μονάδες, ενώ τα άτομα ίδιας ηλικίας κάθε επί πλέον κιλό βάρους, αυξάνει την πίεση κατά 1,077 μονάδες. Η τιμή του σταθερού όρου δεν έχει κανένα νόημα, αφού είναι αδύνατον να υπάρχει στα δεδομένα μας παρατήρηση στο σημείο ($X_1 = 0$, $X_2 = 0$).

Μια άλλη διαδικασία υπολογισμού των \hat{b}_j

Από τον ορισμό του συντελεστή μερικής παλινδρόμησης προκύπτει ότι στην εκτίμηση π.χ του συντελεστή \hat{b}_j θα πρέπει να υπεισέρχονται μόνον οι μεταβολές της X_j . Αυτό σημαίνει

ότι αν η X_j σχετίζεται με τις υπόλοιπες θα πρέπει να υπεισέρχεται μόνον το μέρος της μεταβλητότητας της X_j το οποίο δεν περιέχεται στις υπόλοιπες. Αποδεικνύεται ότι η εκτίμηση ελαχίστων τετραγώνων \hat{b}_j του συντελεστή b_j μετρά την μεταβολή στην εξαρτημένη από μια μοναδιαία μεταβολή της X_j όταν οι γραμμικές επιδράσεις των υπολοίπων ερμηνευτικών μεταβλητών έχουν αφαιρεθεί και από την X_j και από την Ψ . Έτσι π.χ στο μοντέλο με $k=2$ ερμηνευτικές μεταβλητές δηλαδή το

$$Y_i = b_0 + b_1 X_{1i} + b_2 X_{2i} + \varepsilon_i \quad i=1, \dots, n \quad (4.3.1)$$

μπορούμε να εκτιμήσουμε έστω το συντελεστή b_1 ως εξής (ο όρος εκτίμηση αναφέρεται αποκλειστικά στην εκτίμηση ελαχίστων τετραγώνων)

- εκτιμούμε την παλινδρόμηση της Ψ επί της X_1 και έστω $\hat{\varepsilon}_{12}$ τα κατάλοιπα (οι δείκτες για ευκολία παραλείπονται).
- εκτιμούμε την παλινδρόμηση της X_1 επί της X_2 και έστω $\hat{\varepsilon}_{12}$ τα κατάλοιπα.
- εκτιμούμε την παλινδρόμηση των $\hat{\varepsilon}_{12}$ επί των $\hat{\varepsilon}_{12}$

επομένως έχουμε :

$$\hat{b}_1 = \frac{\sum \hat{\varepsilon}_{12} \hat{\varepsilon}_{12}}{\sum \hat{\varepsilon}_{12}^2} \quad (4.3.2)$$

Η τιμή του b_1 που θα παρούμε με τον τρόπο αυτό είναι ίδια με την τιμή που θα πάρουμε αν εφαρμόσουμε τη μέθοδο των ελαχίστων τετραγώνων απ'ευθείας στο μοντέλο (4.3.1).

Αποδεικνύεται εύκολα ότι η σχέση αυτή είναι ισοδύναμη με την:

$$\hat{b}_1 = \frac{\hat{b}_{y_1} - \hat{b}_{y_2} \hat{b}_{12}}{1 - r_{12}^2} \quad (4.3.3)$$

όπου \hat{b}_{y_1} και \hat{b}_{y_2} η εκτίμηση του συντελεστή απλής παλινδρόμησης της Ψ επί της X_1 και Ψ επί της X_2 αντίστοιχα \hat{b}_{12} η εκτίμηση του συντελεστή παλινδρόμησης της X_1 επί της X_2 και r_{12} ο συντελεστής συσχέτισης των παρατηρήσεων X_{1i} και X_{2i} $i=1, \dots, n$. Από τη σχέση (4.3.3) προκύπτουν ορισμένα ενδιαφέροντα συμπεράσματα :

Αν $r_{12} = 0$ οπότε και $\hat{b}_{12} = 0$ τότε $\hat{b}_1 = \hat{b}_{y_1}$

Αν $r_{12} = 1$ τότε \hat{b}_1 είναι απροσδιόριστο.

Αν $r_{12} \cong 1$ οπότε και $\hat{b}_{21} \cong 1$ τότε $\hat{b}_{y_1} \cong \hat{b}_{y_2}$ και το \hat{b}_1 προσδιορίζεται από την αναλογία δύο πολύ μικρών αριθμών. Στην περίπτωση αυτή μικρές μεταβολές στους όρους $\hat{b}_{y_1}, \hat{b}_{y_2}, \hat{b}_{21}, r_{12}^2$ που υπεισέρχονται στον ορισμό του \hat{b}_1 θα είχε ως αποτέλεσμα μεγάλες μεταβολές στην τιμή του \hat{b}_1 . Δηλαδή το \hat{b}_1 είναι αναξιόπιστος εκτιμητής του b_1 . Σημειώνεται ότι όταν $r_{12} \cong 1$ οι X_1, X_2 ονομάζονται ατελώς συγγραμικές.

Συμπεραίνουμε ότι η τιμή του \hat{b}_j που υπολογίζουμε στο πολλαπλό μοντέλο θα διαφέρει, γενικά, από την αντίστοιχη τιμή του στην απλή παλινδρόμηση της Ψ επί της X_j .

4.4 Τα κατάλοιπα της εκτίμησης

Όπως και στο απλό μοντέλο, ορίζουμε τα κατάλοιπα ή σφάλματα εκτίμησης ως τις διαφορές των εκτιμήσεων $\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_{1i} + \dots + \hat{b}_k x_{ki}$ από τις παρατηρήσεις y_i , δηλαδή

$$\begin{aligned}\hat{\varepsilon}_i &= Y_i - \hat{Y}_i \\ &= Y_i - \hat{b}_0 - \hat{b}_1 X_{1i} - \dots - \hat{b}_k X_{ki} \quad i=1, \dots, n\end{aligned}$$

Αν ορίσουμε τα διανύσματα

$$\hat{\varepsilon} = \begin{bmatrix} \hat{\varepsilon}_1 \\ \vdots \\ \hat{\varepsilon}_n \end{bmatrix} \quad \hat{Y} = \begin{bmatrix} \hat{Y}_1 \\ \vdots \\ \hat{Y}_n \end{bmatrix} \quad \hat{b} = \begin{bmatrix} \hat{b}_0 \\ \hat{b}_1 \\ \vdots \\ \hat{b}_k \end{bmatrix}$$

τότε ισχύει :

$$\hat{\varepsilon} = Y - \hat{Y} = Y - X\hat{b}$$

Οι ιδιότητες των καταλοίπων ελαχίστων τετραγώνων τις οποίες είδαμε στο απλό μοντέλο γενικεύονται ως εξής :

- Το διάνυσμα των καταλοίπων είναι ορθογώνιο στο διάνυσμα των τιμών της κάθε ερμηνευτικής μεταβλητής x_j , $j=1, \dots, k$ και άρα και στη μήτρα σχεδιασμού. Δηλαδή ισχύει :

$$\begin{aligned}\sum_{i=1}^n X_{ji} \hat{\varepsilon}_i &= 0 \quad j=1, \dots, k \quad (4.4.1) \\ \iff X^T \hat{\varepsilon} &= 0\end{aligned}$$

Πράγματι έχουμε :

$$\begin{aligned}X^T \hat{\varepsilon} &= X^T (Y - X\hat{b}) = X^T Y - X^T X \hat{b} \\ &= X^T Y - X^T X (X^T X)^{-1} X^T Y \\ &\equiv X^T Y - X^T Y = 0\end{aligned}$$

Έτσι, αν η πρώτη στήλη της μήτρας X είναι το μοναδιαίο διάνυσμα ισχύει

$$\sum_{i=1}^n \hat{\varepsilon}_i = 0$$

- Το διάνυσμα των καταλοίπων είναι ορθογώνιο στο διάνυσμα των εκτιμήσεων \hat{y} δηλαδή ισχύει :

$$\hat{Y}^T \hat{\varepsilon} = \sum_i \hat{Y}_i \hat{\varepsilon}_i = 0 \quad (4.4.2)$$

Πράγματι από την (4.4.1) έχουμε :

$$\hat{Y}^T \hat{\varepsilon} = \hat{b}^T X^T \hat{\varepsilon} = 0$$

Χρήσιμη είναι η εναλλακτική έκφραση του διανύσματος των καταλοίπων

$$\begin{aligned} \hat{\varepsilon} &= Y - \hat{Y} = Y - X \hat{b} \\ &= Y - X (X^T X)^{-1} X^T Y \\ &= [I - X (X^T X)^{-1} X^T] Y \quad (4.4.3) \\ &= [I - P] Y \\ &= H Y \end{aligned}$$

όπου $H=I-P$ και $P=X(X^T X)^{-1} X^T$

Είναι εύκολο να διαπιστωθεί ότι οι μήτρες H και P είναι συμμετρικές, δηλαδή ισχύει $H=H^T$ και $P=P^T$ και ταυτοτικές δηλαδή ισχύει $H H=H$ και $P P=P$. Εξάλλου, ισχύει

$$H X = (I - P) X = X - X (X^T X)^{-1} X^T X = X - X = 0$$

οπότε η (4.4.3) μπορεί να γραφεί ως εξής :

$$\hat{\xi} = HY = H(Xb + \varepsilon) = HXb + H\varepsilon = H\varepsilon \quad (4.4.4)$$

Σημείωση : Επειδή ισχύει $PY = Xb = Y$ δηλαδή η μήτρα P μετασχηματίζει το διάνυσμα των παρατηρήσεων στο διάνυσμα των προβλέψεων, η P ονομάζεται **μήτρα πρόβλεψης**. Εξάλλου, επειδή $HY = \varepsilon$ η μήτρα H ονομάζεται **μήτρα των καταλοίπων** (επειδή μετασχηματίζει τις παρατηρήσεις y στα κατάλοιπα ε).

Η μήτρα διακυμάνσεων - συνδιακυμάνσεων των καταλοίπων ορίζεται ως εξής :

$$E(\hat{\xi}\hat{\xi}^T) = \begin{bmatrix} \text{Var}(\hat{\varepsilon}_1) & \text{Cov}(\hat{\varepsilon}_1, \hat{\varepsilon}_2) & \dots & \text{Cov}(\hat{\varepsilon}_1, \hat{\varepsilon}_n) \\ \text{Cov}(\hat{\varepsilon}_1, \hat{\varepsilon}_2) & \text{Var}(\hat{\varepsilon}_2) & \dots & \text{Cov}(\hat{\varepsilon}_2, \hat{\varepsilon}_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\hat{\varepsilon}_1, \hat{\varepsilon}_n) & \text{Cov}(\hat{\varepsilon}_2, \hat{\varepsilon}_n) & \dots & \text{Var}(\hat{\varepsilon}_n) \end{bmatrix}$$

και από την (4.4.4) προκύπτει ότι

$$\begin{aligned} E(\hat{\xi}\hat{\xi}^T) &= E(H\varepsilon\varepsilon^T H^T) \\ &= HE(\varepsilon\varepsilon^T)H^T \\ &= \sigma^2 H H^T \\ &= \sigma^2 H \text{ αφού } n \text{ H είναι ταυτοτική} \end{aligned}$$

Προκύπτει επομένως ότι :

$$\text{Var}(\hat{\varepsilon}_i) = \sigma^2 h_{ii} = \sigma^2(1 - p_{ii}) \quad i = 1, \dots, n \quad (4.4.5)$$

$$\text{Cov}(\hat{\varepsilon}_i, \hat{\varepsilon}_j) = \sigma^2 h_{ij} = \sigma^2(1 - p_{ij}) \quad i = 1, \dots, n \quad j = 1, \dots, n$$

όπου h_{ii} , p_{ii} το i - διαγώνιο στοιχείο της μήτρας H και P , αντίστοιχα. Επειδή η διακύμανση είναι θετική θα πρέπει να ισχύει $p_{ii} < 1$, $i=1, \dots, n$. Από τις (4.4.5) προκύπτουν τα εξής ενδιαφέροντα συμπεράσματα :

- Η διακύμανση των καταλοίπων είναι μικρότερη από την διακύμανση των διαταρακτικών όρων και εξαρτάται από την μήτρα X .

- Η συνδιακύμανση των καταλοίπων δεν είναι μηδέν ακόμη και όταν οι διαταρακτικοί όροι έχουν μηδενική συνδιακύμανση.

4.5 Ιδιότητες των εκτιμητών ελαχίστων τετραγώνων

Εστω \hat{b} το διάνυσμα των εκτιμητών ελαχίστων τετραγώνων δηλαδή

$$b = (X^T X)^{-1} X^T y$$

Αποδεικνύεται ότι, κάτω από τις συνθήκες Gauss - Markov το διάνυσμα των εκτιμητών b είναι **Γραμμικό** δηλαδή γραμμική συνάρτηση του τυχαίου διανύσματος $y = (y_1, y_2, \dots, y_n)$ και **Αμερόληπτο** δηλαδή ισχύει

$$E(\hat{b}) = b$$

Για να προσδιορίσουμε την μεταβλητότητα του \hat{b} , σημειώνουμε, κατ' αρχήν τη μήτρα.

$$(\hat{b} - b)(\hat{b} - b)^T = \begin{bmatrix} (\hat{b}_0 - b_0)^2 & (\hat{b}_0 - b_0)(\hat{b}_1 - b_1) & \dots & (\hat{b}_0 - b_0)(\hat{b}_k - b_k) \\ (\hat{b}_0 - b_0)(\hat{b}_1 - b_1) & (\hat{b}_1 - b_1)^2 & \dots & (\hat{b}_1 - b_1)(\hat{b}_k - b_k) \\ \vdots & \vdots & \ddots & \vdots \\ (\hat{b}_0 - b_0)(\hat{b}_k - b_k) & (\hat{b}_1 - b_1)(\hat{b}_k - b_k) & \dots & (\hat{b}_k - b_k)^2 \end{bmatrix}$$

Και στη συνέχεια την :

$$\text{Var}(\hat{b}) = E(\hat{b} - b)(\hat{b} - b)^T = \begin{bmatrix} \text{Var}(\hat{b}_0) & \text{Cov}(\hat{b}_0, \hat{b}_1) & \dots & \text{Cov}(\hat{b}_0, \hat{b}_k) \\ \text{Cov}(\hat{b}_0, \hat{b}_1) & \text{Var}(\hat{b}_1) & \dots & \text{Cov}(\hat{b}_1, \hat{b}_k) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\hat{b}_0, \hat{b}_k) & \text{Cov}(\hat{b}_1, \hat{b}_k) & \dots & \text{Var}(\hat{b}_k) \end{bmatrix}$$

δηλαδή :

$$\text{Var}(\hat{b}) = \sigma^2 (X^T X)^{-1} \quad (4.5.1)$$

και αυτή είναι η μικρότερη μεταβλητότητα την οποία μπορεί να έχει οποιοδήποτε διάνυσμα γραμμικών και αμερόληπτων εκτιμητών του b , δηλαδή το διάνυσμα των εκτιμητών ελαχίστων τετραγώνων είναι και **άριστο**. Γι' αυτό όπως αναφέρθηκε ήδη σε προηγούμενο κεφάλαιο, οι εκτιμητές ελαχίστων τετραγώνων χαρακτηρίζονται συνοπτικά ως άριστοι γραμμικοί αμερόληπτοι ή BLUE (Best Linear Unbiased Estimator).

Το θεώρημα των Gauss - Markov γενικεύει τις παραπάνω ιδιότητες σε κάθε γραμμικό συνδυασμό των εκτιμητών ελαχίστων τετραγώνων. Έτσι, π.χ αν $C = 2b_1 + 4b_3 + b_4$ ένας γραμμικός συνδυασμός των παραμέτρων του μοντέλου τότε $C = 2b_1 + 4b_3 + b_4$ είναι άριστος, γραμμικός, αμερόληπτος ή BLUE εκτιμητής του C .

Οι εκτιμητές $\hat{b}_j, j=0, 1, \dots, k$ είναι γραμμικές συναρτήσεις των παρατηρήσεων y_1, y_2, \dots, y_n . Επομένως, αν ισχύει και η συνθήκη $\varepsilon \sim N_n(0, \sigma^2 I)$, δηλαδή αν

$$Y \sim N_n(Xb, \sigma^2 I)$$

τότε και η κατανομή του \hat{b} είναι η k -μεταβλητή κανονική δηλαδή ισχύει

$$\hat{b} \sim N_k(b, \sigma^2 (X^T X)^{-1})$$

Στην πράξη, η σ^2 είναι άγνωστη και την εκτιμούμε με την

$$s^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n - k - 1} = \frac{\hat{\varepsilon}^T \hat{\varepsilon}}{n - k - 1} \quad (4.5.2)$$

όπου $\hat{\varepsilon}^T = (\hat{\varepsilon}_1, \hat{\varepsilon}_2, \dots, \hat{\varepsilon}_n)$ το διάνυσμα των καταλοίπων.

Η (4.5.2) ορίζει ένα αμερόληπτο εκτιμητή της σ^2 όταν ισχύει η συνθήκη $E(\varepsilon \varepsilon^T) = \sigma^2 I$

Επειδή αντικαθιστούμε την σ^2 με την εκτίμηση της s^2 , οι έλεγχοι των υποθέσεων για τις τιμές των $b_j, j=0, \dots, k$ γίνονται με το κριτήριο t , όπως θα δούμε στη συνέχεια.

Οι εκτιμητές μεγίστης Πιθανοφάνειας

Εστω ότι ισχύουν οι συνθήκες $Cov(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j$ και $\epsilon \sim N_n(0, \sigma^2 I)$. Η συνάρτηση πυκνότητας πιθανότητας των $Y = (Y_1, Y_2, \dots, Y_n)$ δίνεται από την

$$L(Y|b, \sigma^2) = (2n\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} (Y-Xb)^T (Y-Xb)\right\}$$

Η L ως συνάρτηση των b και σ^2 ονομάζεται συνάρτηση πιθανοφάνειας και συμβολίζεται με $L(b, \sigma^2/Y)$. Ο νεπερέιος λογάριθμος της συνάρτησης πιθανοφάνειας ισούται με :

$$L = \ln L(b, \sigma^2/Y) = -\frac{n}{2} \ln 2n - \ln \sigma^2 - \frac{(Y-Xb)^T (Y-Xb)}{2\sigma^2} \quad (4.5.3)$$

Επειδή γενικά, ο $\ln z$ είναι αύξουσα συνάρτηση του z , οι τιμές των b και σ^2 που μεγιστοποιούν την L , είναι επομένως εκτιμήσεις μέγιστης πιθανοφάνειας. Είναι φανερό ότι η μεγιστοποίηση της L ισοδυναμεί με ελαχιστοποίηση της $(Y-Xb)^T (Y-Xb)$ και επομένως κάτω από την υπόθεση της κανονικότητας οι εκτιμητές ελαχίστων τετραγώνων \hat{b} είναι και εκτιμητές μέγιστης πιθανοφάνειας των παραμέτρων b του μοντέλου

Αποδεικνύεται ότι ο εκτιμητής μέγιστης πιθανοφάνειας της σ^2 ορίζεται από το τύπο

$$\frac{1}{n} (Y-X\hat{b})^T (Y-X\hat{b}) = \frac{1}{n} \hat{\epsilon}^T \hat{\epsilon} \quad (4.5.4)$$

είναι επομένως μεροληπτικός, αφού ο αμερόληπτος εκτιμητής της σ^2 είναι ο

$$S^2 = \frac{1}{n-k-1} \hat{\epsilon}^T \hat{\epsilon}$$

4.6 Ανάλυση διακυμάνσεως

Είδαμε, στο απλό μοντέλο, ότι ισχύει:

$$SSR = \sum (Y_i - \bar{Y})^2 = \hat{b}_1^2 \sum (X_i - \bar{X})^2 \quad (4.6.1)$$

και επειδη

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

η (4.6.1) γραφεται, ισοδυναμα, ως εξης:

$$\begin{aligned} SSR &= \hat{b}_1 \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= \hat{b}_1 (\sum x_i y_i - n \bar{x} \bar{y}) \\ &= \hat{b}_1 \sum x_i y_i - n \bar{y} (\bar{y} - \hat{b}_0) \\ &= \hat{b}_1 \sum x_i y_i + \hat{b}_0 \sum y_i - n \bar{y}^2 = \begin{bmatrix} \hat{b}_0 & \hat{b}_1 \end{bmatrix} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix} - n \bar{y}^2 \end{aligned}$$

Γενικά αν b είναι το $(k+1) \times 1$ διάνυσμα των εκτιμήσεων ελαχίστων τετραγώνων, Y το $n \times 1$ διάνυσμα των παρατηρήσεων της εξαρτημένης και X η $n \times (k+1)$ μήτρα σχεδιασμού, τότε ισχύει:

$$SSR = \hat{b}^T X^T Y - n \bar{y}^2 \quad (4.6.2)$$

Εξάλλου μπορεί εύκολα να διαπιστωθεί ότι ισχύουν:

$$\begin{aligned} Y^T Y &= \sum y_i^2 \\ \hat{E}^T \hat{E} &= \sum \hat{\varepsilon}_i^2 \\ \hat{Y}^T Y &= \sum \hat{y}_i^2 \end{aligned}$$

οπότε γράφουμε

$$SST = \sum y_i^2 - n \bar{y}^2 = Y^T Y - n \bar{y}^2$$

και

$$\begin{aligned} SSE &= SST - SSR \\ &= Y^T Y - \hat{b}^T X Y \\ &= \sum y_i^2 - \hat{b}_0 \sum y_i - \hat{b}_1 \sum y_i x_{1i} - \dots - \hat{b}_k \sum y_i x_{ki} \end{aligned}$$

Επομένως η γενική μορφή του πίνακα ανάλυσης διακυμάνσεως είναι η εξής :

SOURCE OF VARIATION	SUM OF SQUARES	D.F	MEAN SQUARE
REGRESSION	$SSR = b^T X^T Y - n\bar{y}^2$	K	SSR/K
ERROR	$SSE = Y^T Y - B^T X^T Y$	$n-k-1$	$SSE/n-k-1$
TOTAL	$SST = Y^T Y - n\bar{y}^2$	$n-1$	$SST/n-1$

4.7 Ο συντελεστής R^2

Σε ένα μοντέλο με $k > 1$ ερμηνευτικές μεταβλητές ο συντελεστής R^2 ονομάζεται συντελεστής πολλαπλού προσδιορισμού και ορίζεται όπως και στο απλό μοντέλο παλινδρόμησης δηλαδή ως εξής :

$$R^2 = \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST} \quad (4.7.1)$$

όπου

και ο τελεστής Σ είναι για όλες τις τιμές $i=1, \dots, n$. Δηλαδή ο συντελεστής R^2 είναι η ποσοστιαία μείωση του συνολικού τετραγωνικού σφάλματος εκτίμησης των παρατηρήσεων Y_i , η οποία οφείλεται στην εισαγωγή των μεταβλητών x_1, \dots, x_k στο μοντέλο.

Η θετική τετραγωνική ρίζα του R^2 ονομάζεται συντελεστής πολλαπλής συσχέτισης και ισούται με τον συντελεστή

συσχέτισης των παρατηρήσεων Y_i και των εκτιμήσεων $\hat{Y}_i = \hat{b}_0 + \hat{b}_1 X_{1i} + \dots + \hat{b}_k X_{ki}; i = 1, \dots, n$.
 Πράγματι έχουμε :

$$r_{Y\hat{Y}} = \frac{\sum (Y_i - \bar{Y})(\hat{Y}_i - \bar{Y})}{\sqrt{\sum (Y_i - \bar{Y})^2} \sqrt{\sum (\hat{Y}_i - \bar{Y})^2}} \quad (4.7.2)$$

και επειδή

$$\begin{aligned} \sum (Y_i - \bar{Y})(\hat{Y}_i - \bar{Y}) &= \sum [(\hat{Y}_i - \bar{Y}) + \hat{\epsilon}_i](\hat{Y}_i - \bar{Y}) \\ &= \sum (\hat{Y}_i - \bar{Y})^2 + \sum \hat{\epsilon}_i (\hat{Y}_i - \bar{Y}) \\ &= \sum (\hat{Y}_i - \bar{Y})^2 \end{aligned}$$

η (4.7.2) γράφεται ως εξής :

$$r_{Y\hat{Y}} = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sqrt{\sum (Y_i - \bar{Y})^2} \sqrt{\sum (Y_i - \bar{Y})^2}} = \frac{\sqrt{\sum (\hat{Y}_i - \bar{Y})^2}}{\sqrt{\sum (Y_i - \bar{Y})^2}} = \sqrt{R^2}$$

Σημείωση : Επειδή ο όρος SST και SSE είναι το άθροισμα των τετραγωνικών σφαλμάτων όταν για την εκτίμηση των παρατηρήσεων Y_i χρησιμοποιούμε αντίστοιχα, το βασικό μοντέλο και το μοντέλο με ερμηνευτικές μεταβλητές τις X_1, \dots, X_k μπορούμε να γράψουμε :

$$SST = SSE_0$$

$$SSE = SSE_k$$

Ο δείκτης k στο αντίστοιχο SSE δηλώνει το πλήθος των ερμηνευτικών μεταβλητών του μοντέλου. Έτσι ο συντελεστής προσδιορισμού γράφεται, ισοδύναμα, ως εξής :

$$R^2 = \frac{SSE_0 - SSE_k}{SSE_k} \quad (4.7.2)$$

Ο διορθωμένος συντελεστής προσδιορισμού

Στους δύο όρους που προσδιορίζουν τον R^2 ο $SST = \sum (Y_i - \bar{Y})^2$ είναι σταθερός ενώ ο όρος $SSE = \sum (Y_i - \hat{b}_0 - \hat{b}_1 X_{1i} - \dots - \hat{b}_k X_{ki})^2$ εξαρτάται από το πλήθος των ερμηνευτικών μεταβλητών k .

Στην πραγματικότητα με κάθε επί πλέον ερμηνευτική μεταβλητή στο μοντέλο, ο όρος SSE μειώνεται ή στην χειρότερη περίπτωση δεν μεταβάλλεται. Έτσι, μπορούμε να εκβιάσουμε υψηλή τιμή του R^2 προσθέτοντας απλώς ερμηνευτικές μεταβλητές στο μοντέλο. Αναφέρουμε ως ακραίο παράδειγμα ότι σε 2 ζεύγη τιμών προσαρμόζεται τέλεια ένα πολυώνυμο $n-1$ βαθμού επιτυγχάνεται δηλαδή $R^2 = 1$.

Είναι όμως προφανές ότι ένα τέτοιο μοντέλο δεν έχει καμία πραγματική χρησιμότητα.

Για να πάρουμε υπόψη και το πλήθος των παρατηρήσεων σε σχέση με τον αριθμό των ερμηνευτικών μεταβλητών του μοντέλου, υπολογίζουμε τον διορθωμένο συντελεστή προσδιορισμού R στον οποίο κάθε άθροισμα τετραγώνων διαιρείται με τους αντίστοιχους βαθμούς ελευθερίας. Δηλαδή έχουμε :

$$\bar{R}^2 = 1 - \frac{\sum (Y_i - \hat{Y}_i)^2 / (n - k - 1)}{\sum (Y_i - \bar{Y})^2 / (n - 1)} = 1 - \frac{S^2}{S^2_Y} \quad (4.7.3)$$

όπου S η διακύμανση των παρατηρήσεων \hat{Y}_i γύρω από το εκτιμηθέν μοντέλο και $S_{..}$ η διακύμανση των Y_i γύρω από τη μέση τιμή τους.

Ο διορθωμένος συντελεστής προσδιορισμού μας επιτρέπει να συγκρινούμε διαφορετικά μοντέλα που εκτιμούμε όχι μόνο για τις ίδιες παρατηρήσεις Y_i αλλά και για διαφορετικές. Θα πρέπει πάντως να σημειωθεί ότι

υπάρχουν επιφυλάξεις για την τελευταία χρήση, οι οποίοι θεωρούν ότι στην σύγκριση μοντέλων που εκτιμήθηκαν με διαφορετικά δεδομένα ο \bar{R}^2 μπορεί να χρησιμοποιηθεί μόνο ως ένας αρχικός, χονδρικός δείκτης και προτείνουν άλλα μέτρα της ερμηνευτικής ικανότητας του μοντέλου.

Είναι εύκολο να δειχθεί ότι οι συντελεστές \bar{R}^2 και R^2 συνδέονται.

Πράγματι, αντικαθιστώντας την (4.7.1) στην (4.7.3) παίρνουμε:

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n-1}{n-k-1} \quad (4.7.4)$$

Από τον τύπο αυτό προκύπτει ότι :

- Για $k \geq 1$ ισχύει $\bar{R}^2 \leq R^2$. Η ισότητα ισχύει μόνο όταν υπάρχει τέλεια προσαρμογή, δηλαδή $R^2=1$, ή ασυμπτωματικά, δηλαδή όταν το μέγεθος του δείγματος αυξάνει απεριόριστα.
- Ο \bar{R}^2 μπορεί να πάρει και αρνητικές τιμές (ενώ, θυμίζουμε, ότι ισχύει $0 \leq R^2 \leq 1$). Αυτό συμβαίνει όταν

$$R^2 < \frac{k}{n-1}$$

Αποδεικνύεται ότι αν σε ένα μοντέλο προσθέσουμε μία ερμηνευτική μεταβλητή έστω X_j , ενώ όλα τα υπόλοιπα παραμείνουν σταθερά τότε ο \bar{R}^2 θα αυξηθεί αν η τιμή του συντελεστή \hat{b}_j είναι μεγαλύτερη από ένα τυπικό σφάλμα S_{b_j} .

Ερμηνεία του R^2

Είδαμε ότι ο R^2 είναι η ποσοστιαία μείωση του ολικού σφάλματος εκτίμησης των παρατηρήσεων Y_i που οφείλεται στην εισαγωγή των μεταβλητών X_1, \dots, X_k στο μοντέλο. Επειδή υπολογίζεται εύκολα (με τη βοήθεια του H/Y φυσικά)

συνοδεύει σχεδόν πάντα το εκτιμηθέν μοντέλο μαζί με τα τυπικά σφάλματα των συντελεστών μερικής παλινδρόμησης.

Η ερμηνεία των ακραίων τιμών του R^2 είναι προφανής. Συγκεκριμένα όταν το εκτιμηθέν μοντέλο προσαρμόζεται τέλεια στα δεδομένα, δηλαδή όταν $Y_i = \hat{Y}_i$ $i=1, \dots, n$ τότε $SSE=0$ και $R=1$ (η περίπτωση αυτή είναι, φυσικά, εντελώς απίθανη, στην πράξη). Όταν η εξαρτημένη δεν συνδέεται με τις ερμηνευτικές, δηλαδή ισχύει $b_1 = \dots = b_k = 0$ και στο μοντέλο περιλαμβάνεται μόνο ο σταθερός όρος τότε, καμιά μείωση του ολικού σφάλματος δεν επιτυγχάνεται και ισχύει

$SSE = SST$ οπότε $R^2 = 0$. Οι ενδιάμεσες τιμές ερμηνεύονται αναλόγως. Η ευκολία με την οποία ερμηνεύεται η τιμή του R^2 οδηγεί συχνά σε μια μηχανιστική αξιολόγηση του μοντέλου, με βάση τον συντελεστή προδιορισμού. Μπορούμε όμως να αποφύγουμε πολλές από τις παγίδες που στήνει μια τέτοια πρακτική, έχοντας υποψιν τα εξής :

- Ας υποθέσουμε ότι οι X_1, \dots, X_k δεν συνδέονται με την Y , ισχύει δηλαδή $b_1 = \dots = b_k = 0$ εμείς όμως αυτό δεν το γνωρίζουμε και τις έχουμε περιλάβει στο μοντέλο. Τότε, λόγω της τυχαιότητας των παρατηρήσεων, η τιμή του R μπορεί να μην είναι μηδέν. Αποδεικνύεται ότι, κάτω από ασθενείς υποθέσεις (παρατηρήσεις Y ασυσχέτιστες, με κατανομή συμμετρική) ισχύει:

Έτσι, π.χ για $n = 15$ και $k = 7$ η αναμενόμενη τιμή του R ισούται με 0,5 από καθαρή τύχη ! Γι'αυτό χρήσιμος είναι ο ακόλουθος πρακτικός κανόνας : ερμηνεύουμε μόνο το μέρος της τιμής του R που είναι μεγαλύτερο από την αναμενόμενη τιμή του .

- Ο συντελεστής R^2 είναι μέτρο καλής προσαρμογής του εκτιμηθέντος μοντέλου στα δεδομένα . Δεν μας δίνει όμως καμιά πληροφορία για την προβλεπτική του ικανότητα ιδίως όταν η πρόβλεψη γίνεται για τιμές των ερμηνευτικών μεταβλητών έξω από τα όρια των δεδομένων μας.

- Ο συντελεστής R^2 μπορεί, θεωρητικά, να πάρει τη μέγιστη τιμή του 1 όταν στα δεδομένα δεν υπάρχουν επαναλαμβανόμενες τιμές της ανεξάρτητης τιμής. Όταν έχουμε περισσότερες από μία παρατηρήσεις στο ίδιο σημείο σχεδιασμού (X, \dots, X) τότε η μεταβλητότητα των Y δεν οφείλεται στην μεταβλητότητα των X , είναι δηλαδή καθαρό σφάλμα οπότε το R δεν μπορεί να πάρει την τιμή 1.

- Υψηλές τιμές του R^2 μπορούν να προκύψουν σε περιπτώσεις που παραβιάζεται μια ή περισσότερες από τις συνθήκες Gauss - Markov. Έτσι π.χ σχεδόν σε κάθε μοντέλο παλινδρόμησης το οποίο εκτιμάται με δεδομένα χρονικών σειρών, παρατηρείται εξαιρετικά υψηλός συντελεστής προσδιορισμού ενώ συγχρόνως δεν απορρίπτεται η υπόθεση για ύπαρξη αυτοσυσχέτισης των διαταρακτικών όρων. Το ίδιο μπορεί να συμβαίνει λόγω της ύπαρξης μιας ή περισσότερων ακραίων τιμών στα δεδομένα μας. Μπορούμε να αποφύγουμε λανθασμένα συμπεράσματα, όταν ερμηνεύουμε τον R^2 , αν συγχρόνως αναλύουμε τα κατάλοιπα του μοντέλου.

4.8 Οι συντελεστές μερικού προσδιορισμού

Ο συντελεστής πολλαπλού προσδιορισμού R^2 μετρά την ικανότητα όλων μαζί των ερμηνευτικών μεταβλητών X_1, \dots, X_k να ερμηνεύσουν τη μεταβλητότητα της εξαρτημένης Y . Δεν μας λέει όμως πόση από την ικανότητα αυτή μπορεί να αποδοθεί σε κάθε μια από τις, X_1, \dots, X_k χωριστά. Για το σκοπό αυτό υπολογίζουμε τους συντελεστές μερικού προσδιορισμού.

Ο συντελεστής μερικού προσδιορισμού της Y επί μιας μεταβλητής έστω της X_k ή Τετραγωνικός Συντελεστής Πολλαπλής Συσχέτισης στις Y και X_k και συμβολίζεται με $R^2_{Y|X_k}$ και υπολογίζεται ως εξής :

α. Εκτιμούμε την παλινδρόμηση της Y επί των X_1, \dots, X_{k-1} και έστω

$\hat{Y}_i, i = 1, \dots, n$ τα κατάλοιπα.

β. Εκτιμούμε την παλινδρόμηση της X_k επί των X_1, \dots, X_{k-1} και έστω $\hat{X}_i, i = 1, \dots, n$ τα κατάλοιπα.

γ. Εκτιμούμε την παλινδρόμηση της Y επί της X_k . Ο συντελεστής προσδιορισμού αυτής της εκτίμησης είναι ο $R^2_{Y_k | 1, 2, \dots, k-1}$

Από τον τρόπο υπολογισμού του προκύπτει ότι ο $R^2_{Y_k | 1, 2, \dots, k-1}$ μετρά το ποσοστό της μεταβλητότητας της Y η οποία μένει ανερμήνευτη ή κατάλοιπη από τις X_1, \dots, X_{k-1} και ερμηνεύεται από την κατάλοιπη μεταβλητότητα της X_k .

Ο $R^2_{Y_k | 1, 2, \dots, k-1}$ υπολογίζεται και διαφορετικά ως εξής:

- Εκτιμούμε την παλινδρόμηση της Y επί των X_1, \dots, X_{k-1} και έστω SSE_{k-1} το άθροισμα των τετραγωνικών σφαλμάτων.

- Εκτιμούμε την παλινδρόμηση της Y επί των X_1, \dots, X_{k-1}, X_k και έστω SSE_k το άθροισμα των τετραγωνικών σφαλμάτων. Η διαφορά $SSE_{k-1} -$

$SSE_k = SSR(X_k | X_1, \dots, X_{k-1})$ είναι η μείωση του ολικού σφάλματος εκτίμησης λόγω της επί πλέον εισαγωγής της X_k στο μοντέλο στο οποίο ήδη υπήρχαν οι X_1, \dots, X_{k-1}

- Ο συντελεστής μερικού προσδιορισμού της Y επί της X_k είναι η ποσοστιαία μείωση του ολικού σφάλματος εκτίμησης λόγω της επιπλέον εισαγωγής της X_k στο μοντέλο. Δηλαδή

$$R^2_{Y_k | 1, 2, \dots, k-1} = \frac{SSE_{k-1} - SSE_k}{SSE_{k-1}} = \frac{SSR(X_k | X_1, \dots, X_{k-1})}{SSE_{k-1}} \quad (4.8.1)$$

Παράδειγμα

Στο ακόλουθο μοντέλο Y είναι η ζήτηση χρήματος σε δισ. δρχ., X το ΑΕΠ σε δισ. δρχ. και X το επιτόκιο καταθέσεων ταμειυτηρίου. Η εκτίμηση έγινε με βάση τις 17 ετήσιες παρατηρήσεις της Ελληνικής Οικονομίας για τα έτη 1970-1986 που δημοσιεύτηκαν στο Δελτίο της Τράπεζας της Ελλάδας, το 1987

$$Y = -0,3716 + 0,353X + 5,345X$$

Δίνονται ακόμη

$$SSE = 0,2984$$

$$SST = 1,1146$$

απ' όπου υπολογίζουμε : $SSR = SST - SSE = 1,1146 - 0,2984 = 0,8162$ και

$$R^2 = \frac{0,8162}{1,1146} = 0,7323$$

$$R^2 = 1 - \frac{0,2984/14}{1,1146/16} = 0,694$$

Εκτιμήθηκε ακόμη το απλό μοντέλο παλινδρόμησης με ερμηνευτική μεταβλητή την X_1 ως εξής :

$$Y = -1,033 + 3,54X_1$$

$$\text{με } SSE = 0,392$$

οπότε υπολογίζουμε :

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{0,392}{1,1146} = 0,6487$$

Χρησιμοποιώντας τον συμβολισμό που υιοθετήθηκε προηγουμένως σημειώνουμε :

$$SSE_1 = 0,392$$

$$SSE_2 = 0,2984$$

συντελεστής μερικού προσδιορισμού της Y επί της X_2 υπολογίζεται ως εξής :

$$R^2_{Y|X_2} = \frac{SSE_1 - SSE_2}{SSE_1} = \frac{0,392 - 0,2984}{0,392} = 0,239$$

Το απλό μοντέλο της Y επί της X_2 εκτιμήθηκε ως εξής :

$$Y = -0,2876 + 5,85X_2$$

$$\text{με } SSE = 0,3121$$

Ο μερικός συντελεστής προσδιορισμού της Y επί της X_1 υπολογίζεται ως εξής :

$$R^2_{Y|X_1} = \frac{0,3121 - 0,2984}{0,3121} = 0,044$$

Ο συντελεστής μερικής συσχέτισης

Ο συντελεστής μερικής συσχέτισης των μεταβλητών Y και X_k με σταθερές τις X_1, X_2, \dots, X_{k-1} συμβολίζεται με $r_{Y|X_1, \dots, X_{k-1}}$ και ισούται με τη θετική τετραγωνική ρίζα του αντίστοιχου συντελεστή μερικού προσδιορισμού, δηλαδή :

$$r_{Y|X_1, \dots, X_{k-1}} = \sqrt{R^2_{Y|X_1, \dots, X_{k-1}}}$$

Από τον τρόπο που υπολογίζεται προκύπτει ότι ο μερικός συντελεστής συσχέτισης $r_{Y|X_1, \dots, X_{k-1}}$ μετρά την ένταση της γραμμικής συμμεταβολής ανάμεσα στην εξαρτημένη και την ερμηνευτική μεταβλητή X_k μετά την αφαίρεση των γραμμικών επιδράσεων των ερμηνευτικών X_1, \dots, X_{k-1} και στην Y και στην X_k .

4.9 Έλεγχος των στατιστικών υποθέσεων

Όταν εκτιμούμε ένα μοντέλο παλινδρόμησης θέλουμε να ελέγξουμε αν μια ή περισσότερες από τις ερμηνευτικές μεταβλητές είναι πλεονάζουσες, δηλαδή μπορούν να αφαιρεθούν χωρίς σημαντική

μείωση της ερμηνευτικής του ικανότητας . Υπάρχουν δύο περιπτώσεις που μία ερμηνευτική μεταβλητή, έστω η X_k , πλεονάζει :

α. Όταν η X_k δεν συνδέεται με την Y

β. Όταν η ερμηνευτική ικανότητα της X_j περιέχεται στις υπόλοιπες , οπότε η επί πλέον ερμηνευτική της ικανότητα (επομένως και ο συντελεστής $R^2_{y/k/1,\dots,k-1}$) ισούται με μηδέν ή σχεδόν . Η περίπτωση αυτή θα γίνει ευκολότερα κατανοητή με το ακόλουθο παράδειγμα .

Παράδειγμα

Έστω η συνάρτηση παλινδρόμησης

$$E (Y_i) = 2X_1 + X_2 + 7X_3$$

και ισχύει $X_3 = X_1 + X_2$. Τότε η $E (Y_i)$ μπορεί να γράφει, ισοδύναμα, με έναν από τους ακόλουθους τρόπους :

$$E (Y_i) = 9X_1 + 8X_2$$

$$E (Y_i) = X_1 + 8X_3$$

$$E (Y_i) = 9X_3 - X_2$$

Στο μοντέλο αυτό , οποιοδήποτε ερμηνευτικών μεταβλητών ζεύγος ερμηνευτικών μεταβλητών κάνει την Τρίτη περιττή παρόλο που αυτή συνδέεται με την εξαρτημένη . Οι μεταβλητές που συνδέονται **συγγραμμικές** . Δύο ή περισσότερες συγγραμμικές ερμηνευτικές μεταβλητές κάνουν αδύνατη την εκτίμηση των παραμέτρων με τη μέθοδο των ελαχίστων τετραγώνων . Στην πράξη , συνηθέστερη είναι η περίπτωση που δύο ή περισσότερες ερμηνευτικές μεταβλητές είναι **σχεδόν συγγραμμικές** δηλαδή συνδέονται με μία σχέση σχεδόν γραμμικής με την έννοια ότι ο αντίστοιχος συντελεστής πολλαπλής συσχέτισης είναι κοντά στο 1 . Τα αποτελέσματα της ατελούς συγγραμμικότητας δεν είναι προφανή γι' αυτό και η ανίχνευση της δεν είναι εύκολη . Στη συνέχεια θα παρουσιάσουμε τη διαδικασία ελέγχου της μηδενικής υπόθεσης ότι μια ή περισσότερες από τις ερμηνευτικές μεταβλητές του μοντέλου πλεονάζουν είτε λόγω συγγραμμικότητας είτε διότι δεν έχουν σχέση με την εξαρτημένη .

Η μηδενική και η εναλλακτική υπόθεση του ελέγχου εξειδικεύονται ως εξής :

$$H_0 = b_{p+1} = \dots = b_k = 0$$

$$H_\varepsilon = \text{μια τουλάχιστον από τις } b_j = 0, j = p+1, \dots, k$$

Στην H_ε αντιστοιχεί το μοντέλο

$$Y = b_0 + b_1 X_1 + \dots + b_p X_p + b_{p+1} X_{p+1} + \dots + b_k X_k + \varepsilon$$

το οποίο ονομάζεται **πλήρες ή εναλλακτικό** ενώ στην H_0 αντιστοιχεί το μοντέλο

$$Y = b_0 + b_1 X_1 + \dots + b_p X_p + \varepsilon$$

που ονομάζεται **αναγμένο ή μηδενικό ή μοντέλο κάτω από τον γραμμικό περιορισμό** $b_{p+1} = \dots = b_k = 0$. Ο έλεγχος της μηδενικής υπόθεσης ισοδυναμεί με έλεγχο της υπόθεσης ότι οι μεταβλητές X_{p+1}, \dots, X_k πλεονάζουν, δηλαδή η αφαίρεση τους απ το μοντέλο δεν μειώνει σημαντικά την ερμηνευτική του ικανότητα.

Το κριτήριο απόφασης ορίζεται με βάση τα αθροίσματα των τετραγωνικών σφαλμάτων των δύο μοντέλων ως εξής :

Έστω SSE_k και SSE_p το άθροισμα των τετραγωνικών σφαλμάτων του πλήρους και του μειωμένου μοντέλου αντίστοιχα. Επειδή $p < k$ και κάθε επιπλέον μεταβλητή δεν μπορεί να μειώνει το SSE ισχύει $SSE_p \geq SSE_k$. Η διαφορά των δύο SSE είναι η αύξηση της ερμηνευτικής ικανότητας του μοντέλου με την επιπλέον εισαγωγή των μεταβλητών X_{p+1}, \dots, X_k και συμβολίζουν με $SSR (X_{p+1}, \dots, X_k / X_1, \dots, X_p)$ Δηλαδή έχουμε :

$$SSR (X_{p+1}, \dots, X_k / X_1, \dots, X_p) = SSE_p - SSE_k$$

Αποδεικνύεται ότι, όταν οι διαταρακτικοί όροι ακολουθούν την κανονική κατανομή έχουν μέση τιμή μηδέν και είναι ανά δύο ασυσχέτιστοι (που λόγω της κανονικότητας των ε_i , συνεπάγεται ότι είναι και ανεξάρτητες) τότε το στατιστικό

$$F = \frac{SSE_p - SSE_k / (k - p)}{SSE_k / (n - k - 1)} \quad (4.9.1)$$

ακολουθεί την κατανομή F_{v_1, v_2} με $v_1 = k - p$ και $v_2 = n - k - 1$ βαθμούς ελευθερίας. Έτσι, $F_{v_1, v_2, \alpha}$ είναι η τιμή της F_{v_1, v_2} για την οποία ισχύει :

$$P (F_{v_1, v_2} > F_{v_1, v_2, \alpha}) = \alpha$$

και F_π η τιμή του στατιστικού F που υπολογίζεται με βάση την (4.9.1) εφαρμόζουμε το ακόλουθο κριτήριο απόφασης :

Αν $F_\pi > F_{v_1, v_2, \alpha}$ απορρίπτουμε την H_0 στο $\alpha = 0,05$

4.10 Ο ολικός έλεγχος σημαντικότητας του μοντέλου

Στον έλεγχο αυτό συγκρίνουμε το πλήρες μοντέλο των k ερμηνευτικών μεταβλητών με το μειωμένο μοντέλο στο οποίο υπάρχει μόνο ο σταθερός όρος. Ειδικότερα έχουμε :

Η μηδενική και η εναλλακτική υπόθεση

$$H_0: b_1 = \dots = b_k = 0$$

$$H_e : \text{ένα τουλάχιστον } b_j \neq 0, j=1, \dots, k$$

Το μειωμένο είναι το μοντέλο με μονό το σταθερό όρο ή βασικό μοντέλο με αντίστοιχο ολικό άθροισμα τετραγώνων το $SSE_0 = \sum (Y_i - \bar{Y})^2$ το οποίο συμβολίσαμε ως SST.

Κριτήριο αποφάσεως

$$A_n \quad F_n = \frac{(SST - SSE_k)/k}{SSE_k / (n - k - 1)} = \frac{SSR(x_1, x_2, \dots, x_k / x_0)/k}{SSE_k / (n - k - 1)} > F_{v_1, v_2, \alpha}$$

απορρίπτουμε την H_0 στο επίπεδο σημαντικότητας της α . Σημειώνεται ότι $v_1 = k$ και $v_2 = n - k - 1$.

Όταν απορρίπτουμε την H_0 στον έλεγχο αυτό λέμε ότι εκτιμήσαμε μια στατιστικά σημαντική παλινδρόμηση. Αυτό σημαίνει ότι οι ερμηνευτικές μεταβλητές X_1, \dots, X_k βελτιώνουν σημαντικά τις εκτιμήσεις των παρατηρήσεων Y_i , σε σχέση με το βασικό μοντέλο. Η σημαντικότητα εδώ θα πρέπει να εννοηθεί ως στατιστική σημαντικότητα. Δηλαδή, η αύξηση της ερμηνευτικής ικανότητας του μοντέλου είναι μεγαλύτερη απ' εκείνη που θα αποδίδουμε στην τύχη στο $100(1-\alpha)\%$ των ανάλογων δειγμάτων δεδομένων με το ίδιο n και την ίδια μήτρα X .

Η τιμή του στατιστικού F , σ' αυτό τον έλεγχο, μπορεί να υπολογιστεί απ' ευθείας από το συντελεστή προσδιορισμού R^2 . Σημειώνεται ότι στον ορισμό του συντελεστή R^2 επειδή δεν υπήρχε ανάγκη να διακρίνουμε μεταξύ διαφόρων μοντέλων δεν χρησιμοποιήσαμε δείκτες στα αντίστοιχα αθροίσματα τετραγώνων. Για να κρατήσουμε αντίστοιχα με το συμβολισμό που χρησιμοποιήσαμε στο μέρος αυτό συμβολίζουμε με SSE_k το άθροισμα των τετραγωνικών

σφαλμάτων. Του μοντέλου με $k \geq 1$ ερμηνευτικές μεταβλητές ενώ τη διαφορά $SST - SSE_k$ συμβολίζουμε ως $SSR_k = SSR(x_1, \dots, x_k/x_0)$.

Έτσι έχουμε:

$$\left. \begin{aligned} R^2 &= \frac{SSR_k}{SST} \\ 1 - R^2 &= \frac{SSE_k}{SST} \end{aligned} \right\} \Rightarrow \frac{R^2}{1 - R^2} = \frac{SSR_k}{SSE_k} \Rightarrow \frac{R^2/k}{1 - R^2/n - k - 1} =$$

$$= \frac{SSR_k/k}{SSE_k/n - k - 1} = \frac{(SST - SSE_k)/k}{SSE_k/n - k - 1} = F_n \quad (4.10.1)$$

Πάντως, σπανίως θα χρειαστεί να χρησιμοποιήσουμε την (4.10.1) αφού όλα τα προγράμματα H/Y τυπώνουν, στον πίνακα ανάλυσης διακυμάνσεως, την F_π ως F-value. Μαζί δίνουν την p-value η τιμή της πιθανότητας του ελέγχου που είναι η πιθανότητα, όταν ισχύει η H_0 , να πάρουμε μία τιμή για την F όση η F_π ή μεγαλύτερη. Αν η p-value είναι μικρότερη από το α (συνήθως ίσο με 0,05) απορρίπτουμε την H_0 .

Προβλήματα του ελέγχου

Ο έλεγχος αυτός, όπως και όλοι οι ολικοί έλεγχοι, είναι ελάχιστα χρήσιμος αφού μας δίνει λίγες πληροφορίες για το αν εξειδικεύσαμε σωστά το μοντέλο. Έτσι, όταν απορρίπτουμε την H_0 απλώς δεχόμαστε ότι ένας τουλάχιστον από τους μερικούς συντελεστές παλινδρόμησης είναι στατιστικά σημαντικός αλλά δεν ξέρουμε ποιος. Είναι μάλιστα πιθανό, στον έλεγχο για κάθε b_j χωριστά, κανένας να μην είναι στατιστικά διαφορετικός από μηδέν. Απ'αυτό μπορούμε απλώς να συμπεράνουμε ότι υπάρχει κάποια (ατελής) συγγραμμικότητα μεταξύ των ερμηνευτικών μεταβλητών που κάνει την οριακή συμβολή της κάθε μιας στατιστικά ασήμαντη. Κρατώντας τις αναλογίες, είναι σαν να χρησιμοποιήσουμε π.χ. 4 άτομα ίσης σωματικής ισχύος για να μεταφέρουμε ορισμένο βάρος το οποίο μπορούν να μεταφέρουν 2 άτομα.

Θα πρέπει εδώ να σημειωθεί ακόμη ότι η απόρριψη της H_0 δεν είναι αρκετή για να πάρουμε από το μοντέλο ικανοποιητικές εκτιμήσεις των Y_i . Έχει δειχθεί ότι για να πάρουμε ικανοποιητικές εκτιμήσεις θα πρέπει να ισχύει $F_\pi > 4F_{v_1, v_2, \alpha}$.

ΚΕΦΑΛΑΙΟ ΠΕΜΠΤΟ

Ψευδομεταβλητές

5.1 Εισαγωγή

Σε πολλές πρακτικές εφαρμογές μια ή περισσότερες ποιοτικές μεταβλητές ασκούν σημαντική επίδραση στην εξαρτημένη μεταβλητή του μοντέλου. Έτσι π.χ. οι πωλήσεις ενός προϊόντος μπορεί να επηρεάζονται σοβαρά από την εποχή του έτους, οι εξαγωγές από την ένταξη της χώρας στην ΕΟΚ, η ζήτηση ενός προϊόντος από το αν η παρατήρηση γίνεται σε αγροτική ή μη αγροτική περιοχή.

Μία ποιοτική μεταβλητή με δύο δυνατές τιμές ονομάζεται **δίτιμη** ή **διχοτομική** διαφορετικά, ονομάζεται **πλειότιμη** ή **παράγοντας** και οι τιμές της **επίπεδα**. Έτσι, η ύπαρξη ή απουσία ενός χαρακτηριστικού όταν γίνεται η παρατήρηση Y_i είναι μια διχοτομική μεταβλητή. Η απάντηση σε μία ερώτηση όταν μπορεί να είναι «συμφωνώ», «διαφωνώ», «δεν απαντώ» είναι μία πλειότιμη μεταβλητή. Η ποιότητα του προϊόντος όταν μπορεί να καταταγεί σε μια από τις κατηγορίες «κακή», «μέτρια», «καλή», «εξαιρετική» είναι πλειότιμη μεταβλητή ενώ όταν μπορεί να καταταγεί σε «καλό» ή «ελαττωματικό» είναι δίτιμη μεταβλητή.

Σε ένα μοντέλο για την τυχαία μεταβλητή Y μπορούμε να πάρουμε υπόψη μια δίτιμη ποιοτική μεταβλητή εισάγοντας ως ερμηνευτική μια **ψευδομεταβλητή ή μεταβλητή δείκτη** που ορίζεται ως εξής:

$$D = \begin{cases} 1 & \text{αν η παρατήρηση γίνεται στο } \alpha \\ 0 & \text{διαφορετικά} \end{cases}$$

Ο τρόπος που η D θα εισαχθεί στο μοντέλο εξαρτάται από το μοντέλο αλλά και το είδος των μεταβολών που θέλουμε να συλλάβουμε.

5.2 Μεταβολή του σταθερού όρου στο μοντέλο παλινδρόμησης

Όταν θεωρούμε ότι η μετάβαση από τη μία τιμή της δίτιμης μεταβλητής στην άλλη έχει ως αποτέλεσμα την παράλληλη μετατόπιση

του μοντέλου κατά μία σταθερά, εισάγουμε ως επί πλέον ερμηνευτική μεταβλητή μια ψευδομεταβλητή.

Παράδειγμα

Θέλουμε να εκτιμήσουμε ένα μοντέλο παλινδρόμησης για τη μηνιαία ζήτηση Y ενός αναψυκτικού με ερμηνευτική μεταβλητή τη μηνιαία δαπάνη για διαφήμιση x . Για να πάρουμε υπόψη την επίδραση του καλοκαιριού στην κατανάλωση του αναψυκτικού εισάγουμε στο μοντέλο την ψευδομεταβλητή D η οποία παίρνει την τιμή 1 για τους καλοκαιρινούς μήνες και μηδέν διαφορετικά. Τα δεδομένα είναι μηνιαίες παρατηρήσεις μιας χρονιάς:

Μήνας	Y _i (εκ. δρχ)	X _i (σε χιλ. δρχ)	D _i
I	1,145	312	0
Φ	1,912	344	0
M	1,739	307	0
A	2,092	435	0
M	2,617	307	0
I	9,375	376	1
I	14,012	443	1
A	14,812	548	1
Σ	11,001	568	0
O	3,710	312	0
N	2,561	496	0
Δ	1,912	447	0

Εκτιμούμε το μοντέλο

$$\hat{Y} = -5,363 + 0,022x + 8,16D$$

$(-1,8)$ $(2,9)$ $(5,1)$

$$\mu \epsilon \quad R^2 = 0,842 \quad , \quad SSE = 46,6 \quad S = 2,13$$

και οι αριθμοί στις παρενθέσεις είναι οι t-αναλογίες (b_j / s_{b_j}) για τις αντίστοιχες εκτιμήσεις. Έτσι το μοντέλο της ζήτησης για έναν καλοκαιρινό μήνα ($D=1$) γίνεται

$$\hat{Y} = (-5,363 + 8,16) + 0,022x$$

ή

$$\hat{Y} = 2,797 + 0,022X \quad (5.2.1)$$

ενώ για ένα μη καλοκαιρινό μήνα το μοντέλο ζήτησης είναι το εξής:

$$\hat{Y} = -5,363 + 0,022X \quad (5.2.2)$$

Αγνοώντας την επίδραση του καλοκαιριού στη ζήτηση, εκτιμούμε το μοντέλο

$$\hat{Y} = -8,07 + 0,033 X \quad (5.2.3)$$

$(-1,4) \quad (2,5)$

$$\mu\epsilon \quad R^2 = 0,38 \quad s = 4,3$$

το οποίο έχει μεγαλύτερη κλίση από τα (5.2.1) και (5.2.2).

Συμπαιρένουμε ότι αν δεν ληφθεί υπόψη η επίδραση του καλοκαιριού τότε έχουμε μια υπερεκτίμηση της συμβολής της διαφήμισης στις μεταβολές της ζήτησης. Η διαφορετικά, αυτό έχει ως αποτέλεσμα θετικό σφάλμα μεροληψίας στην εκτίμηση του b_1 . Πάντως το μοντέλο (5.2.3) έχει κακή προσαρμογή στα δεδομένα όπως φαίνεται από το χαμηλό συντελεστή προσδιορισμού (αλλά και την χαμηλή τιμή του στατιστικού Durbin-Watson=0,82, που είναι ένδειξη κακής προσαρμογής όπως θα δούμε στο επόμενο κεφάλαιο).

5.3 Μεταβολή της κλίσης ή των μερικών κλίσεων του μοντέλου

Αν σε ένα μοντέλο παλινδρόμησης εισάγουμε ως επί πλέον ερμηνευτική μεταβλητή το γινόμενο μιας ψευδομεταβλητής D επί μια ερμηνευτική μεταβλητή, έστω x_j , τότε μπορούμε να συλλάβουμε μεταβολές στο συντελεστή (μερικής) παλινδρόμησης της x_j . Η DX_j ονομάζεται πολλαπλαστική ψευδομεταβλητή.

Παράδειγμα

Σε μια έρευνα των παραγόντων που ερμηνεύουν το επίπεδο ανάπτυξης μιας χώρας εκτιμήθηκε το ακόλουθο μοντέλο παλινδρόμησης :

όπου

Y = ο ρυθμός μεγέθυνσης του ΑΕΠ μιας χώρας,

X_1 = ο ρυθμός μεγέθυνσης των εξαγωγών

X_2 = ο ρυθμός μεγέθυνσης του εργατικού δυναμικού

$X_3 = 0$ ο ρυθμός μεγέθυνσης του κεφαλαίου της χώρας

$$D = \begin{cases} 1 & \text{αν το επίπεδο ανάπτυξης της χώρας είναι πάνω από το} \\ & \text{θεωρούμενο κριτικό} \\ 0 & \text{διαφορετικά} \end{cases}$$

Αν η κλίση της πολλαπλασιαστικής ψευδομεταβλητής X_1D είναι στατιστικά σημαντική τότε μπορούμε να συμπεράνουμε ότι :
πάνω από το υποτιθέμενο 'κριτικό' επίπεδο ανάπτυξης, η οριακή συμβολή των εξαγωγών στη μεγέθυνση του ΑΕΠ είναι μεγαλύτερη απ' ότι κάτω απ' αυτό.

5.4 Πλειότιμες ποιοτικές μεταβλητές

Για να πάρουμε υπόψη μια ποιοτική μεταβλητή με δύο τιμές εισάγουμε στο μοντέλο μια ψευδομεταβλητή. Γενικά, για να πάρουμε υπόψη μια ποιοτική μεταβλητή με $\lambda > 2$ τιμές εισάγουμε στο μοντέλο $\lambda - 1$ ψευδομεταβλητές.

Έτσι, π.χ έστω ότι θέλουμε να εκτιμήσουμε το απλό μοντέλο κατανάλωσης

$$Y_i = b_0 + b_1 X_i + \varepsilon_i \quad i = 1, \dots, n$$

όπου $Y =$ η κατανάλωση και το $X =$ το διαθέσιμο εισόδημα μιας οικογένειας. Έστω ακόμη ότι διαθέτουμε διαστρωματικά δεδομένα από 3 διαφορετικές περιοχές, έστω Α, Β, Γ. Αν υποπτευόμαστε ότι το μοντέλο μεταβάλλεται στις τρεις περιοχές μπορούμε να χρησιμοποιήσουμε δύο ψευδομεταβλητές, τις D_1 και D_2 που ορίζονται ως εξής :

$$D_1 = \begin{cases} 1 & \text{αν η παρατήρηση γίνεται στην περιοχή Α} \\ 0 & \text{διαφορετικά} \end{cases}$$

$$D_2 = \begin{cases} 1 & \text{αν η παρατήρηση γίνεται στην περιοχή Β} \\ 0 & \text{διαφορετικά} \end{cases}$$

Διακρίνουμε τις ακόλουθες περιπτώσεις ανάλογα με τις διαφορές του μοντέλου στις τρεις περιοχές.

I. Διαφορετικός σταθερός όρος. Το μοντέλο παλινδρόμησης εξειδικεύεται ως εξής :

$$Y_i = b_0 + b_1 X_i + \gamma_1 D_{1i} + \gamma_2 D_{2i} + \varepsilon_i \quad i = 1, \dots, n$$

το οποίο είναι ισοδύναμο με τρία μοντέλα, ένα για κάθε περιοχή ως εξής:

$$Y_i = (b_0 + \gamma_1) + b_1 X_i + \varepsilon_i \quad \text{για παρατηρήσεις } Y_i \text{ στην περιοχή Α}$$

$$Y_i = (b_0 + \gamma_2) + b_1 X_i + \varepsilon_i$$

για παρατηρήσεις Y_i στην περιοχή Β

$$Y_i = b_0 + b_1 X_i + \varepsilon_i$$

στην περιοχή Γ.

για παρατηρήσεις Y_i στην περιοχή Β

για παρατηρήσεις Y_i στην περιοχή Γ.

II. Διαφορετική κλίση ή στο παράδειγμα μας, διαφορετική οριακή ροπή για κατανάλωση. Τότε θα χρησιμοποιήσουμε δύο πολλαπλασιαστές μεταβλητές οπότε το μοντέλο εξειδικεύεται ως εξής :

$$Y_i = b_0 + b_1 X_i + \delta_1 D_{1i} + \delta_2 D_{2i} + \varepsilon_i \quad i = 1, \dots, n$$

το οποίο γίνεται για τις 3 περιοχές, αντίστοιχα, ως εξής :

$$Y_i = b_0 + (b_1 + \delta_1) X_i + \varepsilon_i \quad \text{για την Α}$$

$$Y_i = b_0 + (b_1 + \delta_2) X_i + \varepsilon_i \quad \text{για την Β}$$

$$Y_i = b_0 + b_1 X_i + \varepsilon_i \quad \text{για την Γ}$$

III. Διαφορετικά και ο σταθερός όρος και η κλίση. Το μοντέλο του παραδείγματος μας παίρνει τη μορφή

$$Y_i = b_0 + b_1 X_i + \gamma_1 D_{1i} + \gamma_2 D_{2i} + \delta_1 D_{1i} X_i + \delta_2 X_i + \varepsilon_i$$

που γίνεται, για τις τρεις περιοχές, αντίστοιχα :

$$Y_i = (b_0 + \gamma_1) + (b_1 + \delta_1) X_i + \varepsilon_i \quad \text{για την Α}$$

$$Y_i = (b_0 + \gamma_2) + (b_1 + \delta_2) X_i + \varepsilon_i \quad \text{για την Β}$$

$$Y_i = b_0 + b_1 X_i + \varepsilon_i \quad \text{για την Γ}$$

Παράδειγμα

Για την ημερήσια ζήτηση του γιαουρτιού μάρκας Φ το οποίο βγαίνει σε 4 τύπους, εκτιμήθηκε το ακόλουθο μοντέλο παλινδρόμησης :

Υ

όπου Υ είναι η ημερήσια ζήτηση του γιαουρτιού Φ, σε χιλιάδες κιβώτια, X είναι ο τηλεοπτικός χρόνος σε sec που διαφημίστηκε η Φ την προηγούμενη μέρα και οι ψευδομεταβλητές D_1 , D_2 , D_3 ορίζονται στον ακόλουθο πίνακα :

Τύπος γιαουρτιού	D_1	D_2	D_3
A	1	0	0
B	0	1	0
Γ	0	0	1
Δ	0	0	0

Έτσι, για τον Δ τύπο γιαουρτιού ($D_1 = D_2 = D_3 = 0$) το μοντέλο είναι το

$$\hat{Y} = 14,716 + 0,673X \quad (5.4.1)$$

Για τον A τύπο το μοντέλο είναι το :

$$\begin{aligned} \hat{Y} &= (14,716 + 0,387) + 0,673X \\ &= 15,103 + 0,673X \quad (5.4.2) \end{aligned}$$

Έτσι, αν η Φ δεν διαφημιστεί από την τηλεόραση την προηγούμενη μέρα εκτιμούμε ότι η ζήτηση για τον A τύπο θα είναι μεγαλύτερη από τη ζήτηση του Δ κατά 0,387 χιλιάδες κιβώτια.

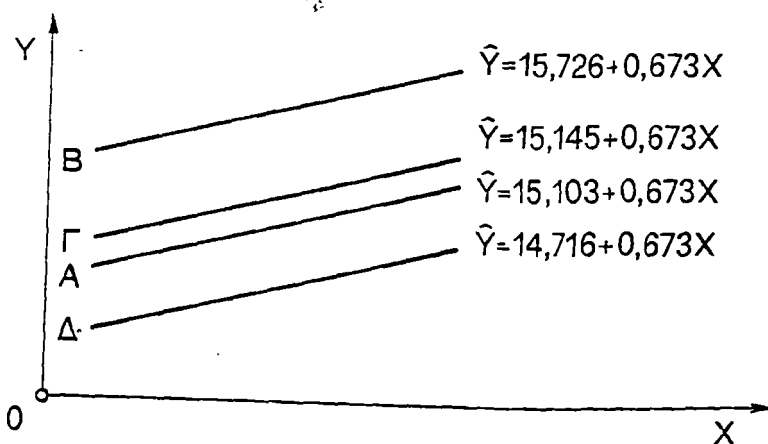
Για τους τύπους Β και Γ το μοντέλο παλινδρόμησης γίνεται αντίστοιχα:

$$\begin{aligned}\hat{Y} &= (14,716 + 1,013) + 0,673 X \\ &= 15,729 + 0,673 X\end{aligned}\quad (5.4.3)$$

και

$$\begin{aligned}\hat{Y} &= (14,716 + 0,429) + 0,673 X \\ &= 15,145 + 0,673 X\end{aligned}\quad (5.4.4)$$

Ο όρος $b_1=0,673$ είναι η προβλεπόμενη αύξηση της ζήτησης για κάθε τύπο γιαουρτιού όταν ο διαφημιστικός χρόνος της προηγούμενης μέρας αυξηθεί κατά 1 sec. Η γραφική παράσταση των μοντέλων (5.4.1)-(5.4.4) δίνεται στο ακόλουθο σχήμα.



5.5 Έλεγχος της σταθερότητας του μοντέλου με ψευδομεταβλητές.

Όταν εκτιμούμε ένα μοντέλο παλινδρόμησης υποθέτουμε ότι οι παράμετροι b_0, b_1, \dots, b_k είναι σταθερές σε όλες τις παρατηρήσεις Y_i . Σε πολλές περιπτώσεις η γραφική παράσταση των καταλοίπων ή πληροφορίες, για τον τρόπο που έχουν συλλεγεί τα δεδομένα μας υπαγορεύουν να εκτιμήσουμε δύο μοντέλα- ένα για τις n_1 παρατηρήσεις και ένα για τις $n_2 = n - 1$

δηλαδή τα εξής:

$$Y_i = b_0' + b_1' X_{1i} + \dots + b_k' X_{ki} + \varepsilon_i, \quad i = 1, \dots, n$$

$$Y_i = b_0^2 + b_1^2 X_{1i} + \dots + b_k^2 X_{ki} + \varepsilon_i, \quad i = n_1 + 1, \dots, n \quad (5.5.1)$$

Με βάση το δείγμα των δεδομένων μας μπορούμε να ελέγξουμε την υπόθεση ότι οι παράμετροι των δύο μοντέλων είναι ίδιες που σημαίνει ότι το μοντέλο είναι κοινό για όλες τις παρατηρήσεις $Y_i, i = 1, \dots, n$. Η μηδενική και εναλλακτική υπόθεση του ελέγχου εξειδικεύονται ως εξής :

$$H_0 : b_j^1 = b_j^2 \quad j = 0, 1, \dots, k$$

$$H_e : b_j^1 \neq b_j^2 \quad j = 0, 1, \dots, k$$

κάτω από την H_0 έχουμε το μοντέλο

$$Y_i = b_0 + b_1 X_{1i} + \dots + b_k X_{ki} + \varepsilon_i, \quad i = 1, \dots, n \quad (5.5.2)$$

Ενώ κάτω από την H_e έχουμε τα δύο μοντέλα (5.5.1) τα οποία μπορούν να γράφουν ενιαία ως εξής :

$$Y_i = b_0 + b_1 X_{1i} + \dots + b_k X_{ki} + \gamma_0 D + \gamma_1 D X_{1i} + \dots + \gamma_k D X_{ki} + \varepsilon_i$$

$$i = 1, \dots, n \quad (5.5.3)$$

όπου

$$D_i = \begin{cases} 1 & \text{αν } i = 1, \dots, n_1 \\ 0 & \text{αν } i = n_1 + 1, \dots, n \end{cases}$$

και

$$\gamma_0 = b_0^1 - b_0^2, \quad \gamma_1 = b_1^1 - b_1^2, \dots, \gamma_k = b_k^1 - b_k^2$$

Για να ελέγξουμε την H_0 έναντι της H_e δηλαδή το (5.5.3) και έστω SSE_e το άθροισμα των τετραγωνικών καταλοίπων το οποίο έχει $n - (k+1)$ βαθμούς ελευθερίας. Η διαφορά $(SSE_0 - SSE_e)$ έχει $[n - (k+1)] - [n - 2(k+1)] = k+1$ βαθμούς ελευθερίας.

Αποδεικνύεται ότι το στατικό

$$F = \frac{(SSE_0 - SSE_e) / (k+1)}{SSE_e / (n-2)(k+1)} \quad (5.5.4)$$

κολουθεί την κατανομή F_{v_1, v_2} με $v_1 = k+1$ και $v_2 = n-2(k+1)$ βαθμούς ελευθερίας.

Επομένως αν η τιμή που υπολογίζουμε από την (5.5.4) είναι μεγαλύτερη από την $F_{v_1, v_2, \alpha}$ απορρίπτουμε την H_0 στο $\alpha = 0,05$.

Ο έλεγχος αυτός είναι γνωστός ως έλεγχος του Chow και μπορεί να εφαρμοστεί μόνο αν $n_1 > k+1$ και $n_2 > k+1$.

Αν $n_1 > k+1$ αλλά $n_2 < k+1$ ο Chow (1960) προτείνει την ακόλουθη παραλαγή του :

- Εκτιμούμε το μοντέλο με τις n_1 παρατηρήσεις και έστω SSE_1 το άθροισμα των τετραγωνικών αποκλίσεων το οποίο έχει $n_1 - (k+1)$ βαθμούς ελευθερίας. Η διαφορά $(SSE_0 - SSE_1)$ έχει n_2 βαθμούς ελευθερίας. Το στατιστικό :

$$F = \frac{(SSE_0 - SSE_1) / n_2}{SSE_1 / (n_1 - k - 1)}$$

συγκρίνεται με την τιμή $F_{v_1, v_2, \alpha}$ της κατανομής F με $v_1 = n_2$ και $v_2 = n_1 - k - 1$ βαθμούς ελευθερίας.

ΚΕΦΑΛΑΙΟ ΕΚΤΟ

Η αυτοσυσχέτιση των σφαλμάτων

6.1 Εισαγωγή

Ο σκοπός της εκτίμησης ενός μοντέλου παλινδρόμησης είναι να συλλάβει όλη τη συστηματικότητα στη συμπεριφορά της εξαρτημένης έτσι ώστε να μείνουν, ως υπόλοιπη μεταβλητότητα, τυχαία και ανεξάρτητα σφάλματα δηλαδή ποσότητες μη προβλέψιμες. Όταν διαδοχικά σφάλματα εκτίμησης συσχετίζονται λέμε ότι υπάρχει σειριακή συσχέτιση ή αυτοσυσχέτιση στα κατάλοιπα. Ειδικότερα, αν θετικά σφάλματα τείνουν να ακολουθούνται από θετικά, η αυτοσυσχέτιση είναι θετική, διαφορετικά, είναι αρνητική. Είναι προφανές ότι στην περίπτωση αυτή μπορούμε να βελτιώσουμε τις εκτιμήσεις του μοντέλου με νέα εξειδίκευση η οποία θα αρεί αυτή τη συστηματικότητα. Η αυτοσυσχέτιση στα σφάλματα εκτίμησης είναι, γενικά, ένδειξη κακής εξειδίκευσης του μοντέλου. Όταν εκτιμούμε ένα μοντέλο με δεδομένα χρονικών σειρών, η αυτοσυσχέτιση των σφαλμάτων μπορεί, σχεδόν πάντα, να αποδοθεί στο ότι λανθασμένα έγινε η υπόθεση των ασυσχετιστών διαταρακτικών ορών. Σε άλλες περιπτώσεις μπορεί να οφείλεται στο ότι δεν έγινε σωστή εξειδίκευση του μαθηματικού τύπου της συνάρτησης παλινδρόμησης ή στο ότι έχει παραλειφθεί μια σημαντική ερμηνευτική μεταβλητή της οποίας οι τιμές αυτοσυσχετίζονται. Τέλος, η αυτοσυσχέτιση των σφαλμάτων της εκτίμησης μπορεί να οφείλεται στο ότι ένα ή περισσότερα ακραία σημεία στα δεδομένα προκαλούν μετατόπιση του εκτιμηθέντος μοντέλου.

6.2 Η αυτοσυσχέτιση των διαταρακτικών όρων.

Μια από τις υποθέσεις του μοντέλου παλινδρόμησης είναι ότι οι διαταρακτικοί όροι είναι ανά δύο ασυσχέτιστοι δηλαδή ότι ισχύει :

$$Cov(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$$

Με άλλα λόγια υποθέτουμε ότι, για δεδομένο σημείο σχεδιασμού, μια παρατήρηση Y_i θα διαφέρει από τη μέση τιμή της κατά ένα ποσό το οποίο είναι ασυσχέτιστο με το μέγεθος της αντίστοιχης διαφοράς σε οποιοδήποτε άλλο σημείο.

Ο διαταρακτικός όρος είναι το αποτέλεσμα της συνδυασμένης δράσης πολλών παραγόντων οι οποίοι είτε δεν είναι αρκετά σημαντικοί για να περιληφθούν στο μοντέλο είτε δεν είναι μετρήσιμοι.

Όταν αναλύσουμε δεδομένα χρονικών σειρών, δηλαδή παρατηρήσεις διατεταγμένες ως προς το χρόνο τότε, συνήθως, η συμπεριφορά των παραγόντων αυτών είναι όμοια με τη συμπεριφορά τους σε προηγούμενες περιόδους. Αυτό έχει ως αποτέλεσμα οι διαταρακτικοί όροι που αντιστοιχούν σε διαφορετικές παρατηρήσεις, συνήθως γειτονικές, να συσχετίζονται. Στην περίπτωση αυτή λέμε ότι οι διαταρακτικοί όροι συσχετίζονται σειριακά ή αυτοσυσχετίζονται.

6.3 Το μοντέλο αυτοσυσχέτισης AR (1)

Η πιο απλή μορφή αυτοσυσχέτισης είναι αυτή στην οποία διαδοχικοί διαταρακτικοί όροι συνδέονται με το ακόλουθο σχήμα:

$$\varepsilon_i = \rho \varepsilon_{i-1} + u_i \quad (6.3.1)$$

όπου u_i ασυσχέτιστες τυχαίες μεταβλητές με μηδενική μέση τιμή και σταθερή διακύμανση σ_u^2 και ανεξάρτητες των $\varepsilon_{i-1}, \varepsilon_{i-2}, \varepsilon_{i-3}, \dots$. Το μοντέλο αυτοσυσχέτισης που ορίζεται από την (6.3.1) ονομάζεται **αυτοπαλινδρομούμενο πρώτης τάξης** και συμβολίζεται ως AR(1). Ονομάζεται έτσι διότι ορίζει ένα μοντέλο παλινδρόμησης του ε_i επί του ε_{i-1} .

Θυμίζουμε ότι στο απλό μοντέλο παλινδρόμησης της Y επί της X ισχύει:

$$b_1 = \rho_{xy} \frac{\sigma_y}{\sigma_x}$$

επομένως και το μοντέλο (6.3.1) ισχύει :

$$b_1 = \rho_1 \frac{\sigma_1^2}{\sigma_1^2 - 1}$$

όπου ρ_1 ο συντελεστής αυτοσυσχέτισης πρώτης τάξης διαδοχικών διαταρακτικών όρων. Αυτά τα u_i είναι ομοσκεδαστικά τότε $\sigma_{\varepsilon_i}^2 = \sigma_{\varepsilon_{i-1}}^2$ και

$$b_1 = \rho_1$$

Δηλαδή ο συντελεστής ρ του μοντέλου (6.3.1) είναι ο συντελεστής αυτοσυσχέτισης πρώτης τάξης των διαταρακτικών όρων.

Είναι εύκολο να δειχθεί ότι ένα σφάλμα στην παρατήρηση i επιστρέφει σε όλες τις επόμενες περιόδους και το μέγεθος της επίδρασης αυτής μειώνεται με το χρόνο. Αρκεί να υπολογίσουμε την συνδιακύμανση του ε_i με το διαταρακτικό όρο οποιαδήποτε προηγούμενης περιόδου :

$$\text{Cov}(\varepsilon_i, \varepsilon_{i-1}) = E(\varepsilon_i \varepsilon_{i-1}) = E(\rho \varepsilon_{i-1} + u_i) \varepsilon_{i-1} = E(\rho \varepsilon_{i-1}^2 + u_i \varepsilon_{i-1})$$

$$\Rightarrow \text{Cov}(\varepsilon_i, \varepsilon_{i-1}) = E(\rho \varepsilon_{i-1} + u_i) \varepsilon_{i-1}$$

$$= E(\rho \varepsilon_{i-1}^2 + u_i \varepsilon_{i-1})$$

$$= \rho E(\varepsilon_{i-1}^2) = \rho \text{Var}(\varepsilon_{i-1}) = \rho \sigma_\varepsilon^2$$

Ομοίως υπολογίζουμε

$$\text{Cov}(\varepsilon_i, \varepsilon_{i-2}) = \rho^2 \sigma_\varepsilon^2$$

$$\text{Cov}(\varepsilon_i, \varepsilon_{i-3}) = \rho^3 \sigma_\varepsilon^2$$

$$\vdots$$

$$\text{Cov}(\varepsilon_i, \varepsilon_{i-\lambda}) = \rho^\lambda \sigma_\varepsilon^2$$

Επομένως για το συντελεστή αυτοσυσχέτισης λ -τάξης, $\lambda=1, 2, 3, \dots$ ισχύει:

$$\rho_\lambda = \frac{\text{Cov}(e_i, e_{i-\lambda})}{\sigma_\varepsilon^2} = \frac{\rho^\lambda \sigma_\varepsilon^2}{\sigma_\varepsilon^2} = \rho^\lambda \quad \lambda = 1, 2, 3, \dots$$

Συμπεραίνουμε ότι στο AR(1) μοντέλο ολόκληρη η δομή της συσχέτισης των διαταρακτικών όρων προσδιορίζεται από μία παράμετρο- τον συντελεστή ρ .

Συνέπειες των AR(1) σφαλμάτων στην εκτίμηση ελάχιστων τετραγώνων

Η μέθοδος των ελαστικών τετραγώνων σε ένα μοντέλο με AR(1) σφάλματα ορίζει εκτιμητές οι οποίοι εξακολουθούν να είναι αμερόληπτοι. Δεν είναι όμως αμερόληπτες οι διακυμάνσεις τους και συνεπώς και τα τυπικά τους σφάλματα. Ειδικότερα αν $\rho > 0$ τότε ισχύει:

$$E(s_{\hat{\beta}_j}^2) < \sigma_{\hat{\beta}_j}^2 \quad j = 0, 1, \dots, k$$

δηλαδή υπάρχει μια συστηματική τάση να υποεκτιμώνται τα τυπικά σφάλματα και συνεπώς να θεωρούνται στατιστικά σημαντικές οι εκτιμήσεις των παραμέτρων ενώ δεν είναι. Αντίστοιχα, όταν $\rho < 0$ τότε τα τυπικά σφάλματα τείνουν να υπερκτιμούνται.

Αν στις ερμηνευτικές μεταβλητές του μοντέλου περιλαμβάνεται και η εξαρτημένη σε υστέρηση τότε τα προβλήματα αυτοσυσχέτισης είναι πιο σοβαρά.

6.4 Το στατιστικό Durbin - Watson.

Έστω το μοντέλο $Y_i = b_0 + b_1 X_{1i} + \dots + b_k X_{ki} + \varepsilon_i \quad i = 1, \dots, n$ του οποίου τις παραμέτρους εκτιμήσαμε με την μέθοδο των ελάχιστων τετραγώνων και $\hat{\varepsilon}_i = Y_i - \hat{Y}_i, i = 1, \dots, n$ τα κατάλοιπα αυτής της εκτίμησης. Το στατιστικό Durbin - Watson υπολογίζεται ως εξής :

$$(6.4.1) \quad D - W \equiv d = \frac{\sum_{i=2}^n (\hat{\varepsilon}_i - \hat{\varepsilon}_{i-1})^2}{\sum_{i=1}^n \hat{\varepsilon}_i^2}$$

Ο συντελεστής αυτοσυσχέτισης πρώτης τάξης των καταλόγων υπολογίζεται, σύμφωνα με τον τύπο (3.12.1) ως εξής :

$$(6.4.2) \quad r_1 = \frac{\sum_{i=2}^n \hat{\varepsilon}_i \hat{\varepsilon}_{i-1}}{\sum_{i=1}^n \hat{\varepsilon}_i^2}$$

Αναπτύσσοντας τον αριθμητή της (6.4.1) έχουμε :

$$d = \frac{\sum \hat{\varepsilon}_i^2}{\sum \varepsilon_i^2} + \frac{\sum \hat{\varepsilon}_{i-1}^2}{\sum \hat{\varepsilon}_i^2} - 2 \frac{\sum \hat{\varepsilon}_i \hat{\varepsilon}_{i-1}}{\sum \hat{\varepsilon}_i^2}$$

Όταν το μέγεθος του δείγματος είναι μεγάλο τότε

$$\sum_{i=2}^n \hat{\varepsilon}_i^2 \cong \sum_{i=1}^{n-1} \hat{\varepsilon}_{i-1}^2 \cong \sum_{i=1}^n \hat{\varepsilon}_i^2$$

οπότε χρησιμοποιώντας και

την (6.4.2) το στατιστικό Durbin - Watson γίνεται

$$d = 2 (1 - \bar{r}_1)$$

Είναι προφανές οι ακόλουθες σχέσεις :

$$r_1 = 0 \iff d = 2$$

$$r_1 = -1 \iff d = 4$$

$$r_1 = 1 \iff d = 0$$

Ο έλεγχος Durbin - Watson για AR(1) σφάλματα

Όταν $n \rightarrow \infty$ οι διαταρακτικοί όροι ή σφάλματα ενός μοντέλου είναι AR(1) είναι συντελεστής αυτοσυσχέτισης πρώτης τάξης ρ είναι η μοναδική παράμετρος που προσδιορίζει τη δομή αυτόσυσχέτισης των διαταρακτικών όρων. Στα μεγάλα δείγματα ο συντελεστής Γ_1 αποτελεί μια ικανοποιητική εκτίμηση του ρ .

Επομένως μια τιμή Γ_1 κοντά στο 1 ή στο -1 και συνεπώς μια τιμή για το $D - W$ κοντά στο 0 ή στο 4 αποτελούν σοβαρή ένδειξη AR(1) διαταρακτικών όρων. Η κατανομή δειγματοληψίας

του d εξαρτάται από τις τιμές των ερμηνευτικών μεταβλητών του μοντέλου. Οι Durbin και Watson (1950) υπολόγισαν, για επιλεγμένα επίπεδα σημαντικότητας δύο κριτικές τιμές d_L και d_U που εξαρτώνται μόνο από το πλήθος k των ερμηνευτικών μεταβλητών και όχι από τις τιμές τους. Έτσι, η μηδενική υπόθεση των ασυσχέτιστων διαταρακτικών όρων ελέγχεται ως προς την εναλλακτική υπόθεση ότι είναι AR(1), με $\rho > 0$ ως εξής :

1. Αν $d < d_L$ η H_0 απορρίπτεται
2. Αν $d > d_U$ η H_0 δεν απορρίπτεται
3. Αν $d_L < d < d_U$ δεν μπορούμε να αποφανθούμε

Για $\rho > 0$ το στατιστικό συγκρίνουμε το στατιστικό $4-d$ με τις κριτικές τιμές d_L και d_U σαν να κάναμε έλεγχο για θετική αυτοσυσχέτιση.

Η περιοχή αβεβαιότητας είναι σοβαρό μειονέκτημα του ελέγχου αυτού. Αποδεικνύεται πάντως ότι, όσο αυξάνει το μέγεθος του δείγματος το εύρος της περιοχής αβεβαιότητας τείνει να μειώνεται. Για τα μοντέλα τα οποία εκτιμώνται με δεδομένα οικονομικών χρονικών σειρών η τακτική που συνήθως ακολουθούνται είναι να θεωρείται και η περιοχή αβεβαιότητας ως περιοχή απόρριψης της H_0 .

6.5 Εκτίμηση ενός μοντέλου με AR(1) διαταρακτικούς τύπους

Για να εφαρμόσουμε τη μέθοδο των Ελαχίστων Τετραγώνων στην εκτίμηση του γραμμικού μοντέλου παλινδρόμησης όταν οι διαταρακτικοί όροι συνδέονται με μια AR(1) διαδικασία πρέπει να εφαρμόσουμε προηγουμένως έναν μετασχηματισμό των δεδομένων, έτσι ώστε να αφαιρεθεί η αυτοσυσχέτιση των διαταρακτικών όρων. Έστω λοιπόν το μοντέλο

$$(6.5.1) \quad Y_i = b_0 + b_1 X_{1i} + \dots + b_k X_{ki} + \varepsilon_i$$

με

$$(6.5.2) \quad \varepsilon_i = \rho \varepsilon_{i-1} + u_i \quad |\rho| < 1 \quad i=2, \dots, n$$

και $u_i (0 - \sigma_u^2)$, $E(u_i u_j) = 0$ για κάθε $j \neq i$ και

$E(u_i \varepsilon_{i-1}) = E(u_i \varepsilon_{i-2}) = \dots = 0$. Το μοντέλο αυτό ισχύει και

για την παρατήρηση Y_{i-1} δηλαδή έχουμε :

$$(6.5.3) \quad Y_{i-1} = b_0 + b_1 X_{1i-1} + \dots + b_k X_{ki-1} + \varepsilon_{i-1}$$

Πολλαπλασιάζουμε με ρ και παίρνουμε :

$$(6.5.4) \quad \rho Y_{i-1} = \rho b_0 + b_1 \rho X_{1i-1} + \dots + b_k \rho X_{ki-1} + \rho \varepsilon_{i-1}$$

Αφαιρούμε την (6.5.4) από την χρησιμοποιούμε την

(6.5.2) οπότε προκύπτει το μοντέλο :

$$(6.5.5) \quad Y_i - \rho Y_{i-1} = b_0(1-\rho) + b_1(X_{1i} - \rho X_{1i-1}) + \dots + b_k(X_{ki} - \rho X_{ki-1}) + u_i$$

στο οποίο οι διαταρακτικοί λογοί είναι ασυσχέτιστοι. Επομένως για να εκτιμήσουμε τις παραμέτρους του αρχικού μοντέλου αρκεί να εκτιμήσουμε τις παραμέτρους του (6.5.5). Αυτό υπαγορεύει τους ακόλουθους μετασχηματισμούς των δεδομένων :

$$(6.5.6) \quad \begin{aligned} Y_i^* &= Y_i - \rho Y_{i-1} \\ X_{1i}^* &= X_{1i} - \rho X_{1i-1} \\ &\vdots \\ X_{ki}^* &= X_{ki} - \rho X_{ki-1} \quad i=2, \dots, n \end{aligned}$$

που ονομάζεται **γενικευμένες Διαφορές**. Για $\rho=1$ οι (6.5.6) γίνονται :

$$Y_i^* = Y_i - Y_{i-1}$$

$$X_{1i}^* = X_{1i} - X_{1i-1}$$

$$\vdots$$

$$X_{ki}^* = X_{ki} - X_{ki-1} \quad i = 2, \dots, n$$

και ονομάζεται **πρώτες Διαφορές** των Y και X_1, \dots, X_k αντίστοιχα. Για $|\rho| < 1$ οι μετασχηματισμοί (6.5.6) ορίζουν **οιονεί Διαφορές**.

Με τους μετασχηματισμούς (6.5.6) χάνεται η πρώτη παρατήρηση γεγονός που μπορεί να είναι πολύ σημαντικό για μικρά δείγματα.

Η λύση είναι να πάρουμε :

$$Y_1^* = \sqrt{1-\rho^2} Y_1 \quad \text{και} \quad X_{j1}^* = \sqrt{1-\rho^2} X_{j1} \quad j = 1, \dots, k$$

Αυτός ο μετασχηματισμός έχει ως αποτέλεσμα η διακύμανση του διαταρακτικού όρου $\varepsilon_1^* = \sqrt{1-\rho^2} \varepsilon_1$ να ισούται με την διακύμανση όλων των άλλων διαταρακτικών όρων. Πράγματι έχουμε :

$$\text{Var}(\varepsilon_1^*) = (1-\rho^2) \text{Var}(\varepsilon_i) = \text{Var}(\varepsilon_i) - \rho^2 \text{Var}(\varepsilon_i) = \sigma_u^2$$

Επομένως, το πρόβλημα της εκτίμησης των παραμέτρων ενός μοντέλου με AR(1) διαταρακτικούς όρους καταλήγει στο πρόβλημα της εκτίμησης του ρ . Όλες οι μέθοδοι εκτίμησης του ρ που χρησιμοποιούνται στην πράξη είναι επαναληπτικές και επομένως δεν μπορούν να εφαρμοστούν παρά μόνο με H/Y. Εδώ θα παρουσιάσουμε τις πιο βασικές.

Μέθοδος Cochrane - Orcutt

Αυτή είναι η πιο παλιά μέθοδος εκτίμησης και περιλαμβάνει τα ακόλουθα στάδια :

- Εκτιμούμε τις παραμέτρους του μοντέλου με τη μέθοδο των Ελάχιστων Τετραγώνων (OLS).

- Υπολογίζουμε τα κατάλοιπα $\varepsilon_i = Y_i - \hat{Y}_i$, $i=1, \dots, n$ και στη συνέχεια τον συντελεστή αυτοσυσχέτισης πρώτης τάξης απ' το τύπο

$$\Gamma_1 = \frac{\sum_{i=2}^n \hat{\varepsilon}_i \hat{\varepsilon}_{i-1}}{\sum_{i=1}^n \varepsilon_i^2}$$

ή απ' ευθείας, από το στατιστικό Durbin - Watson d ως εξής :

$$\Gamma_1 = 1 - \frac{d}{2}$$

- Παίρνουμε τις οιονεί διαφορές των δεδομένων με $\rho = \Gamma_1$ δηλαδή ως εξής :

$$Y_i^* = Y_i - \hat{\rho} Y_{i-1}, X_{1i}^* = X_{1i} - \hat{\rho} X_{1i-1}, \dots, X_{ki}^* = X_{ki} - \hat{\rho} X_{ki-1} \quad i = 1, \dots, n$$

- Εκτιμούμε τις παραμέτρους του μοντέλου χρησιμοποιώντας τα μετασχηματισμένα δεδομένα $(Y_i^*, X_{1i}^*, \dots, X_{ki}^*)$, $i=1, \dots, n$

Η διαδικασία επαναλαμβάνεται μέχρι ότου η τιμή του ρ από δύο διαδοχικές επαναλήψεις $\Delta \rho$ διαφέρει περισσότερο από μια τιμή πολύ μικρή συνήθως ίση με 0,0001.

Να σημειωθεί ότι η διαδικασία αυτή ισοδυναμεί με ελαχιστοποίηση του SSE του μοντέλου (6.5.5) που είναι μη γραμμική συνάρτηση των παραμέτρων $\rho, b_0, b_1, \dots, b_k$ και μπορεί να έχει περισσότερα από ένα ελάχιστα. Έτσι, είναι δυνατόν η μέθοδος Cochrane - Orcutt να οδηγήσει ελάχιστα σε τοπικό και όχι σε ολικό ελάχιστο.

Η μέθοδος Hildred - Lu

Για όλες τις δυνατές τιμές του ρ στο διάστημα $-1 \leq \rho \leq 1$, με μήκος βήματος, συνήθως, ίσο με 0,01 υπολογίζουμε τις οιονεί διαφορές Y με τις οποίες, στη συνέχεια εκτιμούμε το μοντέλο παλινδρόμησης και υπολογίζουμε το SSE.

Τέλος, επιλέγεται η τιμή του ρ για την οποία το SSE είναι ελάχιστο.

Το μειονέκτημα αυτής της μεθόδου είναι το μήκος βήματος που συνήθως επιλέγεται είναι μικρό, διαφορετικά έχει τεράστιο υπολογιστικό κόστος. Αυτό έχει ως αποτέλεσμα ότι μπορεί να προσπερνάται η τιμή του ρ η οποία ελαχιστοποιεί το SSE.

6.6 Αυτοσυσχέτιση των καταλοίπων λόγω κακής εξειδίκευσης της συνάρτησης παλινδρόμησης.

Θα επισημάνουμε ορισμένες χαρακτηριστικές περιπτώσεις κακής εξειδίκευσης που έχει ως αποτέλεσμα την αυτοσυσχέτιση των καταλοίπων και συνακόλουθα μια τιμή για το στατιστικό Durbin - Watson κοντά στο μηδέν ή στο 4.

- Αν υποθέσουμε ότι θέλουμε να ερμηνεύσουμε την Y της οποίας Οι διαδοχικές παρατηρήσεις συσχετίζονται με την X η οποία είναι ασυσχέτιστη με την Y . Είναι προφανές ότι τα κατάλοιπα αυτής της εκτίμησης αυτοσυσχετίζονται και η πληροφορία αυτή ενσωματώνεται στην πολύ μικρή ή πολύ μεγάλη τιμή του στατιστικού Durbin-Watson όπως στο ακόλουθο.

Παράδειγμα: Το μήκος της ανθρώπινης ζωής

Πολλοί άνθρωποι πιστεύουν ότι το μήκος Y της ζωής ενός ανθρώπου που μπορεί να προβλεφθεί από το μήκος X της λεγόμενης "γραμμή της ζωής" στην παλάμη του χεριού του. Στα ακόλουθα Ζεύγη τιμών X_i είναι το μήκος της γραμμής της ζωής σε εκατοστόμετρα, στρογγυλεμένο στο πλησιέστερο 0,15cm Y_i είναι η ηλικία θανάτου του ατόμου στρογγυλεμένη στο πλησιέστερο έτος. Τα δεδομένα δίνονται από τους Wilson M.E-Mather L.E : 'Life expectancy' Journal of the American Medical Association, Vol 229, Nr 11, 1974. Δίνονται επίσης από τους Draper-Smith, 1981, σελ. 67-68.

Από το διάγραμμα διασποράς των (X_i, Y_i) που δίνεται στο σχήμα προκύπτει ότι δεν υπάρχει συστηματική συμμεταβολή των X_i και Y_i .

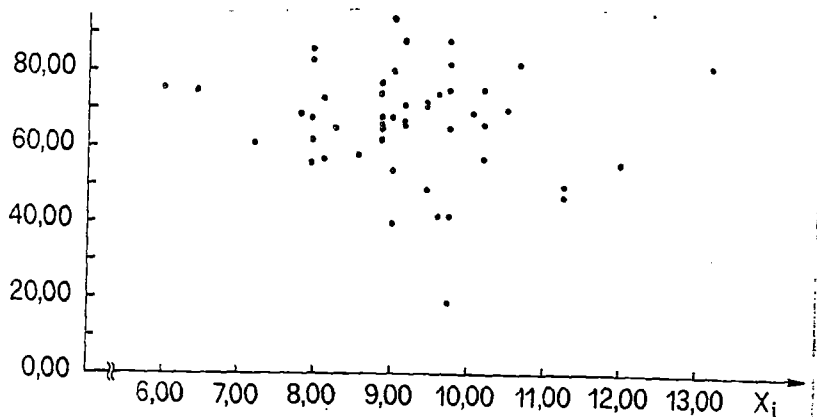
Τα αποτελέσματα της εκτίμησης του απλού γραμμικού μοντέλου με σταθερό όρο είναι τα ακόλουθα (οι αριθμοί στις παρενθέσεις είναι οι t-αναλογίες):

$$Y = 79,23 - 1,37X$$

$$(5,34) \quad (0,86)$$

$$\text{με } R = 0,015 \quad D-W = 0,00796$$

$$SSE = 9608 \quad S = 14,115$$



Παρατηρούμε ότι η εκτίμηση του συντελεστή b_1 δεν είναι στατιστικά σημαντική αφού

$$t_{\pi} = b_1 / S \quad b_1 = -0,86 < 2$$

Έτσι, η μηδενική υπόθεση $b_1 = 0$ δεν μπορεί να απορριφθεί για τα συνηθισμένα επίπεδα σημαντικότητας.

Δοκιμάζουμε ακόμη το μοντέλο χωρίς τον σταθερό όρο και παίρνουμε τα ακόλουθα αποτελέσματα:

$$Y = 7,09 X$$

$$(26,32)$$

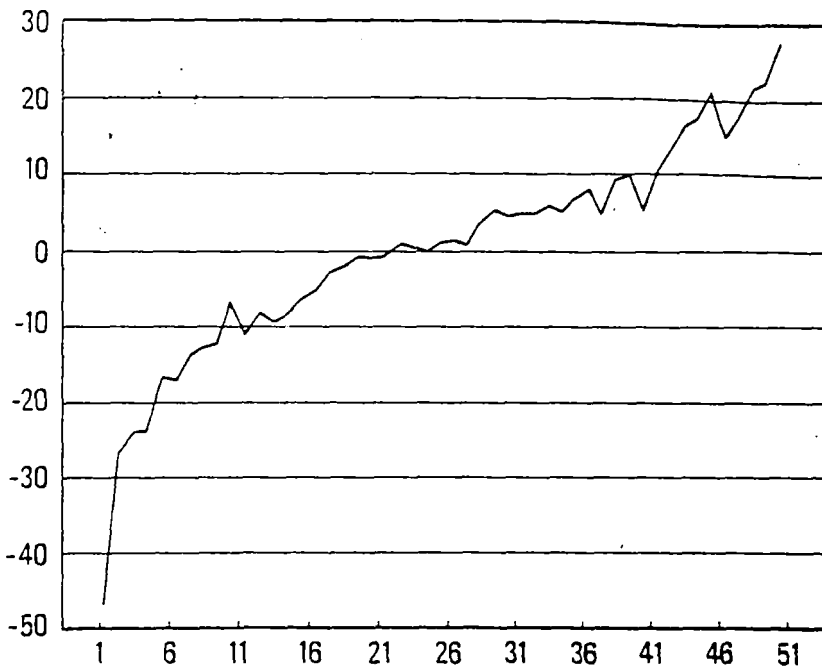
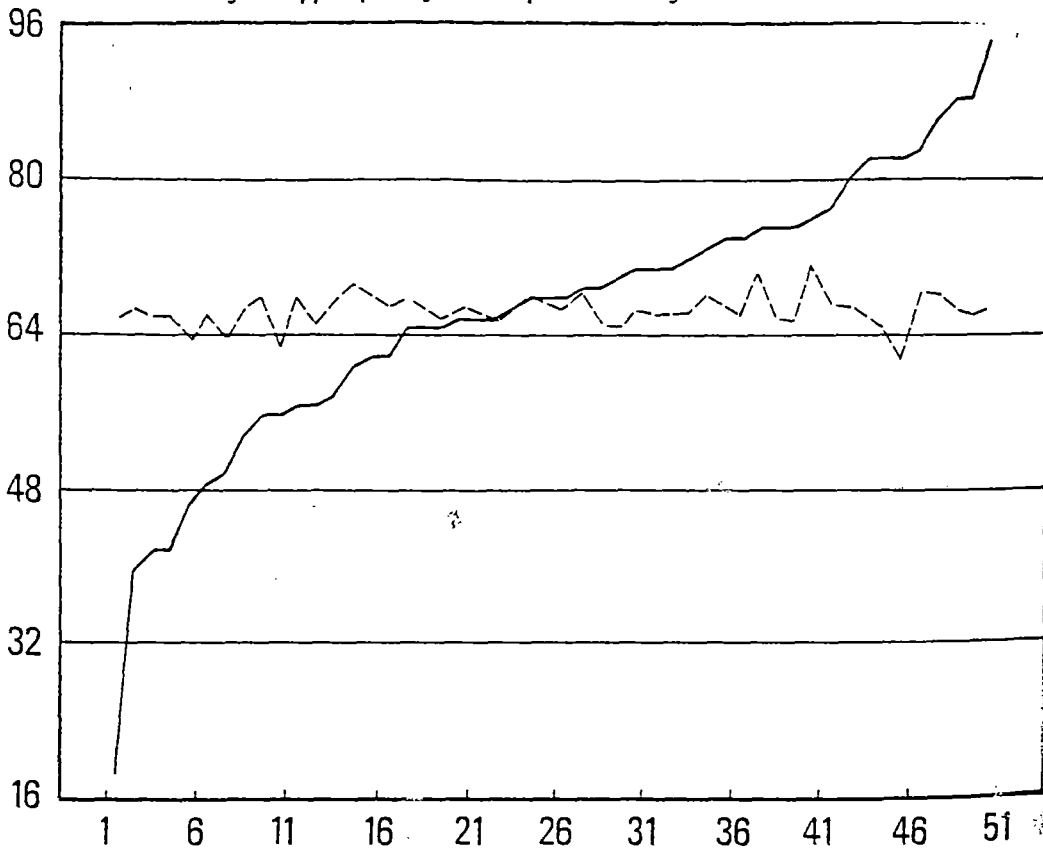
$$\text{με } R^2 = -0,57 \quad D-W = 0,609$$

$$SSE = 15321 \quad s = 17,68$$

Σημειώνεται ότι η αρνητική τιμή του R^2 στο δεύτερο μοντέλο οφείλεται στο ότι εφαρμόσαμε τον ίδιο τύπο για τα δύο μοντέλα.

Στα μοντέλα αυτά η χαμηλή τιμή του Durbin-Watson ερμηνεύεται ως εξής: Οι παρατηρήσεις Y_i αυτοσυσχετίζονται (ειδικότερα $\Gamma_1 = 0,807$) ενώ οι X_i όχι. Εξάλλου, ο συντελεστής συσχέτισης των (X_i, Y_i) είναι χαμηλός επομένως η X δεν ερμηνεύει την συμπεριφορά των παρατηρήσεων X_i . Αυτό έχει ως αποτέλεσμα τα κατάλοιπα της εκτίμησης να αυτοσυσχετίζονται οπότε και η τιμή του $D-W$ είναι μικρή.

Οι ακόλουθες γραφικές παραστάσεις είναι διαφωτιστικός.



Γραφική παράσταση των $\hat{\epsilon}_i$ από το μοντέλο $Y_i = b_0 + b_1 X_i + \epsilon_i$

ΕΦΑΡΜΟΓΗ

Εξετάζουμε την εξής εταιρεία:

V.A.X A.E

Προϊόν : υδρύαλος

Αποτελείται από άμμο (χαλαζιακή) και σόδα (ανθρακική)

και τα οποία αποτελούν τις πρώτες ύλες.Βρίσκονται σε θερμοκρασία κλιβάνου

1200-1500 C⁰ και δημιουργούν μια άμορφη μάζα γυαλιού τη βάζουμε σε κάποια μηχανήματα που λέγονται κλίβανοι υψηλής ατμοσφαιρικής πίεσης,με αποτέλεσμα από τη μεγάλη θερμοκρασία να λιώνει και να γίνεται υγρό.

Αυτό λέγεται υδρύαλος και έχει τις εξείς εφαρμογές :

1. Χρησιμοποιείται σαν κόλλα σε χαρτοβιομηχανίες
2. Χρησιμοποιείται σε βιομηχανίες απορρυπαντικών (δημιουργία πράσινων και μπλέ κόκκων)
- 3.Χρησιμοποιείται σε τεχνικές κατασκευές (φράγματα -τούνελ).

Πωλείται σε βυτία ή σε βαρέλια 25-30 τόννων και η δυναμικότητα της εταιρείας είναι 80-90.000 κιλά την ημέρα .

Δηλαδή η δυναμικότητα του υδρίαλου κυμαίνεται σε 31.025 τόννους ετησίως

(85 χ 365 ημέρες).

Η απορροφητικότητα (ζήτηση) της αγοράς στην Ελλάδα είναι γύρω στους 15.000

τόννους που σημαίνει ότι με βάση τις 20 εργάσιμες ημέρες χ 12 μήνες βγάζουμε με

20.0 τόννους .

Βέβαια ,ανάλογα με τις βάρδιες εξαρτάται και η παραγωγή.

Επίσης ,υπάρχει πάντα ένα στόκ (stock) που είναι 500-1.000 τόννους το στοκ (stock)

είναι 1 ½ εβδομάδας.

Φυσικά ,η παραγωγή εξαρτάται από την ζήτηση .

Το δυναμικό της εταιρείας τα τελευταία 5 χρόνια είναι σταθερό.Στο τομέα της παραγωγής απασχολούνται 7 εργάτες ,οι οποίοι είναι :ένας χημικός ,ένας θερμαστής ,

ένας ηλεκτροτεχνίτης ,ένας χειριστής γερανογέφυρες και τρεις (3) ανειδίκευτοι εργάτες.

Εξετάζουμε την παραγόμενη ποσότητα στα έτη 1991-1995.

Δηλαδή:

1991: Η παραγωγή είναι 14.000 τόννοι

1992: Η παραγωγή είναι 15.000 τόννοι

1993: Η παραγωγή είναι 16.000 τόννοι

1994: Η παραγωγή είναι 16.000 τόννοι

1995: Η παραγωγή είναι 17.000 τόννοι

Η συντήρηση γίνεται από εξωτερικά συνεργεία.

Το Κόστος Παραγωγής αναλύεται ως εξής:

1. Κόστος Α' υλών: 100 κιλά υδρύαλου περιέχουν 40-42 κιλά στέρεη υδρύαλος το οποίο 80% είναι άμμος και 50% είναι σόδα (υπάρχει κάποια φύρα). Το κόστος της σόδας είναι 48 δρχ. το κιλό και της άμμου είναι 12 δρχ. το κιλό .

2. Γ.Β.Ε είναι:

α. Εργατικά = 7 εργάτες / κόστος = 2.200.000 δρχ. το μήνα χ 14 μήνες.

β. Μαζούτ. = κόστος ,57 δρχ. το κιλό /χρειάζεται 12 κιλά μαζούτ ο τόνος του προϊόντος.

γ. Δ.Ε.Η = 400.000 δρχ. το μήνα

δ. Η συντήρηση είναι 200.000 δρχ. το μήνα

3. Η μεταφορά : 2,5 δρχ το κιλό

4. Δοικητικά έξοδα: 1.200.000 δρχ. το μήνα

5. Χρηματοοικονομικά έξοδα: 8 δρχ. το κιλό

6. Αποσβέσεις (βιομηχανικό κόστος): 1,5 δρχ. στα Γ.Β.Ε και 0,5 στα Διαθέσεως.

Η τιμή πώλησης στις βιομηχανίες απορρυπαντικών είναι 27 δρχ. το κιλό, στις κατασκευαστικές εταιρείες (επειδή υπάρχει μεγάλη απόσταση) είναι 35 δρχ. το κιλό και για τις κόλλες σε χαρτοβιομηχανίες είναι 40 δρχ. το κιλό.

Στις βιομηχανίες απορρυπαντικών αποτελεί το 60% της παραγωγής ,στις κατασκευαστικές 30% και στις χαρτοβιομηχανίες (κόλλες) το 10%.

Η τιμή του προϊόντος

Το 1991 ήταν κατά μέσο όρο 40 δρχ. το κιλό

Το 1992 43 δρχ. το κιλό

Το 1993 45 δρχ. το κιλό

Το 1994 40 δρχ. το κιλό

Το 1995 35 δρχ. το κιλό

Στον πίνακα που ακολουθεί δίνονται το κόστος παραγωγής και η τιμή πώλησης σε σταθερές τιμές:

Ζητείται

- i. Να εκτιμηθεί ένα γραμμικό μοντέλο για την τιμή πώλησης με ερμηνευτική μεταβλητή το κόστος παραγωγής.
- ii. Να ερμηνευτούν οι παράμετροι του μοντέλου
- iii. Να υπολογιστούν τα τυπικά σφάλματα των b_0 και b_1 καθώς και ο συντελεστής προσδιορισμού.
- iv. Να ελεγχθεί η υπόθεση ότι οι μεταβλητές X και Y ισούνται γραμμικά με $0,7$ στο επίπεδο σημαντικότητας $\alpha = 0,05$
- v. Να ελεγχθεί η υπόθεση ότι οι μεταβλητές X και Y είναι ίσες στο επίπεδο σημαντικότητας $\alpha = 0,01$
- vi. Να υπολογιστεί ο συντελεστής αυτοσυσχέτισης.

Ε Τ Ο Σ	ΕΤΗΣΙΟ ΚΟΣΤΟΣ ΠΑΡΑΓΩΓΗΣ X_i	ΕΤΗΣΙΑ ΤΙΜΗ ΠΩΛΗΣΗΣ Y_i
1991	146,200	116,529
1992	147,468	121,826
1993	162,485	128,463
1994	174,825	138,954
1995	190,871	149,849

ΛΥΣΗ

i. Συμβολίζουμε με Y την τιμή πώλησης και με X το κόστος παραγωγής και θα εκτιμήσουμε το μοντέλο.

$$Y_i = b_0 + b_1 X_i + E_i$$

Για το σκοπό αυτό υπολογίζουμε:

$$S X_i = 3580,765$$

$$S y_i = 2827,419$$

$$S X_i Y_i = 727848,1$$

$$\bar{X} = 238,718$$

$$\bar{Y} = 188,465$$

$$S X_i^2 - n \bar{X}^2 = 70393,46$$

$$S X_i^2 = 925187,71$$

$$S Y_i^2 = 572826,52$$

$$S Y_i^2 - n \bar{Y}^2 = 39871,04$$

οπότε

$$\hat{b}_1 = \frac{S X_i Y_i - n \bar{X} \bar{Y}}{S X_i^2 - n \bar{X}^2} = \frac{727848,1 - 15(238,718)(188,495)}{70393,46} = 0,75137$$

$$\hat{b}_0 = \bar{Y} - \hat{b}_1 \bar{X} = 188,495 - (0,75137)(238,718) = 9,1297$$

ii. Στο μοντέλο αυτό ο συντελεστής $\hat{b}_1 = 0,75137$ είναι η μέση τιμή πώλησης.

Δηλαδή για κάθε αύξηση του ετήσιου κόστους παραγωγής κατά μία μονάδα ή αναμενόμενη αύξηση της τιμής πώλησης ισούται με 0,75137 ή 751,37 εκ.δρχ του 1993. Ο συντελεστής $\hat{b}_0 = 9,1297$ θα μπορούσε να ερμηνευτεί ως αυτόνομο τμήμα της τιμής της πώλησης δηλαδή το ύψος της όταν το κόστος παραγωγής

ισούτε με μηδέν. Όμως η τιμή $X=0$ είναι πολύ μακριά από το εύρος των δεδομένων μας και επομένως η προέκταση του μοντέλου ως εκεί είναι παρακινδυνευμένη.

iii. Υπολογίζουμε

$$\begin{aligned}
 S^2 &= \frac{1}{n-2} S\hat{\epsilon}_i^2 \\
 &= \frac{1}{n-2} (SY_i^2 - b_0 S Y_i - b_1 S X_i Y_i) \\
 &= \frac{1}{3} (572826,52) - (9,1297) (2827,419) - \\
 &\quad -(0,75137) \cdot (727848,1) \\
 &= \frac{130,98}{3} = 10,075
 \end{aligned}$$

και

$$S = 10,075 = 3,174$$

Επομένως

$$S\hat{b}_1 = S \sqrt{\frac{1}{S X_i^2 - n X^2}} = 3,174 \sqrt{\frac{1}{70393,46}} = 0,01196$$

$$S\hat{b}_0 = S \sqrt{\frac{1}{n} + \frac{X^2}{S X_i - n X^2}} = 3,174 \sqrt{\frac{1}{3} + \frac{(238,718)^2}{70393,46}} = 2,97$$

$$R^2 = 1 - \frac{S \epsilon_i^2}{S Y_i^2 - n Y^2} = 1 - \frac{130,98}{39871,04} = 0,9934$$

Τώρα μπορούμε να παρουσιάσουμε τα αποτελέσματα της εκτίμησης σύμφωνα εεεμε την πάγια τακτική, δηλ. ως εξής:

$$Y = 9,1297 + 0,75137X$$

(3,07) (62,80)

μέ $R^2=9,1297$ $d=1,191597$

$SSE=130,98$ $S=3,17$

Για $n=15$ και $k=1$ βρίσκουμε στον πίνακα 8 στο τέλος της εργασίας τις τιμές $d_l=1,08$ και $d_u=1,36$. Ισχύει:

$$d_l < d < d_u$$

δηλαδή βρισκόμαστε στην περιοχή αβεβαιότητας του ελέγχου. Ακολουθούμε την καθιερωμένη για τις οικονομικές χρονικές στιγμές τακτική και θεωρούμε ότι και αυτή η περιοχή είναι η περιοχή απόρριψης της μηδενικής υπόθεσης $H_0 : \rho=0$ υπέρ της $H_1: \rho>0$. Η μέθοδος Cochran-Orcutt που περιγράψαμε εφαρμόζετε σταδιακά ως εξής:

- Υπολογίζουμε $\hat{\rho}=1-d/2=1-1,191597/2=0,4042$
- Παίρνουμε τις οιονές διαφορές των δεδομένων , δηλαδή τις

$$Y_i^*=Y_i-0,4042Y_{i-1}, X_i^*=X_i-0,4042X_{i-1}, i=2,\dots,5$$

- και έχουμε:

Έτος	$Y_i^* = Y_i - 0,4042Y_{i-1}$	$X_i^* = X_i - 0,4042X_{i-1}$
1992	74.72498	88.3796
1993	79.22093	102.8784
1994	87.02925	109.1486
1995	93.68379	120.2067

- Εκτιμούμε ξανα το μοντέλο χρησιμοποιώντας ως δεδομένα τις μετασχηματισμένες τιμές (Y_i^* , X_i^*) $i=1$

Επαναληπτική εφαρμογή της διαδικασίας αυτής έδωσε τα ακόλουθα αποτελέσματα (οι αριθμοί στις παρενθέσεις είναι t- αναλογίες των αντίστοιχων εκτιμήσεων):

$$Y = 10,88367 + 0,7463073x$$

(2,09) (37,63)

$R^2 = 0,9968$	$\rho = 0,3756$
$SSE = 108,6$	$tp = 1,24$
$s = 3,14$	$D-W = 1,49$

Ο ζητούμενος έλεγχος θα γίνει σταδιακά ως εξής:

$$H_0: b_1 = 0,7$$

$$H_ε: b_1 = 0$$

Επειδή

$$|t_n| = \frac{\hat{b}_1 - 0,7}{Sb_1} = \frac{0,75137 - 0,7}{0,01196} = 4,29 > t_{n-1, \alpha/2} =$$

$=t_{13(0,025)} = 2,16$ απορρίπτουμε την H_0

Στο μοντέλο

$$Y_i = b_0 + b_1 x_i + \varepsilon_i$$

μέ $E(Y_i) = b_0 + b_1 X_i$

η τιμή πώλησης ισούται με b_1 και το ετήσιο κόστος παραγωγής στο σημείο $x = x_i$ είναι η

$$E(Y_i)/X_i = b_0/X_i + b_1$$

Αν $b_0 = 0$ τότε τιμή πώλησης = κόστος παραγωγής = b_1 .

Επομένως ο ζητούμενος έλεγχος είναι ισοδύναμος με τον έλεγχο της μηδενικής υπόθεσης $b_0 = 0$ που γίνεται ως εξής:

$$H_0: b_0 = 0$$

$$H_1: b_0 > 0$$

επειδή

$$t_n = b_0/s_{b_0} = 9,1297/2,97 = 3,072 > t_{13(0,01)} = 2,65$$

η H_0 απορρίπτεται.

Για να υπολογίσουμε τον συντελεστή αυτοσυσχέτισης πρώτης τάξης στις τιμές X_i καταστρώνουμε τον ακόλουθο πίνακα:
 απ' όπου υπολογίζουμε

$$\bar{X} = \sum_{i=1}^5 X_i / 5 = 151,32$$

$$\sum_{i=1}^5 (X_i - \bar{X})^2 = 2959,4$$

$$\sum_{i=2}^5 (X_i - \bar{X})(X_{i-1} - \bar{X}) = 7487,1$$

$$\sum_{i=3}^5 (X_i - \bar{X})(X_{i-2} - \bar{X}) = 1733,59$$

$$\Gamma_1 = 7487,1 / 2959,4 = 0,759$$

οπότε

$$\Gamma_2 = 1733,59 / 2959,4 = 0,586$$

Με τρόπο ανάλογο υπολογίζουμε τους ακόλουθους δύο συντελεστές για την χρονική σειρά Y_i , $i=1, \dots, 5$

ΕΤΟΣ	X_i	X_{i-1}
1992	115,11	100,00
1993	124,00	115,11
1994	128,17	124,00
1995	133,52	128,17

ΕΤΟΣ	Y_i	Y_{i-1}
1992	174,64	156,80
1993	163,38	174,64
1994	173,67	163,38
1995	177,01	173,67

$$\Gamma_1 = 0,792$$

$$\Gamma_2 = 0,624$$

Είναι εύκολο να ελεγχθεί ότι ισχύει :

$$\Gamma_{xy} = 0,83$$

Έτσι, υπολογίσαμε τους συντελεστές αυτοσυσχέτισης πρώτης και δεύτερης τάξης για:

- τον αριθμό της ετήσιας παραγωγής $\Gamma_1 = 0,759$ και $\Gamma_2 = 0,586$
- τον αριθμό της τιμής πώλησης $\Gamma_1 = 0,792$ και $\Gamma_2 = 0,624$

Επειδή $n=5$ το διάστημα μέσα στο οποίο κρίνεται στατιστικά μη σημαντική η τιμή των r_i είναι το

$$(-0,23 \quad 0,23)$$

και άρα όλοι οι συντελεστές που υπολογίσαμε είναι στατιστικά σημαντικοί στο επίπεδο σημαντικότητας 2%.

Συμπερασματικά θα μπορούσαμε να πούμε ότι η πρόβλεψη για τα επόμενα 3-5 έτη, είναι ότι η εταιρία όχι μόνο θα είναι βιώσιμη αλλά θα έχει ανοδική πορεία και αν λάβουμε υπόψη τους Ισολογισμούς των τελευταίων 5 ετών, το καθαρό κέρδος της επιχείρησης θα κυμαίνεται (και στα επόμενα 3-5 έτη) από 57.200.000 δρχ έως 62.350.000 δρχ.

Επομένως, κάτω από αυτούς τους υπολογισμούς χρησιμοποιώντας την στατιστική μέθοδο της συσχέτισης -παλινδρόμησης μπορούμε να μιλάμε για μια υγιή και κερδοφόρα επιχείρηση ορθά στεκούμενη στην ανταγωνιστική μας κοινωνία .

Παρακάτω παραθέτονται οι Ισολογισμοί της εταιρίας ΑΦΟΙ Δ. ΣΗΜΠΗ Α.Ε κατά τα έτη 1991-1995.

ΑΝΟΙΧΤΟ ΔΕΛΤΙΟ ΑΝΕΛΕΥΣΕΩΣ ΚΑΙ ΕΠΙΧΕΙΡΗΣΙΑΚΩΝ ΑΠΟΤΕΛΕΣΜΑΤΩΝ
 ΓΕΩΛΟΓΙΚΟΣ ΣΤΟΙΧΕΙΟΛΟΓΟΣ 31ης ΔΕΚΕΜΒΡΙΟΥ 1990 - 10η ΕΤΑΙΡΙΚΗ ΧΡΗΣΗ
 (1 ΙΑΝΟΥΑΡΙΟΥ - 31 ΔΕΚΕΜΒΡΙΟΥ 1991)
 ΕΣΒΑΡΧΑΤΟ ΒΗΜΑ ΑΡΙΘ. ΜΗΤΡΩΟΥ ΝΟΜΑΡΧΙΑΣ ΒΟΙΩΤΙΑΣ :
 6662 / 11 / 8 - 66 / 0

ΜΗΝΙΑΙΟ	Ποσό προηγούμενης χρήσεως 1990		Ποσό παλαιότερης χρήσεως 1991		ΠΑΡΗΛΗΛΟ	Ποσό προηγ. χροσ. 1990	Ποσό παλαιότερ. χροσ. 1991
	Αξία-ετήσ.	Αποθ.-Αναπόδραστη αξία	αξία-ετήσ.	Αποθ.-Αναπόδραστη αξία			
ΕΣΒΑ ΕΓΧΑΡΤΩΣΕΩΣ					Α. ΖΑΙΑ ΚΕΦΑΛΑΙΑ		
Εξόδα οδοών κ' επί των οδών δασύων δασείων περ/κής ηραϊδίου	164.813.307	49.620.090	115.193.217	117.428.950	1. Μετοχικά κεφάλαια (195.000 μετ. των 1000 δολ)	195.000.000	195.000.000
<u>181.162.291</u>	<u>165.734.672</u>	<u>175.438.719</u>	<u>233.798.034</u>	<u>50.343.240</u>	III. Διεσωτάς ανεπισφραγιστής	70.097.470	70.097.470
					IV. Αποθεματικά κεφάλαια		
ΑΓΙΟ ΕΝΕΡΓΗΤΙΚΟ					1. Τελικό αποθεματικό	2.450.310	2.450.310
Επιδόματα οδοών	59.494.050	-	59.494.050	59.494.050	2. Ειδικό αποθεματικό	40.000	40.000
Εξόδα οδοών	110.423.009	45.467.023	64.935.086	118.953.008	4. Επένδυση αποθεματικά	-	-
Εξόδα οδοών	387.103.749	187.136.913	179.966.836	432.958.014	5. Αποσπαστικά αποθεματικά ειδικών διατάξεων νόμων	60.004.530	60.004.530
Εξόδα οδοών	24.904.708	14.397.574	10.507.134	23.038.578		62.503.849	62.503.849
Εξόδα οδοών	14.091.164	7.504.824	7.186.540	15.430.832			
Εξόδα οδοών κ' καταβολές	1.882.353	225.082	1.656.47	1.882.353			
<u>578.499.033</u>	<u>254.752.516</u>	<u>323.747.617</u>	<u>651.756.045</u>	<u>206.124.411</u>			
					V. Αποτελέσματα εις νέο	(135.927.609)	(55.784.012)
ΟΛΟΚΛΗΡΩΤΟ ΕΝΕΡΓΗΤΙΚΟ			<u>323.747.617</u>		<u>375.632.436</u>		
ΟΛΟΚΛΗΡΩΤΟ ΕΝΕΡΓΗΤΙΚΟ							
Ποσά υπολειμμάτων κ' υπολειμμάτων							
			70.465.555		58.466.082		
			24.589.794		4.741.951		
			9.940.786		4.761.900		
			100.896.135		75.969.933		
			58.214.308		134.927.032		
			30.254.850		4.211.645		
			60.115.048		4.062.040		
			70.469.898		8.874.300		
			3.766.719		625.000		
			40.627.492		47.443.699		
			51.594.529		12.507.401		
			118.955.466		164.374.432		
			6.887.174		6.818.185		
			350.415.686		376.370.765		
			3.770.342		2.525.438		
			543.661		246.393		
			4.314.203		2.771.831		
			55.625.924		499.012.529		
			904.811.660		1.054.090.717		

ΑΝΑΛΥΣΗ ΛΟΓΑΡΙΑΣΜΩΝ ΑΠΟΤΕΛΕΣΜΑΤΑ ΧΡΗΣΕΩΣ

ΠΙΝΑΚΑΣ ΔΙΑΦΕΡΕΣΕΩΣ ΑΠΟΤΕΛΕΣΜΑΤΩΝ

ΑΠΟΤΕΛΕΣΜΑΤΑ ΕΚΜΕΤΑΛΛΕΥΣΕΩΣ	Ποσ. προηγ. χρ. 1990		Ποσ. παλ. χρ. 1991		Ποσ. προηγ. χρ. 1990		Ποσ. παλ. χρ. 1991	
	Ποσ.	Αξ.	Ποσ.	Αξ.	Ποσ.	Αξ.	Ποσ.	Αξ.
Εξόδα οδοών	454.346.616	174.825.749	491.201.025	225.271.759	Καθαρά αποτ/τα (ζημιές) χρήσεως	135.927.609	55.781.014	
Αποτ/τα οδοών	48.921.452	24.721.452	48.921.452	24.721.452	Υπόλοιπα αποτ/των (ζημιών) προηγ. χρήσεων	30.293.514	172.721.123	
Εξόδα οδοών	82.788.292	42.788.292	82.788.292	42.788.292	Υπόλοιπα ζημιών εις νέο	172.721.123	228.501.135	
Εξόδα οδοών	45.677.945	22.677.945	45.677.945	22.677.945				
Εξόδα οδοών	109.312.710	54.312.710	109.312.710	54.312.710				
Αποτ/τα (ζημιές) οδοών	121.295.977	60.295.977	121.295.977	60.295.977				
Εξόδα οδοών	104.084.204	52.084.204	104.084.204	52.084.204				
Αποτ/τα (ζημιές) οδοών	106.704.464	53.704.464	106.704.464	53.704.464				
Αποτ/τα (ζημιές) οδοών	108.030.000	54.030.000	108.030.000	54.030.000				
ΣΥΝΕΚΤΑΚΤΑ ΑΠΟΤΕΛΕΣΜΑΤΑ								
Κατά η ανόργανα χροσ	7.897.203		5.570.606					
κατά η μεταβολή αποτ/τα (ζημιές)	(135.927.609)		(55.784.012)					
Αποτ/τα οδοών	77.028.212		76.289.205					
Αποτ/τα οδοών	77.028.212		76.289.205					
Αποτ/τα οδοών (ΖΗΜΙΕΣ)	(135.927.609)		(55.784.012)					

ΠΡΟΕΔΡΟΣ ΤΟΥ Δ.Σ. ΥΠΑΤΟΣ ΘΗΒΑΣ 28 (ΙΑΝΟΥΑΡΙΟΥ 1992) Ο ΥΠΕΥΘΥΝΟΣ ΛΟΓΙΣΤΗΡΙΟΥ
 ΓΙΑΔΗΣ ΕΜΠΛΗΤΗΣ ΤΟΥ ΔΙΕΥΘΥΝΤΗ ΑΝΔΡΕΑΣ ΜΑΡΚΟΣ ΤΟΥ ΝΙΚΟΛΑΟΥ
 67044/1ΣΤ ΑΛΛΗΝΩΝ Σ 120606/ΚΕΤ ΑΛΛΗΝΩΝ

ΑΦΟΙ Δ.ΣΗΜΙΤΗ ΑΝΩΝΥΜΗ ΕΤΑΙΡΕΙΑ

ΕΔΡΑ: ΥΠΑΤΟ ΘΗΒΑΣ ΑΡΙΘ.ΜΗΤΡ.ΝΟΜΑΡΧΙΑΣ ΒΟΙΩΤΙΑΣ 9992/11/Β/86/0 - ΙΣΟΛΟΓΙΣΜΟΣ 31/12/1995 - 20Η ΕΤΑΙΡΙΚΗ ΧΡΗΣΗ (1/1-31/12/95) (ΠΟΣΑ ΣΕ ΔΡΧ.)

ΚΩΔ	Ποσά Κλειομένης Χρήσεως 1995			Ποσά Προηγούμενης Χρήσεως 1994			ΠΑΘΗΤΙΚΟ	Ποσά Κλειομένης Ποσά Προηγ/νης	
	Αξία Κτίσης	Αποσβέσεις	Ανοπ. Αξία	Αξία Κτίσης	Αποσβέσεις	Ανοποσβ. Αξία		Χρήσεως 1995	Χρήσεως 1994
ΕΓΚΑΤΑΣΤΑΣΕΩΣ							A.ΙΔΙΑ ΚΕΦΑΛΑΙΑ		
ιδιούσεως & πρώτης εγκατ.	16.369.084	11.034.306	5.334.778	16.369.084	11.034.306	5.334.778	1.Μετοχικά Κεφάλαια		
Συνείων κατασκ.περιόδου	1.372.505.582	154.900.777	1.217.604.805	394.553.671	154.900.777	239.652.894	(195.000 μετοχές των 1.000δρχ.)		
	<u>1.388.874.666</u>	<u>165.935.083</u>	<u>1.222.939.583</u>	<u>410.922.752</u>	<u>165.935.083</u>	<u>244.987.672</u>	1.Κατοβλημένο	195.000.000	195.000.000
ΕΝΕΡΓΗΤΙΚΟ							III.Διαφορές αναπροσαρμογής		
αίτες ακινητοποιήσεις							Επιχορηγήσεις επενδύσεων		
Γ-γκικόπεδα	109.072.425	---	109.072.425	109.072.425	---	109.072.425	3.Επιχορηγήσεις επενδύσεων παγίου ενεργ.	73.574.470	73.574.470
& τεχνικά έργα	228.093.248	102.148.479	125.944.769	228.093.248	86.743.368	141.349.880	IV.Αποθεματικά Κεφάλαια		
μητα-τεχν.εγκαταστάσεις	432.958.024	331.366.701	101.591.323	432.958.024	296.569.551	136.388.473	1.Τακτικά αποθεματικά	2.459.319	2.459.319
ορικά μέσα	21.783.517	17.587.632	4.195.885	21.783.517	17.587.632	4.195.855	3.Ειδικά αποθεματικά	40.000	40.000
α & λοιπός εξοπλισμός	15.980.832	14.474.521	1.506.311	15.980.832	14.474.521	1.506.311	5.Αφορολόγητα αποθεματικά ειδ.διατ.νόμων	60.004.530	60.004.530
ομοιότητες υπό εκτέλεση								<u>62.503.849</u>	<u>62.503.849</u>
καταβολές	1.882.353	903.528	978.825	1.882.353	903.528	978.825	V.Αποτελέσματα εις νέο		
ακινήτων (ΓII)	<u>809.770.399</u>	<u>466.480.861</u>	<u>343.289.538</u>	<u>809.770.399</u>	<u>416.278.600</u>	<u>393.491.799</u>	Διαφορά φορολ.ελέγχου	-28.243.366	---
ΦΟΡΟΥΝ ΕΝΕΡΓΗΤΙΚΟ							Υπόλοιπο ζημιών προηγ.χρήσεως	-243.896.019	(292.073.206)
α)α							Μείον κέρδη χρήσεως	+ 51.897.465	48.177.184
α'τα έτοιμα & ημιτελή							Υπόλοιπο ζημιών προηγ.χρήσεως	-214.957.472	243.896.019
αίδια & υπολείμματα			----			167.700	Γ.ΥΠΟΧΡΕΩΣΕΙΣ		
βοηθητικές ύλες-αναλώσιμα			21.694.537			29.700.781	1.Μακροπρόθεσμες Υποχρεώσεις		
νταλικά & είδη ουσκευασίας			<u>21.694.537</u>			<u>29.700.781</u>	2.Δάνεια Τραπεζών (ΕΤΒΑ)	1.616.631.956	708.680.045
εις							II.Βραχυπρόθεσμες Υποχρεώσεις		
εις		73.626.846				55.073.915	1.Προμηθευτές	26.503.219	51.966.665
ατα Εισπρακτέα							2.Γραμμάτια πληρωτέα	5.439.355	15.721.385
υλακίου		240.000				1.000.000	2α.Επιταγές πληρωτέες	---	3.312.824
απτερες για εγγύηση		----				80.000	3.Τράπεζες λβροχρημ.υποχρ.	89.719.760	89.764.204
ατα σε καθυστέρηση			281.004.418			282.789.000	5.Υποχρεώσεις από φόρους - τέλη	36.614.866	16.960.332
εις διάφοροι							6.Ασφαλιστικοί οργανισμοί	13.368.218	27.044.960
αμοιβή διαχειρίσεως			-5.263.478			-5.453.403	11.Πρωτίτες διάφοροι	22.411.073	9.852.141
αβολών & πιστώσεων			<u>349.607.786</u>			<u>333.489.512</u>	Σύνολο Υποχρεώσεων (ΓI+ΓII)	<u>1.810.688.447</u>	<u>923.302.556</u>
αα							Δ.ΜΕΤΑΒΑΤΙΚΟΙ ΛΟΓΑΡΙΑΣΜΟΙ ΠΑΘΗΤΙΚΟΥ		
εις όρους & προθεσμίας		6.208.072				7.258.423	3.Λοιποί μεταβατικοί λογαριασμοί	27.892.240	---
		<u>10.952.918</u>				<u>788.969</u>			
α λοφορσάντης εστρατη πικσ(ΔI - ΔII + ΔIV)		17.160.090				9.647.392			
ΚΩΔΟ ΕΝΕΡΓΗΤΙΚΟΥ (Σ+Γ+Δ)		<u>1.853.891.534</u>				<u>1.010.484.856</u>			
α							ΓΕΝΙΚΟ ΣΥΝΟΛΟ ΠΑΘΗΤΙΚΟΥ (Α+Γ+Δ)	<u>1.954.691.534</u>	<u>1.010.484.856</u>

ΙΣ: 1. Επί των ακινήτων της εταιρείας έχουν εγγραφεί υποθήκες και προσημιώσεις: α)Υπέρ του Ελληνικού Δημοσίου δρχ. 61.974.541 και β) Σε ασφάλεια Τραπεζικών δανείων ποσού δρχ. 510.062.000

ΑΝΑΛΥΣΗ ΛΟΓΑΡΙΑΣΜΟΥ ΑΠΟΤΕΛΕΣΜΑΤΩΝ ΧΡΗΣΕΩΣ		ΠΙΝΑΚΑΣ ΔΙΑΘΕΣΕΩΣ ΑΠΟΤΕΛΕΣΜΑΤΩΝ		
	Ποσά Κλειομένης Χρήσεως 1995	Ποσά Προηγούμενης Χρήσεως 1994	Ποσά Κλειομένης Ποσά Προηγ.	
			Χρήσεως 1995 Χρήσεως 1994	
ατα εκμεταλλεύσεως	511.139.587	480.687.928	-28.243.366	---
ασιών	410.218.715	370.775.864	-243.896.019	-292.073.206
ατος πωλήσεων	100.920.872	1.600.000	57.181.913	48.177.187
αλέσματα (κέρδη) εκμεταλλεύσεως	---	111.512.064	-214.957.472	-243.896.019
α έσοδα εκμεταλλεύσεως	100.920.872			
αδα διοικητικής λειτουργίας	18.887.216	38.554.560		
αδα λειτουργίας διαθέσεως	26.697.872	21.817.285	60.371.845	51.140.219
αλέσματα (ζημιές) εκμεταλλεύσεως	45.585.088			
αεωστικοί τόκοι & συναφή έξοδα	55.335.784			
αλέσματα (ζημιές) εκμεταλλεύσεως	1.674.633			
ατακτα αποτελέσματα	53.661.151			
ακτα & ανάργα έσοδα	3.520.762	(854.820)		
ακτα αποτελέσματα (ζημιές)	57.181.913	48.177.187		
αλο αποσβέσεων παγίων στοιχείων	50.202.261	---		
α: Οι από αυτές ενσωματωμένες		---		
ατα λειτουργικό κόστος	50.202.261	---		
αΤΕΛΕΣΜΑΤΑ (ΖΗΜΙΕΣ) ΧΡΗΣΕΩΣ ΠΡΟ ΦΟΡΩΝ	<u>57.181.913</u>	<u>48.177.187</u>		

ΕΤΑΙΡΙΑ ΜΕΛΟΣ ΤΟΥ Δ.Σ. ΜΙΑΤΠΑΔΗΣ ΣΗΜΙΤΗΣ ΤΟΥ ΔΗΜΗΤΡΙΟΥ Η 678640/ΙΣΤ' ΑΘΗΝΩΝ
 ΕΤΑ ΜΕΛΟΣ ΤΟΥ Δ.Σ. ΔΗΜΟΚΡΙΤΟΣ ΣΗΜΙΤΗΣ ΤΟΥ ΣΠΥΡΙΔΩΝΟΣ Η 055518/ΙΣΤ' ΑΘΗΝΩΝ
 Ο ΠΡΟΪΣΤΑΜΕΝΟΣ ΟΙΚΟΝΟΜΙΚΩΝ ΥΠΗΡΕΣΙΩΝ ΑΝΔΡΕΑΣ ΜΑΡΚΟΣ ΤΟΥ ΝΙΚΟΛΑΟΥ Ξ 120607/ΚΣΤ' ΑΘΗΝΩΝ
 GNN I.T.D. ☎ 8222.848

ΑΦΟΙ Δ.ΣΗΜΙΤΗ ΑΝΩΝΥΜΗ ΕΤΑΙΡΕΙΑ

ΔΡΑ: ΥΠΑΤΟ ΘΗΒΑΣ ΑΡΙΘ.ΜΗΤΡ.ΝΟΜΑΡΧΙΑΣ ΒΟΙΩΤΙΑΣ 9992/11/Β/86/0 - ΙΣΟΛΟΓΙΣΜΟΣ 31/12/1995 - 20Η ΕΤΑΙΡΙΚΗ ΧΡΗΣΗ (1/1-31/12/95) (ΠΟΣΑ ΣΕ ΔΡΧ.)

	Ποσό Κλεισμένης Χρήσεως 1995			Ποσό Προηγούμενης Χρήσεως 1994			ΠΑΘΗΤΙΚΟ	Ποσό Κλεισμένης Χρήσεως 1995	Ποσό Προηγ/νης Χρήσεως 1994
	Αξία Κτήσης	Αποσβέσεις	Αναπ. Αξία	Αξία Κτήσης	Αποσβέσεις	Αναπασθ. Αξία			
ΚΑΤΑΣΤΑΣΕΙΣ							Α.ΙΔΙΑ ΚΕΦΑΛΑΙΑ		
αξίες & πρώτης εγκατ.	16.369.084	11.034.306	5.334.778	16.369.084	11.034.306	5.334.778	I.Μετοχικό Κεφάλαιο		
των κατασκ.περιόδου	<u>1.372.505.582</u>	<u>154.900.777</u>	<u>1.217.604.805</u>	<u>394.553.671</u>	<u>154.900.777</u>	<u>239.652.894</u>	(195.000 μετοχές των 1.000δρχ.)		
	<u>1.388.874.666</u>	<u>165.935.083</u>	<u>1.222.939.583</u>	<u>410.922.755</u>	<u>165.935.083</u>	<u>244.987.672</u>	1.Καταβλημένο	<u>195.000.000</u>	<u>195.000.000</u>
ΡΥΘΙΤΙΚΟ							III.Διαφορές αναπροσαρμογής		
ακίνητοποιήσεις							Επιχορηγήσεις επενδύσεων		
κόπεδα	109.072.425	---	109.072.425	109.072.425	---	109.072.425	3.Επιχορηγήσεις επενδύσεων παγίου ενεργ.	<u>73.574.470</u>	<u>73.574.470</u>
γυμνάσια	228.093.248	102.148.479	125.944.769	228.093.248	88.743.368	141.349.880	IV.Αποθεματικά Κεφάλαια		
α-τεχν.εγκαταστάσεις	432.958.024	331.368.701	101.581.323	432.958.024	298.569.551	136.388.473	1.Τακτικό αποθεματικό	2.459.319	2.459.319
ή μέσα	21.783.517	17.587.832	4.195.885	21.783.517	17.587.832	4.195.885	3.Ειδικό αποθεματικό	40.000	40.000
αποσβ.εξοπλισμός	15.980.832	14.474.521	1.508.311	15.980.832	14.474.521	1.508.311	5.Απορολόγητα αποθεματικά ειδ.διατ.νόμων	<u>60.004.530</u>	<u>60.004.530</u>
πρωτογενή								<u>62.503.849</u>	<u>62.503.849</u>
βελτιωτικά	1.882.353	903.528	978.825	1.882.353	903.528	978.825	V.Αποτελέσματα εις νέο		
αποσβέσεων (ΓII)	<u>809.770.399</u>	<u>466.480.861</u>	<u>343.289.538</u>	<u>809.770.399</u>	<u>418.278.600</u>	<u>393.491.799</u>	Διαφορά φορολ.ελέγχου	-28.243.366	---
ΥΠΟ ΕΝΕΡΓΗΤΙΚΟ							Υπόλοιπα ζημιών προηγ.χρήσεων	-243.896.019	(292.073.208)
οικόμοια & ημιτελή							Μείον κέρδη χρήσεως	<u>+ 51.897.485</u>	<u>48.177.184</u>
1 & υπολειμματα						167.700	Υπόλοιπα ζημιών προηγ.χρήσεων	-214.957.472	243.896.019
παικτικές ύλες-αναλώσιμα							Σύνολο Ιδίων Κεφαλαίων	<u>116.120.847</u>	<u>87.182.300</u>
κά & είδη ασκευασίας			21.694.537			29.700.781	Γ.ΥΠΟΧΡΕΩΣΕΙΣ		
			<u>21.694.537</u>			<u>29.700.781</u>	I.Μακροπρόθεσμες Υποχρεώσεις		
							2.Δάνεια Τραπεζών (ΕΤΒΑ)	<u>1.616.631.956</u>	<u>708.680.045</u>
			73.626.846			55.073.915	II.Βραχυπρόθεσμες Υποχρεώσεις		
							1.Προμηθευτές	26.503.219	51.966.665
							2.Γραμμάτια πληρωτέα	5.439.355	15.721.385
			240.000			1.000.000	2α.Επιταγές πληρωτέες	---	3.312.824
							3.Τράπεζες ληροχρημ.υποχρ.	89.719.760	89.764.204
							5.Υποχρεώσεις από φόρους - τέλη	36.514.866	18.960.332
							6.Ασφαλιστικοί οργανισμοί	13.368.218	27.044.960
							11.Πιστωτές διάφοροι	<u>22.411.073</u>	<u>9.852.141</u>
								194.056.491	214.622.511
							Σύνολο Υποχρεώσεων (ΓI+ΓII)	<u>1.810.688.447</u>	<u>923.302.556</u>
							Δ.ΜΕΤΑΒΑΤΙΚΟΙ ΛΟΓΑΡΙΑΣΜΟΙ ΠΑΘΗΤΙΚΟΥ		
							3.Λοιποί μεταβατικοί λογαριασμοί	<u>27.882.240</u>	---
							ΓΕΝΙΚΟ ΣΥΝΟΛΟ ΠΑΘΗΤΙΚΟΥ (Α+Γ+Δ)	<u>1.954.691.534</u>	<u>1.010.484.856</u>

Τι των ακινήτων της εταιρείας έχουν εγγραφεί υποθήκες και προσημειώσεις : α)Υπέρ του Ελληνικού Δημοσίου σχ. 81.974.541 και β) Σε ασφαλεία Τραπεζικών δανείων ποσού δρχ. 510.062.000

ΑΝΑΛΥΣΗ ΛΟΓΑΡΙΑΣΜΟΥ ΑΠΟΤΕΛΕΣΜΑΤΩΝ ΧΡΗΣΕΩΣ			ΠΙΝΑΚΑΣ ΔΙΑΔΕΞΕΩΣ ΑΠΟΤΕΛΕΣΜΑΤΩΝ		
	Ποσό Κλεισμένης Χρήσεως 1995	Ποσό Προηγούμενης Χρήσεως 1994		Ποσό Κλεισμένης Χρήσεως 1995	Ποσό Προηγ. Χρήσεως 1994
πρεταλλεύσεις	511.139.587	460.687.928	Διαφορά φορολ.ελέγχου	-28.243.366	---
αξιών	<u>410.218.715</u>	<u>370.775.964</u>	Ζημιές προηγούμενης χρήσεως	-243.896.019	(292.073.208)
α (κερδη) εκμεταλλεύσεως	<u>100.920.872</u>	<u>109.912.064</u>	Μείον κέρδη χρήσεως	<u>-57.181.313</u>	<u>48.177.184</u>
εκμεταλλεύσεως		<u>1.600.000</u>	Ζημιές χρήσεων εις νέο	-214.957.472	-243.896.019
	<u>100.920.872</u>	<u>111.512.064</u>			
μηχανική λειτουργίας	18.887.218	38.554.560	Ο ΠΡΟΕΔΡΟΣ ΤΟΥ Δ.Σ.		
αυτοκινήτων διαθέσεως	26.897.872	21.817.295	ΜΙΛΙΑΔΗΣ ΣΗΜΙΤΗΣ ΤΟΥ ΔΗΜΗΤΡΙΟΥ		
α (ζημιές) εκμεταλλεύσεως	55.335.784	51.140.219	Η 678640/ΛΣΤ ΑΘΗΝΩΝ		
α τόκοι & συναφή έξοδα	<u>1.674.633</u>	<u>2.108.212</u>	ΕΝΑ ΜΕΛΟΣ ΤΟΥ Δ.Σ.		
α (ζημιές) εκμεταλλεύσεως	53.661.151	49.932.007	ΔΗΜΟΚΡΙΤΟΣ ΣΗΜΙΤΗΣ ΤΟΥ ΣΠΥΡΙΔΩΝΟΣ		
αποτελέσματα & ανόργανα έσοδα	<u>3.520.762</u>	<u>(854.820)</u>	Η 055518/ΛΣΤ ΑΘΗΝΩΝ		
αποτελέσματα (ζημιές)	57.181.913	48.177.187	Ο ΠΡΟΪΣΤΑΜΕΝΟΣ ΟΙΚΟΝΟΜΙΚΩΝ ΥΠΗΡΕΣΙΩΝ		
αβέσεων παγίων στοιχείων	50.202.281	---	ΑΝΔΡΕΑΣ ΜΑΡΚΟΣ ΤΟΥ ΝΙΚΟΛΑΟΥ		
α αυτές ενσωματωμένες	---	---	Ξ 120607/ΚΣΤ ΑΘΗΝΩΝ		
α λειτουργικά κόστη	50.202.281	---			
ΦΑΝΕΡΑ (ΖΗΜΙΕΣ) ΧΡΗΣΕΩΣ ΠΡΟ ΦΟΡΩΝ	<u>57.181.913</u>	<u>48.177.187</u>			

ΠΕΡΙΕΧΟΜΕΝΑ

ΚΕΦΑΛΑΙΟ ΠΡΩΤΟ : Εισαγωγικές έννοιες	1
1.1 Μοντέλο	1
1.2 Ανάλυση Παλινδρόμησης	2
1.3 Η Γραμμική Συνάρτηση Παλινδρόμησης	3
ΚΕΦΑΛΑΙΟ ΔΕΥΤΕΡΟ : Απλή Γραμμική Παλινδρόμηση	
2.1 Το απλό γραμμικό μοντέλο Παλινδρόμησης	4
2.2 Η μέθοδος των ελαχίστων τετραγώνων	6
2.3 Τα κατάλοιπα ή σφάλματα εκτίμησης	9
2.4 Μέση τιμή και διακύμανση των εκτιμητών ελαχίστων τετραγώνων.....	12
2.5 Πόσο καλοί είναι οι εκτιμητές ελαχίστων τετραγώνων.....	14
2.6 Διαστήματα εμπιστοσύνης και έλεγχος των υποθέσεων	15
2.7 Η ερμηνευτική ικανότητα του μοντέλου	19
2.8 Προβλέψεις	21
2.9 Ανάλυση διακυμάνσεως	25
2.10 Το μοντέλο χωρίς σταθερό όρο	27
ΚΕΦΑΛΑΙΟ ΤΡΙΤΟ : Συσχέτιση	
3.1 Εισαγωγή	35
3.2 Η συνδιακύμανση	36
3.3 Ο συντελεστής συσχέτισης	37
3.4 Ιδιότητες του συντελεστή συσχέτισης	40
3.5 Ερμηνεία του συντελεστή συσχέτισης	42
3.6 Ο συντελεστής συσχέτισης ομαδοποιημένων παρατηρήσεων	44
3.7 Ο θεωρητικός συντελεστής συσχέτισης	47
3.8 Η διμεταβλητή κανονική κατανομή	48
3.9 Έλεγχος των υποθέσεων και διαστήματα εμπιστοσύνης	51
3.10 Παλινδρόμηση και Συσχέτιση	56
3.11 Ο συντελεστής συσχέτισης κατά τάξεις	58
3.12 Οι συντελεστές αυτοσυσχέτισης	62

ΚΕΦΑΛΑΙΟ ΤΕΤΑΡΤΟ : Το πολλαπλό μοντέλο Παλινδρόμησης

.....	
4.1 Εισαγωγή	66
4.2 Η εκτίμηση των παραμέτρων με τη μέθοδο των ελαχίστων τετραγώνων	67
4.3 Οι συντελεστές μερικής παλινδρόμησης	69
4.4 Τα κατάλοιπα της εκτίμησης	72
4.5 Ιδιότητες των εκτιμητών ελαχίστων τετραγώνων	76
4.6 Ανάλυση διακυμάνσεως	78
4.7 Ο συντελεστής R^2	80
4.8 Οι συντελεστές μερικού προσδιορισμού	85
4.9 Έλεγχος των στατιστικών υποθέσεων	88
4.10 Ο ολικός έλεγχος σημαντικότητας του μοντέλου	91

ΚΕΦΑΛΑΙΟ ΠΕΜΠΤΟ : Ψευδομεταβλητές

5.1 Εισαγωγή	93
5.2 Μεταβλητή του σταθερού όρου στο μοντέλο παλινδρόμησης	93
5.3 Μεταβλητή της κλίσης ή των μερικών κλίσεων του μοντέλου	96
5.4 Πλειότιμες ποιοτικές μεταβλητές	97
5.5 Έλεγχος της σταθερότητας του μοντέλου με ψευδομεταβλητές	102

ΚΕΦΑΛΑΙΟ ΕΚΤΟ : Η αυτοσυσχέτιση των σφαλμάτων

6.1 Εισαγωγή	104
6.2 Η αυτοσυσχέτιση των διαταρακτικών όρων	105
6.3 Το μοντέλο αυτοσυσχέτισης AR(I)	105
6.4 Το στατιστικό Durbin - Watson	108
6.5 Εκτίμηση ενός μοντέλου με AR(I) διαταρακτικούς όρους	111
6.6 Αυτοσυσχέτιση των καταλοίπων λόγω κακής εξειδίκευσης της συνάρτησης παλινδρόμησης	114

ΕΦΑΡΜΟΓΗ

ΒΙΒΛΙΟΓΡΑΦΙΑ

ΒΙΒΛΙΟΓΡΑΦΙΑ

Ελληνική

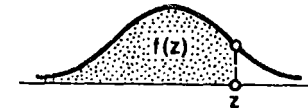
1. Ζαχαροπούλου Χ. (1990), Στατιστική, Μέθοδοι- Εφαρμογές
2. Λουκάκης Μ. (1993), Μαθηματικά Οικονομικών Επιστημών
3. Κεβόρκ Γ. (1991), Συσχέτιση- Παλινδρόμηση, Θεωρία
4. Χρήστου Γ. (1982), Εισαγωγή στην Οικονομετρία

Ξένα

1. Anderson- Sweeney- Williams (1990), Statistics For Business and Economics
2. Anscombe F.J (1973), Graphs in Statistical Analysis
3. Becker W.- Kermedy P. (1992), A Lesson in Least Squares and R Squared

Π Ι Ν Α Κ Ε Σ

ΠΙΝΑΚΑΣ 1: Αθροιστικές πιθανότητες της τυγικής κανονικής κατανομής



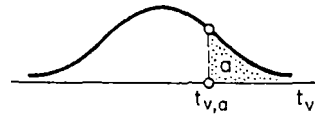
$z/$	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.50000	0.50399	0.50798	0.51197	0.51595	0.51994	0.52392	0.52790	0.53188	0.53586
0.1	0.53983	0.54380	0.54776	0.55172	0.55567	0.55962	0.56356	0.56749	0.57142	0.57535
0.2	0.57926	0.58317	0.58706	0.59095	0.59483	0.59871	0.60257	0.60642	0.61026	0.61409
0.3	0.61791	0.62172	0.62552	0.62930	0.63307	0.63683	0.64058	0.64431	0.64803	0.65173
0.4	0.65542	0.65910	0.66276	0.66640	0.67003	0.67364	0.67724	0.68082	0.68439	0.68793
0.5	0.69146	0.69497	0.69847	0.70194	0.70540	0.70884	0.71226	0.71566	0.71904	0.72240
0.6	0.72575	0.72907	0.73237	0.73565	0.73891	0.74215	0.74537	0.74857	0.75175	0.75490
0.7	0.75804	0.76115	0.76424	0.76730	0.77035	0.77337	0.77637	0.77935	0.78230	0.78524
0.8	0.78814	0.79103	0.79389	0.79673	0.79955	0.80234	0.80511	0.80785	0.81057	0.81327
0.9	0.81594	0.81859	0.82121	0.82381	0.82639	0.82894	0.83147	0.83398	0.83646	0.83891
1.0	0.84134	0.84375	0.84614	0.84849	0.85083	0.85314	0.85543	0.85769	0.85993	0.86214
1.1	0.86433	0.86650	0.86864	0.87076	0.87286	0.87493	0.87698	0.87900	0.88100	0.88298
1.2	0.88493	0.88686	0.88877	0.89065	0.89251	0.89435	0.89617	0.89796	0.89973	0.90147
1.3	0.90320	0.90490	0.90658	0.90824	0.90988	0.91149	0.91309	0.91466	0.91621	0.91774
1.4	0.91924	0.92073	0.92220	0.92364	0.92507	0.92647	0.92785	0.92922	0.93056	0.93189
1.5	0.93319	0.93448	0.93574	0.93699	0.93822	0.93943	0.94062	0.94179	0.94295	0.94408
1.6	0.94520	0.94630	0.94738	0.94845	0.94950	0.95053	0.95154	0.95254	0.95352	0.95449
1.7	0.95543	0.95637	0.95728	0.95818	0.95907	0.95994	0.96080	0.96164	0.96246	0.96327
1.8	0.96407	0.96485	0.96562	0.96638	0.96712	0.96784	0.96856	0.96926	0.96995	0.97062
1.9	0.97128	0.97193	0.97257	0.97320	0.97381	0.97441	0.97500	0.97558	0.97615	0.97670

	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
2.0	0.97725	0.97778	0.97831	0.97882	0.97932	0.97982	0.98030	0.98077	0.98124	0.98169
2.1	0.98214	0.98257	0.98300	0.98341	0.98382	0.98422	0.98461	0.98500	0.98537	0.98574
2.2	0.98610	0.98645	0.98679	0.98713	0.98745	0.98778	0.98809	0.98840	0.98870	0.98899
2.3	0.98928	0.98956	0.98983	0.99010	0.99036	0.99061	0.99086	0.99111	0.99134	0.99158
2.4	0.99180	0.99202	0.99224	0.99245	0.99266	0.99286	0.99305	0.99324	0.99343	0.99361
2.5	0.99379	0.99396	0.99413	0.99430	0.99446	0.99461	0.99477	0.99492	0.99506	0.99520
2.6	0.99534	0.99547	0.99560	0.99573	0.99585	0.99598	0.99609	0.99621	0.99632	0.99643
2.7	0.99653	0.99664	0.99674	0.99683	0.99693	0.99702	0.99711	0.99720	0.99728	0.99736
2.8	0.99744	0.99752	0.99760	0.99767	0.99774	0.99781	0.99788	0.99795	0.99801	0.99807
2.9	0.99813	0.99819	0.99825	0.99831	0.99836	0.99841	0.99846	0.99851	0.99856	0.99861
3.0	0.99865	0.99869	0.99874	0.99878	0.99882	0.99886	0.99889	0.99893	0.99896	0.99900
3.1	0.99903	0.99906	0.99910	0.99913	0.99916	0.99918	0.99921	0.99924	0.99926	0.99929
3.2	0.99931	0.99934	0.99936	0.99938	0.99940	0.99942	0.99944	0.99946	0.99948	0.99950
3.3	0.99952	0.99953	0.99955	0.99957	0.99958	0.99960	0.99961	0.99962	0.99964	0.99965
3.4	0.99966	0.99968	0.99969	0.99970	0.99971	0.99972	0.99973	0.99974	0.99975	0.99976
3.5	0.99977	0.99978	0.99978	0.99979	0.99980	0.99981	0.99981	0.99982	0.99983	0.99983
3.6	0.99984	0.99985	0.99985	0.99986	0.99986	0.99987	0.99987	0.99988	0.99988	0.99989
3.7	0.99989	0.99990	0.99990	0.99990	0.99991	0.99991	0.99992	0.99992	0.99992	0.99992
3.8	0.99993	0.99993	0.99993	0.99994	0.99994	0.99994	0.99994	0.99995	0.99995	0.99995
3.9	0.99995	0.99995	0.99996	0.99996	0.99996	0.99996	0.99996	0.99996	0.99997	0.99997
4.0	0.99997	0.99997	0.99997	0.99997	0.99997	0.99997	0.99998	0.99998	0.99998	0.99998

Πηγή: Olkin I.-Gleser L.-Derman C. (1980), "Probability Models and Applications",

McMillan Publishing Co., Inc.

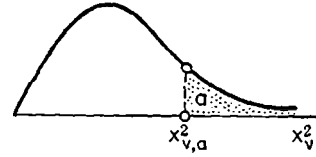
ΠΙΝΑΚΑΣ 2: Κριτικές τιμές της κατανομής t-student



	.100	.050	.025	.010	.005
1	3.078	6.314	12.706	31.821	63.657
2	1.886	2.920	4.303	6.965	9.925
3	1.638	2.353	3.182	4.541	5.841
4	1.533	2.132	2.776	3.747	4.604
5	1.476	2.015	2.571	3.365	4.032
6	1.440	1.943	2.447	3.143	3.707
7	1.415	1.895	2.365	2.998	3.499
8	1.397	1.860	2.306	2.896	3.355
9	1.383	1.833	2.262	2.821	3.250
10	1.372	1.812	2.228	2.764	3.169
11	1.363	1.796	2.201	2.718	3.106
12	1.356	1.782	2.179	2.681	3.055
13	1.350	1.771	2.160	2.650	3.012
14	1.345	1.761	2.145	2.624	2.977
15	1.341	1.753	2.131	2.602	2.947
16	1.337	1.746	2.120	2.583	2.921
17	1.333	1.740	2.110	2.567	2.898
18	1.330	1.734	2.101	2.552	2.878
19	1.328	1.729	2.093	2.539	2.861
20	1.325	1.725	2.086	2.528	2.845
21	1.323	1.721	2.080	2.518	2.831
22	1.321	1.717	2.074	2.508	2.819
23	1.319	1.714	2.069	2.500	2.807
24	1.318	1.711	2.064	2.492	2.797
25	1.316	1.708	2.060	2.485	2.787
26	1.315	1.706	2.056	2.479	2.779
27	1.314	1.703	2.052	2.473	2.771
28	1.313	1.701	2.048	2.467	2.763
29	1.311	1.699	2.045	2.462	2.756
30	1.310	1.697	2.042	2.457	2.750
40	1.303	1.684	2.021	2.423	2.704
60	1.296	1.671	2.000	2.390	2.660
∞	1.282	1.645	1.960	2.326	2.576

Πηγή: P. Newbold (1984), "Statistics for Business and Economics",
Prentice-Hall Inc.

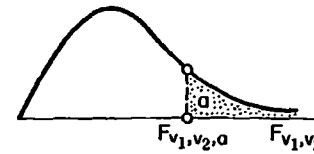
ΠΙΝΑΚΑΣ 3: Κριτικές τιμές της κατανομής χ^2



	.995	.990	.975	.950	.900	.100	.050	.025	.010	.005
1	0.0 ³ 393	0.0 ³ 157	0.0 ³ 982	0.0 ³ 393	0.0158	2.71	3.84	5.02	6.63	7.88
2	0.0100	0.0201	0.0506	0.103	0.211	4.61	5.99	7.38	9.21	10.60
3	0.072	0.115	0.216	0.352	0.584	6.25	7.81	9.35	11.34	12.84
4	0.207	0.297	0.484	0.711	1.064	7.78	9.49	11.14	13.28	14.86
5	0.412	0.554	0.831	1.145	1.61	9.24	11.07	12.83	15.09	16.75
6	0.676	0.872	1.24	1.64	2.20	10.64	12.59	14.45	16.81	18.55
7	0.989	1.24	1.69	2.17	2.83	12.02	14.07	16.01	18.48	20.28
8	1.34	1.65	2.18	2.73	3.49	13.36	15.51	17.53	20.09	21.96
9	1.73	2.09	2.70	3.33	4.17	14.68	16.92	19.02	21.67	23.59
10	2.16	2.56	3.25	3.94	4.87	15.99	18.31	20.48	23.21	25.19
11	2.60	3.05	3.82	4.57	5.58	17.28	19.68	21.92	24.73	26.76
12	3.07	3.57	4.40	5.23	6.30	18.55	21.03	23.34	26.22	28.30
13	3.57	4.11	5.01	5.89	7.04	19.81	22.36	24.74	27.69	29.82
14	4.07	4.66	5.63	6.57	7.79	21.06	23.68	26.12	29.14	31.32
15	4.60	5.23	6.26	7.26	8.55	22.31	25.00	27.49	30.58	32.80
16	5.14	5.81	6.91	7.96	9.31	23.54	26.30	28.85	32.00	34.27
17	5.70	6.41	7.56	8.67	10.09	24.77	27.59	30.19	33.41	35.72
18	6.26	7.01	8.23	9.39	10.86	25.99	28.87	31.53	34.81	37.16
19	6.84	7.63	8.91	10.12	11.65	27.20	30.14	32.85	36.19	38.58
20	7.43	8.26	9.59	10.85	12.44	28.41	31.41	34.17	37.57	40.00
21	8.03	8.90	10.28	11.59	13.24	29.62	32.67	35.48	38.93	41.40
22	8.64	9.54	10.98	12.34	14.04	30.81	33.92	36.78	40.29	42.80
23	9.26	10.20	11.69	13.09	14.85	32.01	35.17	38.08	41.64	44.18
24	9.89	10.86	12.40	13.85	15.66	33.20	36.42	39.36	42.98	45.56
25	10.52	11.52	13.12	14.61	16.47	34.38	37.65	40.65	44.31	46.93
26	11.16	12.20	13.84	15.38	17.29	35.56	38.89	41.92	45.64	48.29
27	11.81	12.88	14.57	16.15	18.11	36.74	40.11	43.19	46.96	49.64
28	12.46	13.56	15.31	16.93	18.94	37.92	41.34	44.46	48.28	50.99
29	13.12	14.26	16.05	17.71	19.77	39.09	42.56	45.72	49.59	52.34
30	13.79	14.95	16.79	18.49	20.60	40.26	43.77	46.98	50.89	53.67
40	20.71	22.16	24.43	26.51	29.05	51.81	55.76	59.34	63.69	66.77
50	27.99	29.71	32.36	34.76	37.69	63.17	67.50	71.42	76.15	79.49
60	35.53	37.48	40.48	43.19	46.46	74.40	79.08	83.30	88.38	91.95
70	43.28	45.44	48.76	51.74	55.33	85.53	90.53	95.02	100.4	104.2
80	51.17	53.54	57.15	60.39	64.28	96.58	101.9	106.6	112.3	116.3
90	59.20	61.75	65.65	69.13	73.29	107.6	113.1	118.1	124.1	128.3
100	67.33	70.06	74.22	77.93	82.36	118.5	124.3	129.6	135.8	140.2

Πηγή: P. Newbold (1984), "Statistics for Business and Economics",
Prentice-Hall Inc.

ΠΙΝΑΚΑΣ 4: Κριτικές τιμές της κατανομής F για $\alpha=0,05$



$v_2 \backslash v_1$	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	40	60	120	∞
1	161	200	216	225	230	234	237	239	241	242	244	246	248	249	250	251	252	253	254
2	18.5	19.0	19.2	19.2	19.3	19.3	19.4	19.4	19.4	19.4	19.4	19.4	19.4	19.5	19.5	19.5	19.5	19.5	19.5
3	10.1	9.55	9.28	9.12	9.01	8.94	8.89	8.85	8.81	8.79	8.74	8.70	8.66	8.64	8.62	8.59	8.57	8.55	8.53
4	7.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04	6.00	5.96	5.91	5.86	5.80	5.77	5.75	5.72	5.69	5.66	5.63
5	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82	4.77	4.74	4.68	4.62	4.56	4.53	4.50	4.46	4.43	4.40	4.37
6	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15	4.10	4.06	4.00	3.94	3.87	3.84	3.81	3.77	3.74	3.70	3.67
7	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73	3.68	3.64	3.57	3.51	3.44	3.41	3.38	3.34	3.30	3.27	3.23
8	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44	3.39	3.35	3.28	3.22	3.15	3.12	3.08	3.04	3.01	2.97	2.93
9	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23	3.18	3.14	3.07	3.01	2.94	2.90	2.86	2.83	2.79	2.75	2.71
10	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07	3.02	2.98	2.91	2.85	2.77	2.74	2.70	2.66	2.62	2.58	2.54
11	4.84	3.98	3.59	3.36	3.20	3.09	3.01	2.95	2.90	2.85	2.79	2.72	2.65	2.61	2.57	2.53	2.49	2.45	2.40
12	4.75	3.89	3.49	3.26	3.11	3.00	2.91	2.85	2.80	2.75	2.69	2.62	2.54	2.51	2.47	2.43	2.38	2.34	2.30
13	4.67	3.81	3.41	3.18	3.03	2.92	2.83	2.77	2.71	2.67	2.60	2.53	2.46	2.42	2.38	2.34	2.30	2.25	2.21
14	4.60	3.74	3.34	3.11	2.96	2.85	2.76	2.70	2.65	2.60	2.53	2.46	2.39	2.35	2.31	2.27	2.22	2.18	2.13
15	4.54	3.68	3.29	3.06	2.90	2.79	2.71	2.64	2.59	2.54	2.48	2.40	2.33	2.29	2.25	2.20	2.16	2.11	2.07

16	4.49	3.63	3.24	3.01	2.85	2.74	2.66	2.59	2.54	2.49	2.42	2.35	2.28	2.24	2.19	2.15	2.11	2.06	2.01
17	4.45	3.59	3.20	2.96	2.81	2.70	2.61	2.55	2.49	2.45	2.38	2.31	2.23	2.19	2.15	2.10	2.06	2.01	1.96
18	4.41	3.55	3.16	2.93	2.77	2.66	2.58	2.51	2.46	2.41	2.34	2.27	2.19	2.15	2.11	2.06	2.02	1.97	1.92
19	4.38	3.52	3.13	2.90	2.74	2.63	2.54	2.48	2.42	2.38	2.31	2.23	2.16	2.11	2.07	2.03	1.98	1.93	1.88
20	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45	2.39	2.35	2.28	2.20	2.12	2.08	2.04	1.99	1.95	1.90	1.84
21	4.32	3.47	3.07	2.84	2.68	2.57	2.49	2.42	2.37	2.32	2.25	2.18	2.10	2.05	2.01	1.96	1.92	1.87	1.81
22	4.30	3.44	3.05	2.82	2.66	2.55	2.46	2.40	2.34	2.30	2.23	2.15	2.07	2.03	1.98	1.94	1.89	1.84	1.78
23	4.28	3.42	3.03	2.80	2.64	2.53	2.44	2.37	2.32	2.27	2.20	2.13	2.05	2.01	1.96	1.91	1.86	1.81	1.76
24	4.26	3.40	3.01	2.78	2.62	2.51	2.42	2.36	2.30	2.25	2.18	2.11	2.03	1.98	1.94	1.89	1.84	1.79	1.73
25	4.24	3.39	2.99	2.76	2.60	2.49	2.40	2.34	2.28	2.24	2.16	2.09	2.01	1.96	1.92	1.87	1.82	1.77	1.71
30	4.17	3.32	2.92	2.69	2.53	2.42	2.33	2.27	2.21	2.16	2.09	2.01	1.93	1.89	1.84	1.79	1.74	1.68	1.62
40	4.08	3.23	2.84	2.61	2.45	2.34	2.25	2.18	2.12	2.08	2.00	1.92	1.84	1.79	1.74	1.69	1.64	1.58	1.51
60	4.00	3.15	2.76	2.53	2.37	2.25	2.17	2.10	2.04	1.99	1.92	1.84	1.75	1.70	1.65	1.59	1.53	1.47	1.39
120	3.92	3.07	2.68	2.45	2.29	2.18	2.09	2.02	1.96	1.91	1.83	1.75	1.66	1.61	1.55	1.50	1.43	1.35	1.25
∞	3.84	3.00	2.60	2.37	2.21	2.10	2.01	1.94	1.88	1.83	1.75	1.67	1.57	1.52	1.46	1.39	1.32	1.22	1.00

Πηγή: Maddala G.S. (1992), "Introduction to econometrics, second Edition", McMillan.

ΠΙΝΑΚΑΣ 5: Κριτικές τιμές της κατανομής F για $\alpha=0,01$

$v_2 \backslash v_1$	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	40	60	120	∞
1	4.052	5.000	5.403	5.625	5.764	5.859	5.928	5.982	6.023	6.056	6.106	6.157	6.209	6.235	6.261	6.287	6.313	6.339	6.366
2	98.5	99.0	99.2	99.3	99.3	99.3	99.4	99.4	99.4	99.4	99.4	99.4	99.5	99.5	99.5	99.5	99.5	99.5	99.5
3	34.1	30.8	29.5	28.7	28.2	27.9	27.7	27.5	27.3	27.2	27.1	26.9	26.7	26.6	26.5	26.4	26.3	26.2	26.1
4	21.2	18.0	16.7	16.0	15.5	15.2	15.0	14.8	14.7	14.5	14.4	14.2	14.0	13.9	13.8	13.7	13.7	13.6	13.5
5	16.3	13.3	12.1	11.4	11.0	10.7	10.5	10.3	10.2	10.1	9.89	9.72	9.55	9.47	9.38	9.29	9.20	9.11	9.02
6	13.7	10.9	9.78	9.15	8.75	8.47	8.26	8.10	7.98	7.87	7.72	7.56	7.40	7.31	7.23	7.14	7.06	6.97	6.88
7	12.2	9.55	8.45	7.85	7.46	7.19	6.99	6.84	6.72	6.62	6.47	6.31	6.16	6.07	5.99	5.91	5.82	5.74	5.65
8	11.3	8.65	7.59	7.01	6.63	6.37	6.18	6.03	5.91	5.81	5.67	5.52	5.36	5.28	5.20	5.12	5.03	4.95	4.86
9	10.6	8.02	6.99	6.42	6.06	5.80	5.61	5.47	5.35	5.26	5.11	4.96	4.81	4.73	4.65	4.57	4.48	4.40	4.31
10	10.0	7.56	6.55	5.99	5.64	5.39	5.20	5.06	4.94	4.85	4.71	4.56	4.41	4.33	4.25	4.17	4.08	4.00	3.91
11	9.65	7.21	6.22	5.67	5.32	5.07	4.89	4.74	4.63	4.54	4.40	4.25	4.10	4.02	3.94	3.86	3.78	3.69	3.60
12	9.33	6.93	5.95	5.41	5.06	4.82	4.64	4.50	4.39	4.30	4.16	4.01	3.86	3.78	3.70	3.62	3.54	3.45	3.36
13	9.07	6.70	5.74	5.21	4.86	4.62	4.44	4.30	4.19	4.10	3.96	3.82	3.66	3.59	3.51	3.43	3.34	3.25	3.17
14	8.86	6.51	5.56	5.04	4.70	4.46	4.28	4.14	4.03	3.94	3.80	3.66	3.51	3.43	3.35	3.27	3.18	3.09	3.00
15	8.68	6.36	5.42	4.89	4.56	4.32	4.14	4.00	3.89	3.80	3.67	3.52	3.37	3.29	3.21	3.13	3.05	2.96	2.87

16	8.53	6.23	5.29	4.77	4.44	4.20	4.03	3.89	3.78	3.69	3.55	3.41	3.26	3.18	3.10	3.02	2.93	2.84	2.75
17	8.40	6.11	5.19	4.67	4.34	4.10	3.93	3.79	3.68	3.59	3.46	3.31	3.16	3.08	3.00	2.92	2.83	2.75	2.65
18	8.29	6.01	5.09	4.58	4.25	4.01	3.84	3.71	3.60	3.51	3.37	3.23	3.08	3.00	2.92	2.84	2.75	2.66	2.57
19	8.19	5.93	5.01	4.50	4.17	3.94	3.77	3.63	3.52	3.43	3.30	3.15	3.00	2.92	2.84	2.76	2.67	2.58	2.49
20	8.10	5.85	4.94	4.43	4.10	3.87	3.70	3.56	3.46	3.37	3.23	3.09	2.94	2.86	2.78	2.69	2.61	2.52	2.42
21	8.02	5.78	4.87	4.37	4.04	3.81	3.64	3.51	3.40	3.31	3.17	3.03	2.88	2.80	2.72	2.64	2.55	2.46	2.36
22	7.95	5.72	4.82	4.31	3.99	3.76	3.59	3.45	3.35	3.26	3.12	2.98	2.83	2.75	2.67	2.58	2.50	2.40	2.31
23	7.88	5.66	4.76	4.26	3.94	3.71	3.54	3.41	3.30	3.21	3.07	2.93	2.78	2.70	2.62	2.54	2.45	2.35	2.26
24	7.82	5.61	4.72	4.22	3.90	3.67	3.50	3.36	3.26	3.17	3.03	2.89	2.74	2.66	2.58	2.49	2.40	2.31	2.21
25	7.77	5.57	4.68	4.18	3.86	3.63	3.46	3.32	3.22	3.13	2.99	2.85	2.70	2.62	2.53	2.45	2.36	2.27	2.17
30	7.56	5.39	4.51	4.02	3.70	3.47	3.30	3.17	3.07	2.98	2.84	2.70	2.55	2.47	2.39	2.30	2.21	2.11	2.01
40	7.31	5.18	4.31	3.83	3.51	3.29	3.12	2.99	2.89	2.80	2.66	2.52	2.37	2.29	2.20	2.11	2.02	1.92	1.80
60	7.08	4.98	4.13	3.65	3.34	3.12	2.95	2.82	2.72	2.63	2.50	2.35	2.20	2.12	2.03	1.94	1.84	1.73	1.60
120	6.85	4.79	3.95	3.48	3.17	2.96	2.79	2.66	2.56	2.47	2.34	2.19	2.03	1.95	1.86	1.76	1.66	1.53	1.38
∞	6.63	4.61	3.78	3.32	3.02	2.80	2.64	2.51	2.41	2.32	2.18	2.04	1.88	1.79	1.70	1.59	1.47	1.32	1.00

Πηγή: Maddala G.S. (1992), "Introduction to econometrics, second Edition", McMillan.

ΠΙΝΑΚΑΣ 6: Κριτικές τιμές του συντελεστή συσχέτισης κατά τάξεις του Spearman

Μέγεθος δείγματος n	α 10%	α 5%	α 1%
4	1.000		
5	0.900	1.000	
6	0.829	0.886	1.000
7	0.714	0.786	0.929
8	0.643	0.738	0.885
9	0.600	0.683	0.845
10	0.564	0.648	0.808
11	0.520	0.620	0.773
12	0.496	0.591	0.741
13	0.475	0.566	0.716
14	0.456	0.544	0.695
15	0.438	0.524	0.676
16	0.425	0.506	0.657
17	0.411	0.490	0.639
18	0.399	0.475	0.622
19	0.388	0.462	0.606
20	0.377	0.450	0.591
21	0.368	0.438	0.576
22	0.359	0.428	0.562
23	0.351	0.418	0.549
24	0.343	0.409	0.537
25	0.336	0.400	0.526
26	0.329	0.392	0.515
27	0.323	0.384	0.505
28	0.317	0.377	0.496
29	0.311	0.370	0.487
30	0.305	0.364	0.478

Πηγή: Κουνιάς, Κολυβά-Μαχαίρα, Μπαγιάνης, Μπόρα-Σέντα (1985),
 "Εισαγωγή στην Στατιστική", Εκδ. Υπηρεσίας Δημοσιευμάτων ΑΠΘ,
 Θεσσαλονίκη.

ΠΙΝΑΚΑΣ 7: Ο Μετασχηματισμός $W = \frac{1}{2} \log_e \left(\frac{1+r_{XY}}{1-r_{XY}} \right)$ του συντελεστή συσχέτισης

<u>r</u>	<u>z</u>	<u>r</u>	<u>z</u>	<u>r</u>	<u>z</u>
0.00	0.0000	0.45	0.4847	0.90	1.4722
0.01	0.0100	0.46	0.4973	0.91	1.5275
0.02	0.0200	0.47	0.5101	0.92	1.5890
0.03	0.0300	0.48	0.5230	0.93	1.6584
0.04	0.0400	0.49	0.5361	0.94	1.7380
0.05	0.0500	0.50	0.5493	0.95	1.8318
0.06	0.0601	0.51	0.5627	0.96	1.9459
0.07	0.0701	0.52	0.5763	0.961	1.9588
0.08	0.0802	0.53	0.5901	0.962	1.9721
0.09	0.0902	0.54	0.6042	0.963	1.9857
0.10	0.1003	0.55	0.6184	0.964	1.9996
0.11	0.1104	0.56	0.6328	0.965	2.0139
0.12	0.1206	0.57	0.6475	0.966	2.0287
0.13	0.1307	0.58	0.6625	0.967	2.0439
0.14	0.1409	0.59	0.6777	0.968	2.0595
0.15	0.1511	0.60	0.6931	0.969	2.0756
0.16	0.1614	0.61	0.7089	0.970	2.0923
0.17	0.1717	0.62	0.7250	0.971	2.1095
0.18	0.1820	0.63	0.7414	0.972	2.1273
0.19	0.1923	0.64	0.7582	0.973	2.1457
0.20	0.2027	0.65	0.7753	0.974	2.1649
0.21	0.2132	0.66	0.7928	0.975	2.1847
0.22	0.2237	0.67	0.8107	0.976	2.2054
0.23	0.2342	0.68	0.8291	0.977	2.2269
0.24	0.2448	0.69	0.8480	0.978	2.2494
0.25	0.2554	0.70	0.8673	0.979	2.2729
0.26	0.2661	0.71	0.8872	0.980	2.2976
0.27	0.2769	0.72	0.9076	0.981	2.3235
0.28	0.2877	0.73	0.9287	0.982	2.3507
0.29	0.2986	0.74	0.9505	0.983	2.3796
0.30	0.3095	0.75	0.9730	0.984	2.4101
0.31	0.3205	0.76	0.9962	0.985	2.4427
0.32	0.3316	0.77	1.0203	0.986	2.4774
0.33	0.3428	0.78	1.0454	0.987	2.5147
0.34	0.3541	0.79	1.0714	0.988	2.5550
0.35	0.3654	0.80	1.0986	0.989	2.5987
0.36	0.3769	0.81	1.1270	0.990	2.6467
0.37	0.3884	0.82	1.1568	0.991	2.6996
0.38	0.4001	0.83	1.1881	0.992	2.7587
0.39	0.4118	0.84	1.2212	0.993	2.8257
0.40	0.4236	0.85	1.2562	0.994	2.9031
0.41	0.4356	0.86	1.2933	0.995	2.9945
0.42	0.4477	0.87	1.3331	0.996	3.1063
0.43	0.4599	0.88	1.3758	0.997	3.2504
0.44	0.4722	0.89	1.4219	0.998	3.4534

Πηγή: Owen B. (1962), "Handbook of Statistical tables", Addison-Wesley.

ΠΙΝΑΚΑΣ 8: Όρια του στατιστικού Durbin-Watson

n	k = 1		k = 2		k = 3		k = 4		k = 5	
	d _L	d _U	d _L	d _U	d _L	d _U	d _L	d _U	d _L	d _U
15	1.08	1.36	0.95	1.54	0.82	1.75	0.69	1.97	0.56	2.21
16	1.10	1.37	0.98	1.54	0.86	1.73	0.74	1.93	0.62	2.15
17	1.13	1.38	1.02	1.54	0.90	1.71	0.78	1.90	0.67	2.10
18	1.16	1.39	1.05	1.53	0.93	1.69	0.82	1.87	0.71	2.06
19	1.18	1.40	1.08	1.53	0.97	1.68	0.86	1.85	0.75	2.02
20	1.20	1.41	1.10	1.54	1.00	1.68	0.90	1.83	0.79	1.99
21	1.22	1.42	1.13	1.54	1.03	1.67	0.93	1.81	0.83	1.96
22	1.24	1.43	1.15	1.54	1.05	1.66	0.96	1.80	0.86	1.94
23	1.26	1.44	1.17	1.54	1.08	1.66	0.99	1.79	0.90	1.92
24	1.27	1.45	1.19	1.55	1.10	1.66	1.01	1.78	0.93	1.90
25	1.29	1.45	1.21	1.55	1.12	1.66	1.04	1.77	0.95	1.89
26	1.30	1.46	1.22	1.55	1.14	1.65	1.06	1.76	0.98	1.88
27	1.32	1.47	1.24	1.56	1.16	1.65	1.08	1.76	1.01	1.86
28	1.33	1.48	1.26	1.56	1.18	1.65	1.10	1.75	1.03	1.85
29	1.34	1.48	1.27	1.56	1.20	1.65	1.12	1.74	1.05	1.84
30	1.35	1.49	1.28	1.57	1.21	1.65	1.14	1.74	1.07	1.83
31	1.36	1.50	1.30	1.57	1.23	1.65	1.16	1.74	1.09	1.83
32	1.37	1.50	1.31	1.57	1.24	1.65	1.18	1.73	1.11	1.82
33	1.38	1.51	1.32	1.58	1.26	1.65	1.19	1.73	1.13	1.81
34	1.39	1.51	1.33	1.58	1.27	1.65	1.21	1.73	1.15	1.81
36	1.40	1.52	1.34	1.58	1.28	1.65	1.22	1.73	1.16	1.80
36	1.41	1.52	1.35	1.59	1.29	1.65	1.24	1.73	1.18	1.80
37	1.42	1.53	1.36	1.59	1.31	1.66	1.25	1.72	1.19	1.80
38	1.43	1.54	1.37	1.59	1.32	1.66	1.26	1.72	1.21	1.79
39	1.43	1.54	1.38	1.60	1.33	1.66	1.27	1.72	1.22	1.79
40	1.44	1.54	1.39	1.60	1.34	1.66	1.29	1.72	1.23	1.79
45	1.48	1.57	1.43	1.62	1.38	1.67	1.34	1.72	1.29	1.78
50	1.60	1.59	1.46	1.63	1.42	1.67	1.38	1.72	1.34	1.77
55	1.53	1.60	1.49	1.64	1.45	1.68	1.41	1.72	1.38	1.77
60	1.55	1.62	1.51	1.65	1.48	1.69	1.44	1.73	1.41	1.77
65	1.57	1.63	1.54	1.66	1.50	1.70	1.47	1.73	1.44	1.77
70	1.58	1.64	1.55	1.67	1.52	1.70	1.49	1.74	1.46	1.77
75	1.60	1.65	1.57	1.68	1.54	1.71	1.51	1.74	1.49	1.77
80	1.61	1.66	1.59	1.69	1.56	1.72	1.53	1.74	1.51	1.77
85	1.62	1.67	1.60	1.70	1.57	1.72	1.55	1.75	1.52	1.77
90	1.63	1.68	1.61	1.70	1.59	1.73	1.57	1.75	1.54	1.78
95	1.64	1.69	1.62	1.71	1.60	1.73	1.58	1.75	1.56	1.78
100	1.65	1.69	1.63	1.72	1.61	1.74	1.59	1.76	1.57	1.78

Πηγή: Sen-Srivastava (1990), "Regression Analysis", Springer-Verlag.