

Μαρία Ι. Διαμαντοπούλου

# Δασική Στατιστική

Θεωρία και Δασοβιομετρικές Εφαρμογές  
με χρήση IBM-SPSS και R

Για επικοινωνία με τη συγγραφέα:

<https://mdiamantopoulou.gr/>

<https://www.for.auth.gr/προσωπικο/μελη/δεπ/διαμαντοπουλου-μαρια>

ISBN 978-960-456-621-1

© Copyright Σεπτέμβριος 2024, Εκδόσεις ΖΗΤΗ, Μαρία Ι. Διαμαντοπούλου

---

*Το παρόν έργο πνευματικής ιδιοκτησίας προστατεύεται κατά τις διατάξεις του ελληνικού νόμου (Ν.2121/1993 όπως έχει τροποποιηθεί και ισχύει σήμερα) και τις διεθνείς συμβάσεις περί πνευματικής ιδιοκτησίας. Απαγορεύεται απολύτως η άνευ γραπτής άδειας του εκδότη κατά οποιοδήποτε τρόπο ή μέσο αντιγραφή, φωτοανατύπωση και εν γένει αναπαραγωγή, εκμίσθωση ή δανεισμός, μετάφραση, διασκευή, αναμετάδοση στο κοινό σε οποιαδήποτε μορφή (ηλεκτρονική, μηχανική ή άλλη) και η εν γένει εκμετάλλευση του συνόλου ή μέρους του έργου.*

---

<b>Φωτοστοιχειοθεσία</b>	<b>Π. ΖΗΤΗ &amp; Σια ΙΚΕ</b>
<b>Εκτύπωση</b>	18ο χλμ Θεσ/νίκης-Περαίας
<b>Βιβλιοδεσία</b>	Τ.Θ. 4171 • Περαία Θεσσαλονίκης • Τ.Κ. 570 19
	Τηλ.: 2392.072.222 - Fax: 2392.072.229 • e-mail: info@ziti.gr



**ΒΙΒΛΙΟΠΩΛΕΙΟ ΘΕΣΣΑΛΟΝΙΚΗΣ - ΚΕΝΤΡΙΚΗ ΔΙΑΘΕΣΗ:**  
Αρμενοπούλου 27, 546 35 Θεσσαλονίκη  
Τηλ.: 2310.203.720, Fax: 2310.211.305 • e-mail: sales@ziti.gr

**ΒΙΒΛΙΟΠΩΛΕΙΟ ΑΘΗΝΩΝ - ΠΩΛΗΣΗ ΛΙΑΝΙΚΗ-ΧΟΝΔΡΙΚΗ:**  
Χαριλάου Τρικούπη 22, 106 79 Αθήνα  
Τηλ.-Fax: 210.3816.650 • e-mail: athina@ziti.gr

**ΗΛΕΚΤΡΟΝΙΚΟ ΒΙΒΛΙΟΠΩΛΕΙΟ:** [www.ziti.gr](http://www.ziti.gr)

## Πρόλογος

Ο σκοπός συγγραφής αυτού του βιβλίου ήταν η κάλυψη των αναγκών της διδασκαλίας του μαθήματος της «Εφαρμοσμένης Στατιστικής», του Τμήματος Δασολογίας και Φυσικού Περιβάλλοντος, του Α.Π.Θ. Η κάλυψη των αναγκών αυτών είχε διττό ρόλο. Αφενός την θεωρητική ανάπτυξη και κατανόηση βασικών μεθοδολογιών της εφαρμοσμένης στατιστικής σε βιολογικά και μάλιστα δασικά δεδομένα, αφετέρου την υλοποίηση των μεθοδολογιών αυτών, με χρήση των αυτοματοποιημένων ηλεκτρονικών διαδικασιών, οι οποίες προσφέρονται σήμερα μέσω στατιστικών προγραμμάτων και γλωσσών προγραμματισμού. Τέλος, δόθηκε ιδιαίτερη βαρύτητα στην κατανόηση των μεθοδολογιών οι οποίες αναπτύσσονται, κάνοντας χρήση του απαραίτητου βασικού μαθηματικού υποβάθρου αλλά και στην πρακτική κατανόηση της τελικής στατιστικής ερμηνείας και συμπερασματολογίας.

Κατ' αυτό τον τρόπο, πιστεύεται ότι το βιβλίο θα συμβάλλει έτσι ώστε η Δασική Στατιστική να αποτελέσει βάση και κατανοητό αρωγό στις ερευνητικές προσπάθειες αλλά και στην καθημερινή πρακτική των φοιτητριών/-ων του Τμήματος Δασολογίας και Φυσικού Περιβάλλοντος του Α.Π.Θ., αλλά και άλλων σχετικών Τμημάτων. Επιπλέον, φιλοδοξεί να αποτελέσει ένα κατανοητό βοήθημα για οποιονδήποτε επιστήμονα αντιμετωπίζει ζητήματα τα οποία άπτονται του πεδίου της Εφαρμοσμένης Στατιστικής.

Η ανάπτυξη των βιομετρικών μεθόδων από την εποχή του K. Pearson και του R. Fisher στις αρχές του 20<sup>ου</sup> αιώνα, συνετέλεσε στη διαμόρφωση και ανάπτυξη της Δασικής Βιομετρίας ως τη βασική επιστήμη στη δασική έρευνα και πράξη, η οποία συνδυάζει μετρήσεις πεδίου με προηγμένες στατιστικές και υπολογιστικές τεχνικές, με σκοπό να παρέχει πληροφορίες για τη δυναμική των δασών και να καθοδηγεί τις διαδικασίες λήψης αποφάσεων. Ως αναπόσπαστο και βασικό μέρος της Δασικής Βιομετρίας, η Δασική Στατιστική δίνει έμφαση στην προσαρμογή των μεθόδων της εφαρμοσμένης στατιστικής, στη δασική έρευνα και πράξη. Για το λόγο αυτό, αποφασίστηκε η συγκεκριμένη επιλογή του τίτλου του βιβλίου ως: «*Δασική Στατιστική – Θεωρία και Δασοβιομετρικές εφαρμογές με χρήση IBM-SPSS και R*».

Η ύλη του βιβλίου βασίζεται στο ερευνητικό και συγγραφικό έργο της συγγραφέα και ενσωματώνει τη σύγχρονη επιστημονική γνώση σε αντικείμενα σχετικά με την Εφαρμοσμένη Δασική Στατιστική. Δόθηκε έμφαση, ώστε το βιβλίο να έχει έναν εφαρμοσμένο προσανατολισμό, αποφεύγοντας τις θεωρητικές αποδείξεις, παραπέμποντας τον ενδιαφερόμενο σε εξειδικευμένα βιβλία και επιστημονικά άρθρα, εφόσον επιθυμεί βαθύτερη μαθηματική αναζήτηση.

Επίσης, δόθηκε έμφαση, ώστε σε όλο το βιβλίο, η παρουσίαση της θεωρίας να συνοδεύεται από πληθώρα παραδειγμάτων, κυρίως δασοβιομετρικών μεταβλητών, για να είναι ευκολότερη η κατανόησή της. Όπου επίσης θεωρείται απαραίτητο, η επίλυση των παραδειγμάτων γίνεται με χρήση και αναλυτική περιγραφή της διαδικασίας, τόσο του στατιστικού πακέτου IBM-SPSS, όσο και της γλώσσας προγραμματισμού R.

Το βιβλίο αποτελείται από ένδεκα Κεφάλαια. Σε όλα τα κεφάλαια, δίνεται η επίλυση δασοβιομετρικών κυρίως παραδειγμάτων με χρήση και αναλυτική περιγραφή της διαδικασίας, τόσο του πακέτου IBM-SPSS, όσο και της γλώσσας προγραμματισμού R.

**Στο πρώτο κεφάλαιο** της Εισαγωγής περιγράφεται η έννοια και η διαχρονική διαδρομή της Στατιστικής, της Δασικής Βιομετρίας ως κλάδο της εφαρμοσμένης Στατιστικής, δίνονται εισαγωγικοί στατιστικοί ορισμοί και περιγράφονται οι εισαγωγικές βασικές έννοιες χρήσης της στατιστικής γλώσσας προγραμματισμού R και του στατιστικού πακέτου IBM® SPSS® Statistics (SPSS).

**Στο δεύτερο κεφάλαιο** περιγράφονται οι μέθοδοι και τα μέσα συλλογής στατιστικών στοιχείων, η προ-επεξεργασία των δεδομένων, οι διαδικασίες ελέγχου των πρωτογενών στοιχείων μέσω της διερευνητικής ανάλυσης δεδομένων, οι μέθοδοι παρουσίασης στατιστικών δεδομένων και αποτελεσμάτων στατιστικής ανάλυσης.

**Στο τρίτο κεφάλαιο** περιγράφονται εμπειρικές κατανομές συχνοτήτων και συγκεκριμένα η εμπειρική κατανομή συχνοτήτων ασυνεχούς μεταβλητής και η εμπειρική κατανομή συχνοτήτων συνεχούς μεταβλητής.

**Στο τέταρτο κεφάλαιο** περιγράφονται τα μέτρα κεντρικής τάσης, τα οποία περιλαμβάνουν τον αριθμητικό μέσο όρο και τις ιδιότητές του, τον σταθμισμένο αριθμητικό μέσο όρο, τον αποκομμένο αριθμητικό μέσο όρο, τον γεωμετρικό μέσο, τον αρμονικό μέσο, τον τετραγωνικό μέσο, τη σχέση μεταξύ αριθμητικού, γεωμετρικού, αρμονικού και τετραγωνικού μέσου, τους M-εκτιμητές, τη διάμεσο, τον υπολογισμό της διαμέσου σε μη ομαδοποιημένα στοιχεία και σε ομαδοποιημένα στοιχεία, την επικρατούσα τιμή, τον υπολογισμό της επικρατούσας τιμής σε μη ομαδοποιημένα στοιχεία και σε ομαδοποιημένα στοιχεία, τη σχέση μεταξύ αριθμητικού μέσου, διαμέσου και επικρατούσας τιμής, τα ποσοστημόρια και τα τεταρτημόρια.

**Στο πέμπτο κεφάλαιο** περιγράφονται τα μέτρα διασποράς, τα οποία περιλαμβάνουν το εύρος, το ενδοτεταρτημοριακό εύρος, τη μέση απόλυτη απόκλιση, τη διακύμανση, τις βασικές ιδιότητες της διακύμανσης, την τυπική απόκλιση, τις ιδιότητες της τυπικής απόκλισης, το συντελεστή μεταβλητότητας, το συντελεστή κύμανσης τεταρτημορίου και την καμπύλη του Lorenz.

**Στο έκτο κεφάλαιο** περιγράφονται τα μέτρα ασυμμετρίας και κύρτωσης, τα οποία περιλαμβάνουν τους συντελεστές ασυμμετρίας και κύρτωσης.

**Στο έβδομο κεφάλαιο** περιγράφονται τα στοιχεία πιθανοτήτων και οι θεωρητικές κατανομές πιθανοτήτων. Περιλαμβάνονται οι βασικοί ορισμοί και έννοιες της θεωρίας των πιθανοτήτων και οι αντίστοιχοι συμβολισμοί, η οπτική αναπαράσταση των συνόλων και των πράξεων αυτών, η έννοια και ο ορισμός της πιθανότητας, η στοχαστική ή στατιστική ανεξαρτησία, η δεσμευμένη πιθανότητα, το θεώρημα του Bayes, τα βασικά στοιχεία συνδυαστικής ανάλυσης, οι βασικές μέθοδοι προσδιορισμού του δειγματοχώρου σε σύνθετα πειράματα, η τυχαία ή στοχαστική μεταβλητή, η διωνυμική κατανομή, η κατανομή Poisson, η προσαρμογή εμπειρικής κατανομής δεδομένων δείγματος στις θεωρητικές συχνότητες της διωνυμικής κατανομής και της κατανομής Poisson, η κανονική και τυπική κανονική κατανομή, η προσαρμογή εμπειρικής κατανομής δεδομένων δείγματος σε κανονική κατανομή, η  $X^2$ -κατανομή, η  $t$ -κατανομή και η  $F$ -κατανομή.

**Στο όγδοο κεφάλαιο** περιγράφονται στοιχεία επαγωγικής στατιστικής και ο έλεγχος υποθέσεων. Περιλαμβάνονται οι ιδιότητες εκτιμητή, το διάστημα εμπιστοσύνης εκτιμητή, το διάστημα εμπιστοσύνης αριθμητικού μέσου όρου, το διάστημα εμπιστοσύνης της διαφοράς δύο αριθμητικών μέσων όρων, ο έλεγχος υποθέσεων αριθμητικού μέσου όρου, ο έλεγχος υποθέσεων περισσότερων των δύο πληθυσμών (ANOVA), ο έλεγχος υπόθεσης κατά έναν παράγοντα (one-way ANOVA), ο έλεγχος υπόθεσης κατά δύο παράγοντες (two-way ANOVA) και ο παραγοντικός έλεγχος υπόθεσης (factorial ANOVA).

**Στο ένατο κεφάλαιο** περιγράφονται η συσχέτιση και η παλινδρόμηση. Περιλαμβάνει την συνδιακύμανση, την συσχέτιση, την μερική συσχέτιση, την παλινδρόμηση, την απλή γραμμική παλινδρόμηση, τον έλεγχο υπόθεσης των συντελεστών παλινδρόμησης, τον έλεγχο ανάλυσης διακύμανσης για τη διαπίστωση της γραμμικότητας της εξίσωσης, τα μέτρα αξιολόγησης της προσαρμοσμένης εξίσωσης, την πολλαπλή γραμμική παλινδρόμηση, τη μη-γραμμική παλινδρόμηση, την παλινδρόμηση Ridge και την παλινδρόμηση Lasso.

**Στο δέκατο κεφάλαιο** περιγράφεται η διαχείριση ποιοτικών δεδομένων και περιλαμβάνονται τα κατάλληλα περιγραφικά στατιστικά για κατηγορικές μεταβλητές, ο έλεγχος καλής προσαρμογής (goodness of fit test), ο έλεγχος ανεξαρτησίας (chi-square test for independence/of association), ο έλεγχος ομοιογένειας

(test of homogeneity), η διωνυμική λογιστική παλινδρόμηση (binomial logistic regression), η ανάλυση συστάδων (cluster analysis), η τεχνική  $K$ -modes ( $K$ -modes), η ιεραρχική ταξινόμηση (hierarchical clustering), η παραγοντική ανάλυση μικτών μεταβλητών (Factorial Analysis of Mixed Data, FAMD), η εγκυρότητα (validity) και η αξιοπιστία (reliability) ερωτηματολογίου.

**Στο ενδέκατο κεφάλαιο** περιγράφονται οι μη-παραμετρικοί έλεγχοι και συγκεκριμένα, η δοκιμασία Wilcoxon (Wilcoxon Signed-Rank Test), η δοκιμασία Mann-Whitney  $U$  (Wilcoxon Rank-Sum Test), η δοκιμασία Kolmogorov-Smirnov ( $K-S$ ) για ένα δείγμα, η δοκιμασία Kruskal-Wallis  $H$  (Kruskal-Wallis  $H$  Test) και η δοκιμασία Friedman (Friedman Test).

Για τη συγγραφή αυτού του βιβλίου, η συγγραφέας άντλησε γνώση, εκτός των άλλων βιβλιογραφικών πηγών, και από τα διδακτικά συγγράμματα τα οποία και η ίδια διδάχθηκε, σε αντικείμενα Δασικής Βιομετρίας, του συναδέλφου Ομότιμου Καθηγητή του Εργαστηρίου Δασικής Βιομετρίας, του Τμήματος Δασολογίας και Φυσικού Περιβάλλοντος, ΑΠΘ, κ. Μάτη Κωνσταντίνου, τον οποίο ως καθηγήτη μου, θερμά ευχαριστώ.

Επίσης, η συγγραφέας άντλησε γνώση και εμπειρία από την πολυετή διδασκαλία που κλήθηκε να δώσει σε αντικείμενα Δασικής Βιομετρίας και Εφαρμοσμένης Στατιστικής, σε διάφορα Τμήματα της Τριτοβάθμιας εκπαίδευσης, συνειδητοποιώντας κατ' αυτό τον τρόπο τις υφιστάμενες ανάγκες σε παροχή έκτασης και βάθους γνώσης. Επίσης, γνώση αντλήθηκε από τη διαθέσιμη εκτενή βιβλιογραφία καθώς και από τις επιστημονικές δημοσιεύσεις και διδακτικές σημειώσεις που κατά καιρούς συνέγραψε η συγγραφέας, οι οποίες και αναφέρονται στην Βιβλιογραφία του παρόντος βιβλίου.

Στην ολοκλήρωση του βιβλίου αυτού θεωρώ χρέος και χαρά μου να ευχαριστήσω θερμά την οικογένειά μου η οποία μου συμπαραστάθηκε με κατανόηση, υπομονή και εμπύχωση, προκειμένου να ολοκληρωθεί η προσπάθεια.

Το βιβλίο αυτό αποτελεί πρώτη έκδοση. Η προσπάθεια για τη συγγραφή του είναι πιθανό να περιέχει λάθη ή ελλείψεις, οι οποίες προφανώς δεν κατέστη δυνατό να διαπιστωθούν στην παρούσα φάση. Για το λόγο αυτό η παροχή σχετικών υποδείξεων είναι ιδιαιτέρως ευπρόσδεκτη.

Θεσσαλονίκη, 2024

Μαρία Ι. Διαμαντοπούλου,  
Επίκουρη Καθηγήτρια Α.Π.Θ.

## Περιεχόμενα

### Κεφάλαιο 1

#### Εισαγωγή 1

1.1. Η Στατιστική ως έννοια και επιστήμη διαχρονικά .....	1
1.2. Η Δασική Βιομετρία ως κλάδος της εφαρμοσμένης Στατιστικής σήμερα .....	3
1.3. Εισαγωγή στη στατιστική γλώσσα προγραμματισμού R .....	5
1.4. Εισαγωγή στο στατιστικό πακέτο IBM-SPSS Statistics (SPSS) .....	13
1.5. Εισαγωγικοί στατιστικοί ορισμοί .....	16
Παράδειγμα 1.1 .....	16
Παράδειγμα 1.2 .....	17

### Κεφάλαιο 2

#### Στατιστικά στοιχεία – Προ-επεξεργασία (pre-processing) 21

2.1. Μέθοδοι και μέσα συλλογής στατιστικών στοιχείων .....	21
2.2. Προ-επεξεργασία (pre-processing process) δεδομένων .....	25
2.2.1. Διερευνητική Ανάλυση δεδομένων (Exploratory Data Analysis, EDA) .....	25
2.3. Μέθοδοι παρουσίασης στατιστικών δεδομένων και αποτελεσμάτων στατιστικής ανάλυσης .....	32
2.3.1. Πίνακες .....	32
2.3.2. Γραφήματα .....	33
Παράδειγμα 2.1 .....	35
Παράδειγμα 2.2 .....	38

### Κεφάλαιο 3

#### Εμπειρικές κατανομές συχνοτήτων 41

3.1. Εμπειρική κατανομή συχνοτήτων ασυνεχούς μεταβλητής .....	41
Παράδειγμα 3.1 .....	42

3.2. Εμπειρική κατανομή συχνοτήτων συνεχούς μεταβλητής .....	46
Παράδειγμα 3.2 .....	48

## Κεφάλαιο 4

<b>Μέτρα κεντρικής τάσης (location measures-central tendency measures)</b> .....	<b>51</b>
4.1. Μέσοι όροι .....	51
4.1.1. Αριθμητικός μέσος όρος ή μέση τιμή (arithmetic mean or average) .....	51
Παράδειγμα 4.1 .....	52
Παράδειγμα 4.2 .....	54
4.1.2. Ιδιότητες αριθμητικού μέσου όρου .....	55
Παράδειγμα 4.3 .....	56
Παράδειγμα 4.4 .....	56
Παράδειγμα 4.5 .....	57
Παράδειγμα 4.6 .....	58
4.1.3. Σταθμισμένος αριθμητικός μέσος όρος (weighted arithmetic mean) .....	59
Παράδειγμα 4.7 .....	59
4.1.4. Αποκομμένος μέσος όρος (trimmed mean) .....	61
Παράδειγμα 4.8 .....	62
4.1.5. Γεωμετρικός μέσος (geometric mean) .....	64
Παράδειγμα 4.9 .....	65
4.1.6. Αρμονικός μέσος (harmonic mean) .....	68
Παράδειγμα 4.10 .....	68
4.1.7. Τετραγωνικός μέσος (quadratic mean) .....	71
Παράδειγμα 4.11 .....	71
4.1.8. Σχέση μεταξύ αριθμητικού, γεωμετρικού, αρμονικού και τετραγωνικού μέσου .....	73
Παράδειγμα 4.12 .....	73
4.1.9. M-εκτιμητές (M-estimators) .....	74
4.2. Διάμεσος (median) .....	74
4.2.1. Υπολογισμός διαμέσου σε μη ομαδοποιημένα στοιχεία .....	75
Παράδειγμα 4.13 .....	75
Παράδειγμα 4.14 .....	77
4.2.2. Υπολογισμός διαμέσου σε ομαδοποιημένα στοιχεία .....	78
Παράδειγμα 4.15 .....	78
4.3. Τύπος ή επικρατούσα τιμή (mode) .....	80
4.3.1. Υπολογισμός της επικρατούσας τιμής σε μη ομαδοποιημένα στοιχεία .....	80
Παράδειγμα 4.16 .....	80



4.3.2. Υπολογισμός επικρατούσας τιμής σε ομαδοποιημένα στοιχεία .....	82
Παράδειγμα 4.17 .....	83
4.3.3 Σχέση μεταξύ αριθμητικού μέσου – διαμέσου - επικρατούσας τιμής .....	83
4.4. Ποσοστημόρια (quantiles) .....	84
4.4.1. Τεταρτημόρια (quartiles) .....	85
Παράδειγμα 4.18 .....	86
Παράδειγμα 4.19 .....	88

## Κεφάλαιο 5

### Μέτρα διασποράς (measures of dispersion) 91

5.1. Εύρος (range) .....	91
Παράδειγμα 5.1 .....	92
5.2. Ενδοτεταρτημοριακό εύρος (inter-quartile range) .....	94
5.3. Μέση απόλυτη απόκλιση (mean absolute deviation) .....	95
Παράδειγμα 5.2 .....	96
Παράδειγμα 5.3 .....	97
5.4. Διασπορά ή διακύμανση (variance) .....	98
Παράδειγμα 5.4 .....	100
5.4.1. Βασικές ιδιότητες της διακύμανσης .....	102
Παράδειγμα 5.5 .....	103
5.5. Τυπική απόκλιση (standard deviation) .....	104
Παράδειγμα 5.6 .....	105
Παράδειγμα 5.7 .....	106
5.5.1 Ιδιότητες τυπικής απόκλισης .....	107
5.6. Συντελεστής μεταβλητότητας ή κύμανσης (coefficient of variation) .....	109
Παράδειγμα 5.8 .....	110
5.7. Συντελεστής κύμανσης τεταρτημορίου (coefficient of quartile variation) .....	111
Παράδειγμα 5.9 .....	112
5.8. Καμπύλη του Lorenz (Lorenz curve) .....	113
Παράδειγμα 5.10 .....	114

## Κεφάλαιο 6

### Μέτρα ασυμμετρίας και κύρτωσης (measures of skewness and kurtosis) 117

6.1. Ασυμμετρία (skewness) .....	117
Παράδειγμα 6.1 .....	120

6.2. Κύρτωση (Kurtosis) .....	124
Παράδειγμα 6.2 .....	126
6.3. Υπολογισμός ασυμμετρίας και κύρτωσης με τη γλώσσα προγραμματισμού R .....	129
Παράδειγμα 6.3 .....	130

## Κεφάλαιο 7

### Στοιχεία πιθανοτήτων - Θεωρητικές κατανομές πιθανοτήτων (probability distributions)

133

7.1. Βασικοί ορισμοί και έννοιες της θεωρίας των πιθανοτήτων – βασικοί συμβολισμοί .....	134
Παράδειγμα 7.1 .....	135
Παράδειγμα 7.2 .....	135
Παράδειγμα 7.3 .....	136
7.2. Οπτική αναπαράσταση των συνόλων και πράξεων μεταξύ τους .....	137
Παράδειγμα 7.4 .....	137
7.3. Έννοια και ορισμός της πιθανότητας .....	138
Παράδειγμα 7.5 .....	139
Παράδειγμα 7.6 .....	142
Παράδειγμα 7.7 .....	143
Παράδειγμα 7.8 .....	143
7.4. Στοχαστική ή στατιστική ανεξαρτησία γεγονότων (statistically independent events) .....	144
Παράδειγμα 7.9 .....	145
7.5. Πιθανότητα υπό συνθήκη (δεσμευμένη πιθανότητα) (conditional probability) .....	146
Παράδειγμα 7.10 .....	146
Παράδειγμα 7.11 .....	147
7.6. Θεώρημα του Bayes (Bayes' theorem or Bayes' law or Bayes' rule) .....	149
Παράδειγμα 7.12 .....	150
7.7. Βασικά στοιχεία συνδυαστικής ανάλυσης (Combinatorial analysis) .....	151
Παράδειγμα 7.13 .....	152
Παράδειγμα 7.14 .....	152
Παράδειγμα 7.15 .....	153
Παράδειγμα 7.16 .....	154
7.8. Βασικές μέθοδοι προσδιορισμού του δειγματοχώρου σε σύνθετα πειράματα .....	155

Παράδειγμα 7.17 .....	156
Παράδειγμα 7.18 .....	157
7.9. Τυχαία ή στοχαστική μεταβλητή (random or stochastic variable) .....	157
7.10. Διωνυμική κατανομή (binomial distribution) .....	159
Παράδειγμα 7.19 .....	161
Παράδειγμα 7.20 .....	163
Παράδειγμα 7.21 .....	164
7.11. Κατανομή Poisson (Poisson distribution) .....	167
Παράδειγμα 7.22 .....	169
7.12. Προσαρμογή εμπειρικής κατανομής δεδομένων δείγματος, στις θεωρητικές συχνότητες της διωνυμικής κατανομής και της κατανομής Poisson .....	173
Παράδειγμα 7.23 .....	174
Παράδειγμα 7.24 .....	176
7.13. Κανονική και τυπική κανονική κατανομή (normal and standard normal distribution) .....	178
Παράδειγμα 7.25 .....	181
Παράδειγμα 7.26 .....	183
Παράδειγμα 7.27 .....	184
7.14. Προσαρμογή εμπειρικής κατανομής δεδομένων δείγματος σε κανονική κατανομή .....	187
Παράδειγμα 7.28 .....	188
7.15. $\chi^2$ - κατανομή (chi-squared distribution) .....	192
Παράδειγμα 7.29 .....	193
7.16. t-κατανομή ή Student – κατανομή (t-distribution or Student- distribution) .....	194
Παράδειγμα 7.30 .....	196
7.17. F-κατανομή (F-distribution or Snedecor's F-distribution or Fisher-Snedecor distribution) .....	198
Παράδειγμα 7.31 .....	199

## Κεφάλαιο 8

<b>Στοιχεία επαγωγικής στατιστικής (statistical inference) - Έλεγχος υποθέσεων (hypothesis testing)</b> .....	<b>201</b>
---	------------

8.1. Ιδιότητες εκτιμητή .....	202
8.2. Διάστημα εμπιστοσύνης εκτιμητή .....	204

8.3. Διάστημα εμπιστοσύνης αριθμητικού μέσου όρου .....	204
Παράδειγμα 8.1 .....	207
Παράδειγμα 8.2 .....	211
8.4. Διάστημα εμπιστοσύνης της διαφοράς δύο αριθμητικών μέσων όρων ....	215
Παράδειγμα 8.3 .....	218
Παράδειγμα 8.4 .....	222
Παράδειγμα 8.5 .....	223
8.5. Έλεγχος υποθέσεων ενός πληθυσμού .....	226
8.6. Έλεγχος υποθέσεων αριθμητικού μέσου όρου .....	229
Παράδειγμα 8.6 .....	231
Παράδειγμα 8.7 .....	234
8.7. Έλεγχος υποθέσεων περισσότερων του ενός πληθυσμών .....	240
8.7.1. Έλεγχος υποθέσεων δύο πληθυσμών .....	240
Παράδειγμα 8.8 .....	242
8.7.2. Έλεγχος υποθέσεων περισσότερων των δύο πληθυσμών (ANOVA) .....	244
8.7.2.1. Έλεγχος υπόθεσης κατά έναν παράγοντα (one-way ANOVA) .....	245
Παράδειγμα 8.9 .....	248
8.7.2.2. Έλεγχος υπόθεσης κατά δύο παράγοντες (two-way ANOVA) .....	255
Παράδειγμα 8.10 .....	259
Παράδειγμα 8.11 .....	268
8.7.2.3. Παραγοντικός έλεγχος υπόθεσης (factorial ANOVA) .....	275
Παράδειγμα 8.12 .....	277

## Κεφάλαιο 9

### Συσχέτιση (correlation) – Παλινδρόμηση (regression) 279

9.1. Συνδιακύμανση (covariance) .....	280
Παράδειγμα 9.1 .....	281
Παράδειγμα 9.2 .....	284
9.2. Συσχέτιση (correlation) .....	285
Παράδειγμα 9.3 .....	287
9.2.1 Μερική συσχέτιση (partial correlation) .....	292
Παράδειγμα 9.4 .....	293
9.3. Παλινδρόμηση (regression) .....	298
9.3.1. Απλή γραμμική παλινδρόμηση (simple linear regression, SLR) .....	299
9.3.1.1. Έλεγχος υπόθεσης των συντελεστών παλινδρόμησης .....	306
9.3.1.2. Έλεγχος ανάλυσης διακύμανσης για τη διαπίστωση της γραμμικότητας της εξίσωσης .....	308

9.3.1.3. Μέτρα αξιολόγησης της προσαρμοσμένης εξίσωσης .....	311
Παράδειγμα 9.5 .....	313
9.3.2. Πολλαπλή γραμμική παλινδρόμηση (multiple linear regression, MLR) .....	324
Παράδειγμα 9.6 .....	328
9.3.3. Μη-γραμμική παλινδρόμηση (nonlinear regression, NLR) .....	335
Παράδειγμα 9.7 .....	337
9.3.4. Παλινδρόμηση Ridge (Ridge regression analysis, RRA) .....	342
Παράδειγμα 9.8 .....	344
9.3.5. Παλινδρόμηση Lasso (Lasso regression) .....	349

## Κεφάλαιο 10

<b>Διαχείριση ποιοτικών δεδομένων (categorical data)</b> .....	<b>353</b>
10.1. Περιγραφικά στατιστικά κατάλληλα για κατηγορικές μεταβλητές .....	355
Παράδειγμα 10.1 .....	358
10.2. Έλεγχος καλής προσαρμογής (goodness of fit test) .....	364
Παράδειγμα 10.2 .....	366
Παράδειγμα 10.3 .....	371
Παράδειγμα 10.4 .....	373
10.3. Έλεγχος ανεξαρτησίας (Chi-square test for independence/of association) .....	375
Παράδειγμα 10.5 .....	376
Παράδειγμα 10.6 .....	382
10.4. Έλεγχος ομοιογένειας (test of Homogeneity) .....	387
Παράδειγμα 10.7 .....	387
10.5. Διωνυμική λογιστική παλινδρόμηση (binomial logistic regression) .....	388
Παράδειγμα 10.8 .....	390
10.6. Ανάλυση συστάδων (cluster analysis) .....	396
10.6.1. Τεχνική K-modes (K-modes) .....	397
Παράδειγμα 10.9 .....	399
10.6.2. Ιεραρχική ταξινόμηση (hierarchical clustering) .....	400
Παράδειγμα 10.10 .....	401
Παράδειγμα 10.11 .....	404
10.7. Παραγοντική ανάλυση μικτών μεταβλητών (Factorial Analysis of Mixed Data, FAMD) .....	410
Παράδειγμα 10.12 .....	413
Παράδειγμα 10.13 .....	415
Παράδειγμα 10.14 .....	417
10.8. Εγκυρότητα (Validity) και αξιοπιστία (Reliability) ερωτηματολογίου .....	422

10.8.1. Εγκυρότητα (Validity) ερωτηματολογίου .....	423
Παράδειγμα 10.15 .....	424
10.8.2. Αξιοπιστία (Reliability) ερωτηματολογίου .....	425
Παράδειγμα 10.16 .....	427
Παράδειγμα 10.17 .....	428

## Κεφάλαιο 11

### **Μη παραμετρικοί έλεγχοι (non-parametric tests)** 435

11.1. Δοκιμασία Wilcoxon (Wilcoxon Signed-Rank Test) .....	436
Παράδειγμα 11.1 .....	437
11.2. Δοκιμασία Mann-Whitney U (Wilcoxon Rank-Sum Test) .....	444
Παράδειγμα 11.2 .....	446
11.3. Δοκιμασία Kolmogorov-Smirnov (K-S) για ένα δείγμα .....	454
Παράδειγμα 11.3 .....	455
11.4. Δοκιμασία Kruskal-Wallis H (Kruskal-Wallis H Test) .....	459
Παράδειγμα 11.4 .....	460
11.5. Δοκιμασία Friedman (Friedman test) .....	468
Παράδειγμα 11.5 .....	469

### ΒΙΒΛΙΟΓΡΑΦΙΑ

A. Ελληνική .....	475
B. Ξενόγλωσση .....	479
Γ. Διαδικτυακή .....	490

ΠΑΡΑΡΤΗΜΑ Ι .....	493
-------------------	-----

ΕΥΡΕΤΗΡΙΟ ΟΡΩΝ .....	521
----------------------	-----

«Τύχην νόμιζε» (να λαμβάνουμε υπ' όψη «το τυχαίο»)

ΔΕΛΦΙΚΑ ΠΑΡΑΓΓΕΛΜΑΤΑ

«Νουν ηγεμόνα ποιεί» (Το μυαλό να κυριαρχεί στις αποφάσεις)

ΣΟΛΩΝ ο ΑΘΗΝΑΙΟΣ

### 1.1 Η Στατιστική ως έννοια και επιστήμη διαχρονικά

Ο όρος στατιστική εικάζεται ότι προέρχεται είτε από την αρχαία ελληνική λέξη ίστημι και του εξ αυτού παραγώγου ρήματος στατίζω, που σημαίνει τοποθετώ, ταξινομώ, συμπεραίνω, είτε από την λατινική λέξη status, που σημαίνει πολιτεία (κράτος). Από αρχαιοτάτων χρόνων, ήταν γνωστή η στατιστική, ως τρόπος σκέψης και ενέργειας, χωρίς αυτό να γίνεται συνειδητά, αρκετές προ Χριστού χιλιετίες. Υπάρχουν δε αναφορές συλλογής στοιχείων που παραπέμπουν σε στατιστική από πολλούς φιλοσόφους και συγγραφείς της αρχαιότητας, όπως στον Κομφούκιο (551-479 π.Χ.), στον Ξενοφώντα (430-355 π.Χ.), Ηρόδοτο (484-425 π.Χ.), Αριστοτέλη (384-322 π.Χ.). Στη β' ραψωδία της Ιλιάδας του Ομήρου κατά τον 8<sup>ο</sup>-9<sup>ο</sup> π.Χ. αιώνα διασώζεται ο κατάλογος των πλοίων των Αχαιών που εκστράτευσαν κατά της Τροίας (Τρωικός πόλεμος, 13<sup>ο</sup>-12<sup>ο</sup> αιώνας π.Χ.). Επίσης, είναι γνωστό το γεγονός της απογραφής του Ιουδαϊκού πληθυσμού την εποχή της γέννησης του Χριστού.

Ως πατέρας της Στατιστικής, αναφέρεται ο Γερμανός Gottfried Anchenwall (1719-1772), ο οποίος εργαζόμενος ως οικονομολόγος και αργότερα ως διδάσκαλος, για πρώτη φορά εισήγαγε τη λέξη στατιστική (Statistik), προκειμένου να αναφερθεί και να περιγράψει αριθμητικά δεδομένα τα οποία συγκέντρωνε, αλλά ήταν και ο πρώτος που χρησιμοποίησε πίνακες και γραφήματα γι' αυτόν το σκο-

πό. Υπάρχουν θεωρίες αμφισβητήσεων επ' αυτού, οι οποίες θεωρούν ότι αντίστοιχη δουλειά η οποία έγινε από τον Άγγλο William Petty (1623-1687), δεν έτυχε της ίδιας αναγνώρισης, γι' αυτό και δεν αναφέρεται.

Το 1654, ο Γάλλος Blaise Pascal (1623-1662) έθεσε τις βάσεις για τη Συνδυαστική Ανάλυση και το Λογισμό των Πιθανοτήτων. Ο Ελβετός Jacob Bernoulli, (1655-1705) ήταν εκείνος ο οποίος με δημοσίευσμά του το οποίο κυκλοφόρησε μετά το θάνατό του από συγγενή του, περιέγραψε την εφαρμογή της θεωρίας πιθανοτήτων σε τυχερά παιχνίδια και πρότεινε το θεώρημα του ισχυρού νόμου των μεγάλων αριθμών.

Στη συνέχεια, άρχισε ο συσχετισμός κάποιων στοιχείων, όπως η σχέση μεταξύ ηλικίας και θανάτου, όπου παρατηρήθηκε η κανονικότητα της θνησιμότητας των διαφόρων ηλικιών. Πρώτος ο Γάλλος Abraham de Moivre, το 1773, έθεσε τη βάση στην ανακάλυψη της κανονικής κατανομής, που αποτελεί τη βάση του μεγαλύτερου μέρους της στατιστικής, η οποία ονομάζεται και κατανομή του Gauss, προς τιμή του Γερμανού Karl Friedrich Gauss (1776-1855), που την ίδια σχεδόν εποχή κατέληξε στα ίδια συμπεράσματα. Λίγο αργότερα, ο Βέλγος A. Quetelet (1796-1874), επηρεαζόμενος από τους Laplace και Fourier, προσπάθησε να ορίσει τον "μέσο άνθρωπο" με βάση σωματομετρικά μεγέθη, δηλ. ουσιαστικά τη μέση τιμή, η οποία σαν έννοια είχε χρησιμοποιηθεί και από τον Αρχιμήδη ως "κέντρο βάρους". Επίσης, ήταν ο πρώτος που αντιμετώπισε την κανονική κατανομή, ως το νόμο των σφαλμάτων και παράλληλα το 1853 οργάνωσε το πρώτο διεθνές στατιστικό συνέδριο, του οποίου η συμβολή τόσο στη στατιστική θεωρία όσο και στην εφαρμογή ήταν εξαιρετικά σημαντική.

Αργότερα, η ωριμότητα της σκέψης των ανθρώπων που ασχολήθηκαν με την επιστήμη της στατιστικής, έδωσε φοβερή ώθηση σ' αυτό τον κλάδο της επιστήμης. Ενδεικτικά, ο Άγγλος Francis Galton (1822-1911) ήταν αυτός που έθεσε τις βάσεις της εφαρμογής της στατιστικής επιστήμης στα έμβια όντα. Λίγο αργότερα, ο Άγγλος Karl (Carl) Pearson (1857-1936), επηρεασμένος από τις μελέτες του Francis Galton, ήταν αυτός ο οποίος έθεσε τις βάσεις των στατιστικών μεθόδων που εδράζονται στη μέθοδο των ελαχίστων τετραγώνων, όπως είναι η θεωρία της παλινδρόμησης. Τέλος, ο Βρετανός Sir Ronald Aylmer Fisher (1890-1962), έθεσε τις βάσεις της θεωρίας του σχεδιασμού πειραμάτων σε αγροτικές καλλιέργειες.

Όσον αφορά την κατάρτιση σε θέματα Εφαρμοσμένης Στατιστικής των φοιτητριών/φοιτητών του Τμήματος Δασολογίας και Φυσικού Περιβάλλοντος, σύμφωνα με όσα αναφέρει ο Μάτης (2003), στο Τμήμα Δασολογίας και Φυσικού Περι-



βάλλοντος του Α.Π.Θ. πριν το ακαδημαϊκό έτος 1971-72 δεν διδασκόταν το μάθημα της Στατιστικής ως ξεχωριστό μάθημα. Στοιχεία Εφαρμοσμένης Στατιστικής διδασκόταν εντός της ύλης του μαθήματος της Δεντρομετρίας. Στη συνέχεια, για τα επόμενα δύο χρόνια η Στατιστική διδασκόταν ως μάθημα επιλογής στο δεύτερο έτος σπουδών. Από το ακαδημαϊκό έτος 1974-75 συμπεριλήφθηκε στο πρόγραμμα σπουδών του Τμήματος ως υποχρεωτικό μάθημα, διδασκόμενο στο πρώτο εξάμηνο του δεύτερου έτους σπουδών, ως Δασική Βιομετρία Ι, μέχρι και την τελευταία αναθεώρηση προπτυχιακού προγράμματος σπουδών.

Από το πανεπιστημιακό έτος 2021-2022 για τις/τους φοιτήτριες/φοιτητές με έτος εισαγωγής το έτος 2020, ο τίτλος του μαθήματος τροποποιήθηκε και έγινε Εφαρμοσμένη Στατιστική και διδάσκεται ως υποχρεωτικό μάθημα στο 1ο εξάμηνο του πρώτου έτους σπουδών του Τμήματος.

## 1.2 Η Δασική Βιομετρία ως κλάδος της εφαρμοσμένης Στατιστικής σήμερα

Σε γενικές γραμμές μπορεί να λεχθεί ότι **εφαρμοσμένη στατιστική** είναι εκείνος ο κλάδος της Στατιστικής επιστήμης, που ασχολείται με εμπειρικά δεδομένα τα οποία μπορεί να προέρχονται είτε από μετρήσεις/παρατηρήσεις γεγονότων που συμβαίνουν στον περιβάλλοντα χώρο είτε από καταγραφή αποτελεσμάτων πειραμάτων, με σκοπό να εξάγει βασικά επιστημονικά συμπεράσματα, λαμβάνοντας υπόψη το στοιχείο αβεβαιότητας που μπορεί να περικλείουν αυτά τα δεδομένα. Η ιδιαιτερότητά της ως κλάδος της γενικότερης μαθηματικής επιστήμης, έγκειται στο γεγονός ότι ενώ η γλώσσα των αριθμών εμπεριέχει έναν χαρακτήρα απόλυτης ακρίβειας, η εφαρμοσμένη στατιστική περιλαμβάνει ανακρίβειες, οι οποίες οφείλονται στα ίδια αυτά τα μεγέθη που μελετά, υπό περιορισμούς και οριοθετήσεις μέσω στατιστικών και μαθηματικών νόμων. Αποτελεί ως επιστήμη, τον καλύτερο δυνατό τρόπο με τον οποίο μπορούν να αναλυθούν και να αντιμετωπιστούν όλες εκείνες οι καταστάσεις όπου τα «καθαρά» μαθηματικά, εξαιτίας του απόλυτου των αριθμητικών μεγεθών τους, δεν είναι ευέλικτα. Έτσι λοιπόν η εφαρμοσμένη στατιστική, ως απαραίτητο και μοναδικό εργαλείο, έχει πάρα πολλά πεδία εφαρμογής και ανάλογα με την γνωστική περιοχή στην οποία εφαρμόζεται, παίρνει και το αντίστοιχο χαρακτηριστικό όνομα καθώς και τις αντίστοιχες θεωρητικές και πρακτικές εξειδικεύσεις που οφείλονται στη φύση του εκάστοτε γνωστικού αντικειμένου.

Η **Δασική Στατιστική**, ως κλάδος της εφαρμοσμένης Στατιστικής είναι η επιστήμη της μέτρησης (ποσοτικοποίησης) των χαρακτηριστικών (μεταβλητών) του δασικού οικοσυστήματος. Ο Prodan (1968) έδωσε τον ορισμό της Δασικής Βιομετρίας, ως το πεδίο της δασικής επιστήμης το οποίο περιλαμβάνει τις μεθόδους μαθηματικής στατιστικής και βιομετρίας που είναι σημαντικές για το δασικό οικοσύστημα. Δηλαδή, η Δασική Βιομετρία περιλαμβάνει την ποσοτικοποίηση, μέσα σε κάποιο εύλογο επίπεδο ακρίβειας και συνέπειας, των βιολογικών και φυσικών χαρακτηριστικών ενός δασικού οικοσυστήματος, μέσω της συγκέντρωσης, της παρουσίασης, της επεξεργασίας, της ανάλυσης και ερμηνείας αυτών των δεδομένων, προκειμένου να βρεθούν λύσεις ή απαντήσεις σε αναζητήσεις της δασικής επιστήμης. Ο ορισμός ο οποίος δόθηκε από τον Μ. Prodan υπογραμμίζει ότι η δασική Βιομετρία ασχολείται θεμελιωδώς με την εφαρμογή ποσοτικών μεθόδων για την κατανόηση και την περιγραφή των δασικών οικοσυστημάτων. Επίσης, σύμφωνα με τον ορισμό που δίνεται από το Ερευνητικό Ινστιτούτο Δασικής Στατιστικής του Portland, USA, η Δασική Στατιστική/Βιομετρία είναι η επιστήμη της μέτρησης (-μετρία) των δασών (βιο-). Περιλαμβάνει τον ποσοτικό προσδιορισμό των βιολογικών και φυσικών χαρακτηριστικών των δέντρων και της σχετικής βλάστησης, των εντόμων, των ασθενειών, της άγριας ζωής, της τοπογραφίας, των εδαφών και του κλίματος, μεμονωμένα και συλλογικά. Τα χαρακτηριστικά αυτά περιλαμβάνουν όλα τα ποσοτικώς προσδιορίσιμα χαρακτηριστικά εντός του δάσους, τόσο χρονικά όσο και χωρικά.

Η Δασική Στατιστική γενικά αποτελεί απαραίτητο εργαλείο για να αναλυθεί και να κριθεί μια κατάσταση, καθώς και να προβλεφθεί η πορεία κάποιου φαινομένου στο μέλλον. Πολλές φορές, κάποιες καταστάσεις δεν οδηγούν σε αυτονόητα ή αυταπόδεικτα συμπεράσματα, πχ. ποιός είναι ο συνολικός ξυλώδης όγκος που μπορεί ένα δασικό οικοσύστημα να παράγει, έτσι ώστε να γίνει προγραμματισμός ανάλογων υλοτομιών ή ποιά είναι το ποσοστό φυτρωτικότητας κάποιων σπόρων, κλπ. Σ' αυτή την περίπτωση, αλλά και σε πάρα πολλές άλλες, η χρήση της Στατιστικής, μπορεί να κρίνει την κατάσταση και να οδηγήσει ασφαλώς σε συμπεράσματα και αποφάσεις. Επίσης, τις περισσότερες φορές καλούμαστε να περιγράψουμε και να αναλύσουμε πληθυσμούς των οποίων τα μέλη είναι τόσα πολλά πχ. όλα τα δέντρα που απαρτίζουν ένα δασικό οικοσύστημα, ώστε να είναι αδύνατο να πάρουμε τις απαραίτητες πληροφορίες απ' όλα. Και σ' αυτές τις περιπτώσεις, η εφαρμογή των μεθόδων της Στατιστικής, δίνει τη δυνατότητα λύσεων. Κατ' αυτό τον τρόπο, η Δασική Στατιστική, εφαρμόζει τις αρχές και μεθόδους της εφαρμοσμένης στατιστικής χρησιμοποιώντας στοιχεία που αφορούν το δασικό οικοσύστημα. Περιλαμβάνει δε τρία βασικά μεθοδολογικά κεφάλαια. Το πρώ-

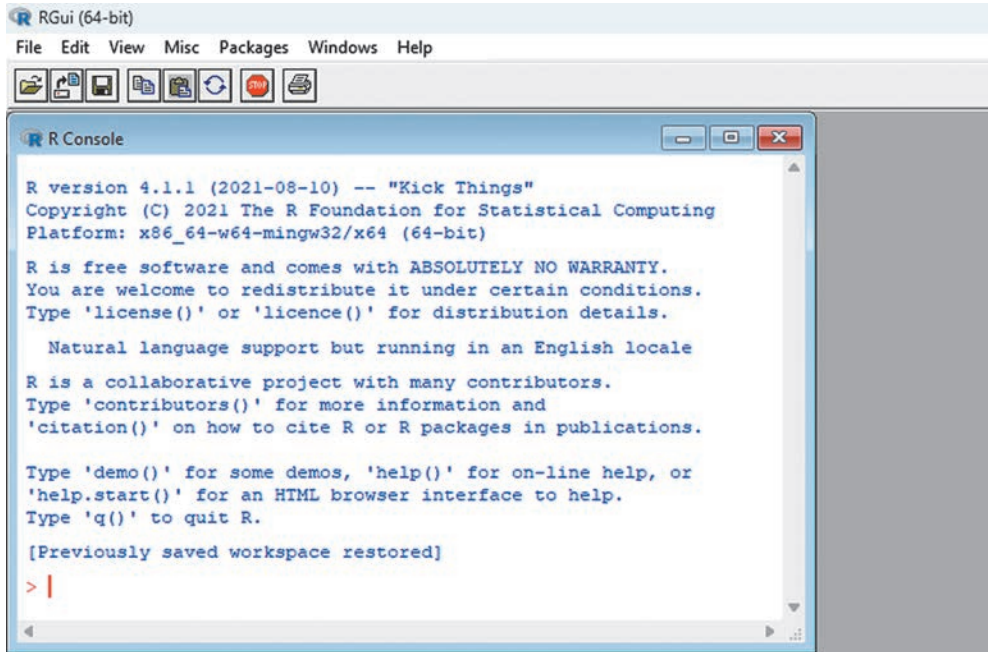
το, αφορά το σχεδιασμό και την εφαρμογή της διαδικασίας της συλλογής των πρωτογενών στοιχείων, είτε στο δασικό οικοσύστημα (πεδίο), είτε στο εργαστήριο μέσω εφαρμογής των αρχών του κλάδου του σχεδιασμού πειραμάτων. Η λήψη αυτή των πρωτογενών στοιχείων, όταν πρόκειται για λήψη στοιχείων δείγματος, ακολουθεί επιπλέον τις μεθοδολογίες και τις αρχές που αναπτύσσονται στον κλάδο της δειγματοληψίας. Το δεύτερο, αφορά την περιγραφή και διερευνητική ανάλυση των δεδομένων, μέσω κυρίως περιγραφικών στατιστικών μεθόδων, ενώ το τρίτο μεθοδολογικό κεφάλαιο αφορά την ανάλυση των δεδομένων μέσω εφαρμογής στατιστικών τεχνικών και διαδικασιών, οι οποίες οδηγούν τέλος στην στατιστική συμπερασματολογία.

Η γνώση της εφαρμογής των μεθόδων της εφαρμοσμένης στατιστικής στη δασική έρευνα είναι απαραίτητο να συνοδεύεται από τουλάχιστον τη βασική γνώση προγραμματισμού και γνώση χρήσης στατιστικών πακέτων (statistical softwares), προκειμένου να είναι δυνατή η επεξεργασία και ανάλυση πολλών αριθμητικά δεδομένων, αλλά και να υπάρχει η δυνατότητα εφαρμογής προηγμένων στατιστικών μεθόδων και τεχνικών, οι οποίες μπορούν να οδηγήσουν σε τελικά ασφαλή συμπεράσματα. Στα πλαίσια αυτού του βιβλίου, θα δοθούν στοιχεία χρήσης της γλώσσας προγραμματισμού R και η εφαρμογή της στην επίλυση στατιστικών διαδικασιών, οι οποίες θα αναπτυχθούν στα επόμενα κεφάλαια, καθώς και η αντίστοιχη χρήση του στατιστικού πακέτου IBM-SPSS (από εδώ και στο εξής η αναφορά στο πακέτο θα γίνεται ως SPSS) για την επίλυση στατιστικών προβλημάτων, με χρήση δασικών δεδομένων.

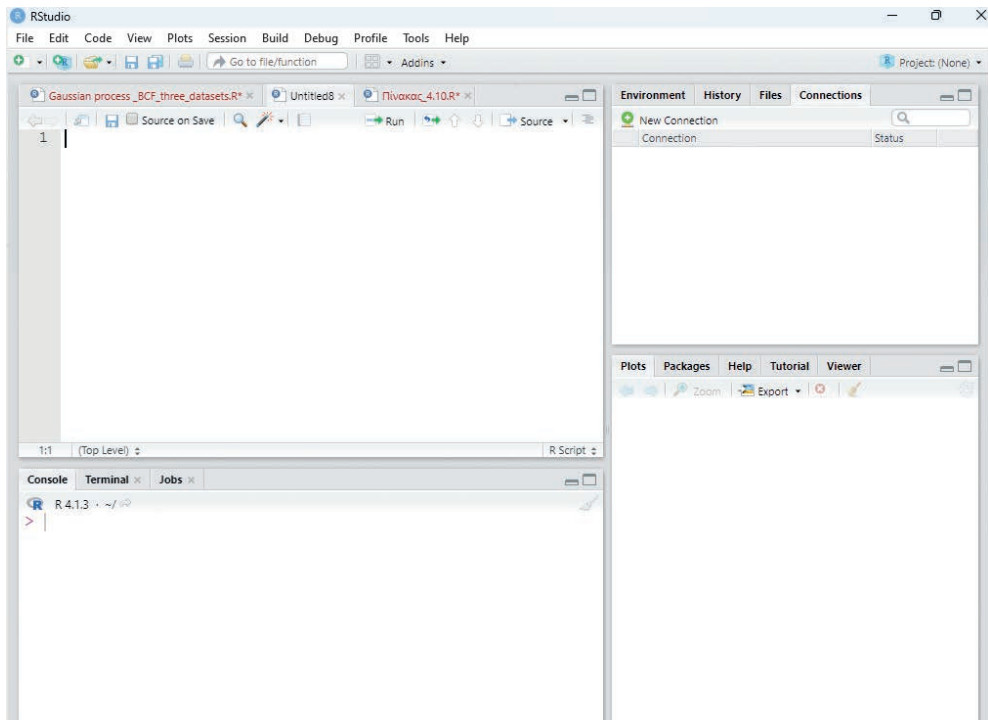
### 1.3 Εισαγωγή στη στατιστική γλώσσα προγραμματισμού R

Η γλώσσα προγραμματισμού R, αποτελεί ένα ελεύθερο λογισμικό, το οποίο διατίθεται για διάφορα λειτουργικά συστήματα (Windows, Linux, κλπ) από την ηλεκτρονική διεύθυνση <http://www.r-project.org/>. Στο συγκεκριμένο site δίνονται όλες οι απαραίτητες οδηγίες εγκατάστασης. Επίσης, η R μπορεί να αποκτηθεί και από ένα άλλο δίκτυο διανομής, το οποίο είναι ο πρότυπος καθρέφτης (mirror) του CRAN (Comprehensive R Archive), όπου η διαδικτυακή του διεύθυνση είναι <http://cran.r-project.org>.

Είναι δυνατό να “δουλέψει” στη βασική της έκδοση, αλλά και μέσα από το περιβάλλον RStudio (Εικόνα 1.1), το οποίο διανέμεται ελεύθερα από τη διεύθυνση <http://www.rstudio.org>.



(α)



(β)

Εικόνα 1.1. Περιβάλλον εργασίας της R α) σε βασική έκδοση και β) σε περιβάλλον RStudio.

Λεπτομερής ανάπτυξη της χρήσης της R, της σύνταξης των εντολών και της δημιουργίας γραφημάτων μπορεί ο αναγνώστης να βρει στην ιστοσελίδα της R και συγκεκριμένα στη διεύθυνση: <https://cran.r-project.org/doc/manuals/r-release/R-intro.pdf>. Η εκμάθηση της γλώσσας R δεν αποτελεί αντικείμενο και σκοπό του βιβλίου αυτού. Στο πλαίσιο των στατιστικών διαδικασιών οι οποίες αναπτύσσονται και αποτελούν το βασικό αντικείμενο του βιβλίου, θα δίνονται τμήματα κώδικα της R, έτσι ώστε να είναι εύκολη η υλοποίηση των στατιστικών διαδικασιών και η εξαγωγή των σχετικών αποτελεσμάτων. Παρ' όλα αυτά θα αναφερθούν κάποια πολύ βασικά στοιχεία της χρήσης της R, προκειμένου να είναι σε θέση ο αναγνώστης, ακόμη και αν δεν γνωρίζει προγραμματισμό στην R, να χρησιμοποιήσει τον κώδικα και να πάρει τα σχετικά αποτελέσματα.

Όπως φαίνεται στην Εικόνα 1.1, το βασικό περιβάλλον εργασίας (Εικόνα 1,α) δίνει ως βασικό μενού επιλογών (πρώτη μπάρα επιλογών) τα: File, Edit, View, Misc, Packages, Windows και Help, ενώ το περιβάλλον εργασίας RStudio (Εικόνα 1,β) δίνει ως βασικό μενού επιλογών (πρώτη μπάρα επιλογών) τα: File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools και Help. Επιπλέον, υπάρχει η δυνατότητα διαίρεσης της κεντρικής οθόνης σε 4 τμήματα, γεγονός το οποίο διευκολύνει την εγγραφή εντολών κώδικα, την παρακολούθηση της «εκτέλεσής» τους και την εμφάνιση των αποτελεσμάτων της «εκτέλεσης» του κώδικα. Γι' αυτό το λόγο, προτείνεται η χρήση της R μέσα από το περιβάλλον του Rstudio, ως ευκολότερη και σε μεγαλύτερο βαθμό κατανοητή η χρήση της, όσον αφορά τουλάχιστον τους μη έμπειρους χρήστες. Στη συνέχεια, η περιγραφή και ανάπτυξη του κώδικα και των διαδικασιών που θα υλοποιούνται με τη γλώσσα προγραμματισμού R, θα αφορούν το περιβάλλον RStudio. Τα αρχεία εντολών προγραμματισμού τα οποία δημιουργούνται στην R (scripts), είναι της μορφής όνομα\_αρχείου.R

Όπως φαίνεται στο κάτω αριστερό (όπως το βλέπει κάποιος στην οθόνη) ορθογώνιο του RStudio υπάρχει η κονσόλα στην οποία εγγράφονται εντολές μετά από το «>». Εφόσον ολοκληρωθεί η εγγραφή τους, πατώντας «enter», εκτελείται η εντολή και φαίνεται το αποτέλεσμα της. Οι βασικές αριθμητικές πράξεις εγγράφονται σύμφωνα με τον Πίνακα 1.1 με το πάτημα του πλήκτρου «enter». Επίσης κατά τον ίδιο τρόπο (Πίνακας 1.1) εγγράφονται και οι βασικές λογικές εκφράσεις οι οποίες έχουν ως αποτέλεσμα έκφραση αλήθειας (TRUE) ή λάθους (FALSE).



5. Ενσωματώνεται στο γράφημα η γραμμή των  $45^\circ$ , η οποία αντιπροσωπεύει την πλήρη εξίσωση των τιμών των ζευγών  $(x,y)$ , η οποία αποτελεί την ιδανική κατάσταση.

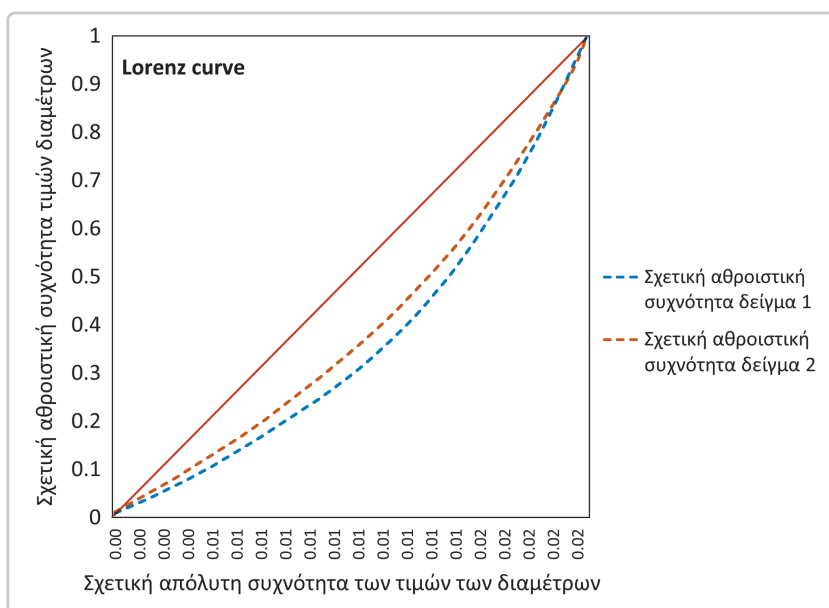
Όσο πιο μακριά βρίσκεται η καμπύλη του Lorenz των δεδομένων από την ιδανική κατάσταση, τόσο οι ανισότητες είναι μεγαλύτερες και η διασπορά μεγαλύτερη. Το συγκεκριμένο γράφημα είναι ιδανικό προκειμένου να συγκριθούν δύο ή και περισσότερες κατανομές συχνοτήτων μεταξύ τους. Η καμπύλη του Lorenz η οποία απέχει περισσότερο από την ιδανική κατάσταση, έχει τις μεγαλύτερες ανισότητες και τη μεγαλύτερη διακύμανση τιμών.

### Παράδειγμα 5.10

Για τις διαμέτρους του **Παραδείγματος 5.9** να γίνει γραφικός έλεγχος με την καμπύλη του Lorenz. Ποια από τις δύο εμπειρικές κατανομές εμφανίζει μεγαλύτερες ανισότητες και τελικά ποιο από τα δύο δείγματα είναι πιο ομοιογενές;

#### Απάντηση

Σύμφωνα με όσα αναπτύχθηκαν, δημιουργούνται οι καμπύλες του Lorenz για τα δύο δείγματα στο ίδιο γράφημα, προκειμένου να είναι δυνατή η σύγκρισή τους (Σχήμα 5.4).



Σχήμα 5.4. Η καμπύλη του Lorenz για τη σύγκριση των δύο εμπειρικών κατανομών συχνοτήτων των δειγμάτων 1 και 2.

Όπως φαίνεται στο Σχήμα 5.4, η εμπειρική κατανομή του δείγματος 1 εμφανίζει τη μεγαλύτερη απόκλιση από την ευθεία  $45^\circ$  γραμμή, σε σχέση με την αντίστοιχη κατανομή του δείγματος 2. Άρα επιβεβαιώνεται το συμπέρασμα το οποίο έχει εξαχθεί μέσω του υπολογισμού των συντελεστών κύμανσης, ότι δηλαδή το δείγμα 1 εμφανίζει τη μεγαλύτερη μεταβλητότητα.

Το γράφημα της καμπύλης προσφέρεται και από τις βιβλιοθήκες της R. Το σχετικό script σε γλώσσα προγραμματισμού R και το παραγόμενο αποτέλεσμα για το πρώτο δείγμα των τιμών διαμέτρων (δείγμα 1), χρησιμοποιώντας μόνο τις σχετικές αθροιστικές συχνότητες δίνεται στον Πίνακα 5.14. Όπως μπορεί κανείς να παρατηρήσει από το Σχήμα 5.4 και το παραγόμενο αποτέλεσμα του Πίνακα 5.14, είτε χρησιμοποιηθούν στον άξονα των  $x$  οι τιμές της σχετικής απόλυτης συχνότητας είτε οι τιμές της σχετικής αθροιστικής συχνότητας, το οπτικό αποτέλεσμα παραμένει παρόμοιο.

**Πίνακας 5.14.** Γράφημα της καμπύλης Lorenz στην R.

#### Κώδικας R (R-script)

```
install.packages("ineq")
library(ineq)
# Create data,
x<-c(12,10,8,12,10,8,9,10,8,8,9,12,8,12,11,12,11,8,12,9,12,11,8,.....,
     12,14,21,21,10,9)
# Lorenz Curve steps
# Sort the data in ascending order
sorted_data <- sort(x)
# Calculate the cumulative proportions
cumulative_proportions <- cumsum(sorted_data) / sum(sorted_data)
# Create a vector of cumulative proportions
cumulative_proportions_vector <- c(0, cumulative_proportions)
# Create a vector of the perfect equality line
perfect_equality_line <- seq(0, 1, length.out =
  length(cumulative_proportions_vector))
# Plot the Lorenz curve
plot(perfect_equality_line, cumulative_proportions_vector, type = "l",
     xlab = "Cumulative Proportion", ylab = "Cumulative Proportion",
     main = "Lorenz Curve")
# Add a diagonal reference line (perfect equality)
abline(0, 1, col = "red", lty = 2)
```



```

#Παράδειγμα 7.20
dbinom(2, size=10, prob=0.4)
> #Παράδειγμα 7.20
> dbinom(2, size=10, prob=0.4)
[1] 0.1209324
dbinom(7, size=10, prob=0.4)
> dbinom(7, size=10, prob=0.4)
[1] 0.04246733
>
#Παράδειγμα 7.21
dbinom(40, size=90, prob=0.33)
> #Παράδειγμα 7.21
> dbinom(40, size=90, prob=0.33)
[1] 0.0066300010

```

## 7.11 Κατανομή Poisson (Poisson distribution)

Η **κατανομή Poisson** (Poisson distribution), αποτελεί μια πολύ σημαντική διακριτή κατανομή, η οποία βρίσκει εφαρμογή σε πολλά φυσικά φαινόμενα. Ονομάζεται αλλιώς και ως κατανομή των σπάνιων γεγονότων, επειδή η εμφάνιση των γεγονότων είναι πιθανότερη στις μικρές τιμές της μεταβλητής  $X$ , ενώ οι μεγαλύτερες τιμές της μεταβλητής είναι σπάνιο να εμφανιστούν. Γενικά, φαινόμενα των οποίων η πιθανότητα εμφάνισής τους περιγράφεται από την κατανομή αυτή, είναι εκείνα τα οποία συμβαίνουν (εμφάνιση) σε ένα σταθερό διάστημα χρόνου ή/και χώρου, συμβαίνουν με ένα γνωστό μέσο ρυθμό, ενώ η εμφάνιση του επόμενου γεγονότος, είναι ανεξάρτητη από το διάστημα εμφάνισης του προηγούμενου. Παράδειγμα τέτοιων γεγονότων μπορεί να είναι ο αριθμός δασικών πυρκαγιών σε μια περιοχή, μέσα σε ορισμένο χρονικό διάστημα, ή ο αριθμός των τηλεφωνικών κλήσεων που φθάνουν σε ένα δασαρχείο, κλπ.

Η κατανομή Poisson, αποτελεί μια ειδική περίπτωση διωνυμικής κατανομής, η οποία εμφανίζεται όταν συμβαίνουν ταυτόχρονα οι εξής προϋποθέσεις: α) η πιθανότητα επιτυχίας του πειράματος είναι μικρότερη από 10% και β) ο μέσος όρος της κατανομής βρίσκεται μεταξύ του 0 και του 10. Τότε η πιθανότητα εμφάνισης του γεγονότος ακολουθεί την κατανομή Poisson, με μέσο όρο  $\lambda$ . Όταν μια μεταβλητή  $X$  ακολουθεί την κατανομή Poisson, με μέσο όρο  $\lambda$ , τότε συμβολίζεται ως  $X \sim P(\lambda)$  και η ποσότητα  $\lambda$  είναι η παράμετρος της κατανομής. Ο αριθμητικός μέσος της κατανομής είναι ίσος με

$$\lambda = n \cdot p,$$

ενώ η διακύμανσή της ισούται με

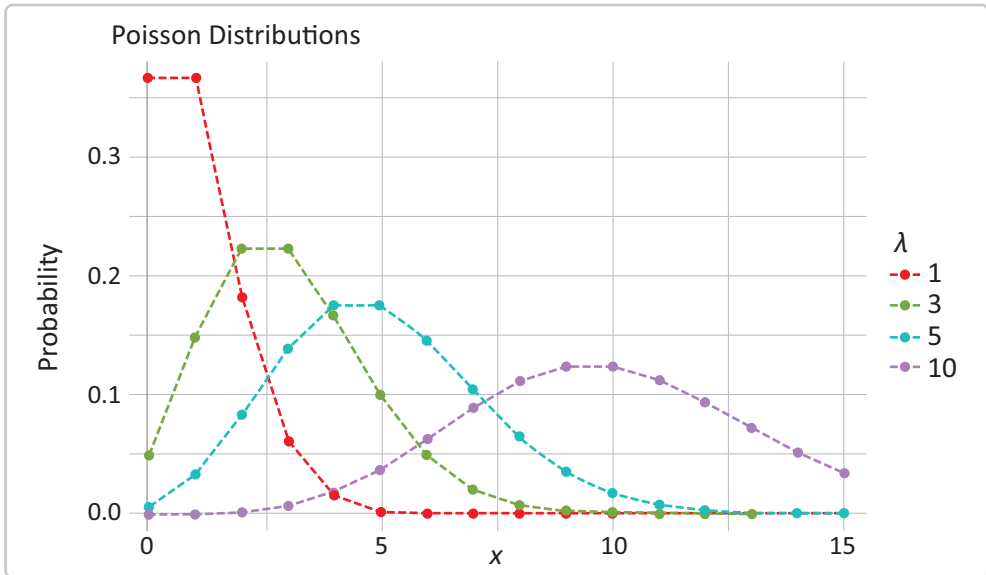
$$\sigma^2 = \lambda .$$

Δηλαδή, ο μέσος όρος της κατανομής είναι ίσος με την τιμή της διασποράς της.

Επειδή η τιμή της διακύμανσης μπορεί να διαφοροποιηθεί εξαιτίας της διακύμανσης της δειγματοληψίας, θεωρούμε ότι αν  $\sigma^2 = \lambda \pm 0,2 \cdot \lambda$ , τότε ισχύει κατά προσέγγιση η ισότητα διακύμανσης με το μέσο όρο της κατανομής. Η συνάρτηση πιθανότητας, η οποία περιγράφει την κατανομή Poisson για τη μεταβλητή  $X \sim P(\lambda)$ , η οποία παίρνει την τιμή  $x$ , είναι:

$$P(X=x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!} \quad (7.30)$$

Η γραφική παράσταση διαφόρων κατανομών Poisson, με διαφορετικούς μέσους όρους ( $\lambda$ ), δίνεται στο Σχήμα 7.4, όπου η ένωση των σημείων με διακεκομμένη γραμμή έχει γίνει μόνο για λόγους καλύτερης κατανόησης των διαφορετικών μορφών που παίρνει η κατανομή, για διαφορετικές τιμές του μέσου όρου της και για τιμές της μεταβλητής  $X$  από 0 έως 15. Πρόκειται για μία διακριτή κατανομή και οι τιμές μεταξύ των σημείων δεν έχουν νόημα.



Σχήμα 7.4. Γραφική παράσταση κατανομής Poisson, για διαφορετικές τιμές  $\lambda$ .

Για τον υπολογισμό των πιθανοτήτων της κατανομής Poisson, όπως και στην περίπτωση της διωνυμικής κατανομής, μπορούν να δημιουργηθούν πίνακες είτε απλοί, είτε αθροιστικής κατανομής. Προκειμένου να δημιουργηθεί πίνακας απλών πιθανοτήτων της κατανομής Poisson, έτσι ώστε η τυχαία μεταβλητή  $X$  να πάρει κάποια συγκεκριμένη τιμή, η οποία εκφράζει την πιθανότητα επιτυχίας, αρκεί να χρησιμοποιηθεί η σχέση (7.30). Η εύρεση της αθροιστικής πιθανότητας,

## 7.13 Κανονική και τυπική κανονική κατανομή (normal and standard normal distribution)

Η **κανονική κατανομή** θεωρείται ως η σπουδαιότερη συνεχής κατανομή πιθανοτήτων, γιατί βρίσκει πάρα πολλές εφαρμογές σε προβλήματα της πράξης, προσεγγίζοντας πολλά φυσικά φαινόμενα.

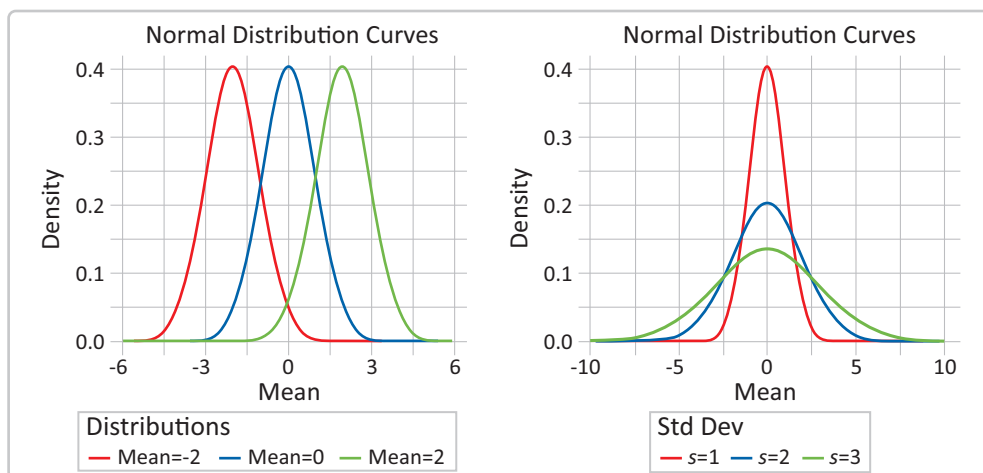
Πολλοί ερευνητές ασχολήθηκαν με την κανονική κατανομή. Πρώτος τη διατύπωσε ο μαθηματικός De Moivre γύρω στο 1733, ο οποίος διαπίστωσε ότι στη διωνυμική κατανομή, καθώς το μέγεθος του δείγματος τείνει στο άπειρο, το διωνυμικό ανάπτυσμα  $(p+q)^n$  τείνει να προσεγγίσει την κανονική κατανομή, ενώ αργότερα ξαναμελετήθηκε από τους Gauss και Laplace γύρω στο 1812, ως κατανομή των σφαλμάτων.

Πρόκειται για μια συμμετρική κατανομή, κωδωνοειδούς μορφής, με ασύμπτωτο τον άξονα των τετμημένων, η οποία έχει δύο παραμέτρους και συμβολίζεται με:  $X \sim N(\mu, \sigma^2)$ . Δηλαδή, η τυχαία μεταβλητή  $X$  η οποία είναι συνεχής, ακολουθεί την κανονική κατανομή με μέσο όρο  $\mu$  και διασπορά  $\sigma^2$ .

Η συνάρτηση πυκνότητας-πιθανότητας της κατανομής δίνεται από την σχέση:

$$f(x) = \frac{1}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot e^{-0,5 \cdot \left(\frac{x-\mu}{\sigma}\right)^2} \quad (7.32)$$

Στο Σχήμα 7.7, δίνεται η γραφική παράσταση διαφορετικών κανονικών κατανομών, οι οποίες έχουν διαφορετικούς μέσους όρους και ίδια τιμή διακύμανσης (Σχήμα 7.7α) και διαφορετικές τιμές διακύμανσης και ίσους μέσους όρους (Σχήμα 7.7β).

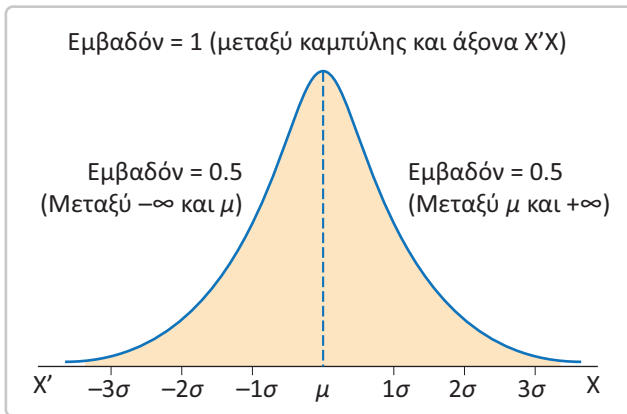


Σχήμα 7.7. Διαφορετικές μορφές κανονικών κατανομών.

Όπως διαπιστώνεται, η κανονική κατανομή ανεξάρτητα από την τιμή του μέσου όρου της και της διακύμανσής της είναι κωδωνοειδούς μορφής, συμμετρική γύρω από το μέσο όρο της. Η διαφορετική τιμή μέσου όρου χαρακτηρίζει τη θέση της κατανομής ως προς τον άξονα των  $X'$ , ενώ η διαφορετική τιμή της διακύμανσης της κατανομής, άρα και της τυπικής απόκλισης, καθορίζει τη μορφή της κατανομής.

Θεωρείται ότι το εμβαδόν μεταξύ της συμμετρικής κωδωνοειδούς καμπύλης της κανονικής κατανομής και του άξονα  $X'$ , είναι ίσο με ένα (Σχήμα 7.8), το οποίο μοιράζεται εξίσου εκατέρωθεν του κάθετου άξονα στον  $X'$ , ο οποίος περνάει από το μέσο όρο της κατανομής. Αυτό μαθηματικά εκφράζεται με το αόριστο ολοκλήρωμα της συνάρτησης πυκνότητας πιθανότητας:

$$\int_{-\infty}^{+\infty} f(x) d(x) = \int_{-\infty}^{+\infty} \frac{1}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot e^{-0.5 \cdot \left(\frac{x-\mu}{\sigma}\right)^2} d(x) = 1 \quad (7.33)$$



Σχήμα 7.8. Εμβαδόν μεταξύ κανονικής κατανομής και άξονα  $x'$ .

Οι αθροιστικές πιθανότητες της κατανομής  $F(x)$ , δηλαδή η πιθανότητα η τυχαία μεταβλητή  $X \sim N(\mu, \sigma^2)$  να πάρει τιμές σε εύρος μικρότερο από μία συγκεκριμένη τιμή  $x$ , δηλ.  $P(X \leq x)$ , αποδίδεται από τη γραφική παράσταση της αθροιστικής συνάρτησης της κανονικής κατανομής, η οποία φαίνεται στο Σχήμα 7.9. Στο Σχήμα 7.9α φαίνεται η γραφική παράσταση αθροιστικών κανονικών κατανομών με διαφορετικούς μέσους όρους, αλλά ίση διακύμανση, ενώ στο Σχήμα 7.9β φαίνεται η γραφική παράσταση αθροιστικών κανονικών κατανομών με ίσους μέσους όρους, αλλά διαφορετικές διακυμάνσεις.

**Coefficients:**

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	<b>1.892500</b>	<b>0.092774</b>	<b>20.40</b>	9.02e-07 ***
X	0.075167	0.003674	20.46	8.87e-07 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1191 on 6 degrees of freedom

**Multiple R-squared: 0.9859**, Adjusted R-squared: 0.9835**F-statistic: 418.5 on 1 and 6 DF, p-value: 8.873e-07**

Από τη γραφική παράσταση του παραγόμενου αποτελέσματος του Πίνακα 9.9, φαίνεται η άριστη προσαρμογή της εξίσωσης ευθείας γραμμής στα μετρημένα δεδομένα, τα οποία αναπαρίστανται με στικτό διάγραμμα.

### 9.3.2 Πολλαπλή γραμμική παλινδρόμηση (multiple linear regression, MLR)

Η πολλαπλή γραμμική παλινδρόμηση αποτελεί μια προέκταση της απλής γραμμικής παλινδρόμησης. Δηλαδή επιλύει τα γραμμικά εκείνα μοντέλα, τα οποία έχουν περισσότερες από μία ανεξάρτητες μεταβλητές. Και σ' αυτή την περίπτωση, όπως και στην απλή γραμμική παλινδρόμηση, υφίστανται οι περιορισμοί της ύψωσης σε πρώτη δύναμη των ανεξάρτητων μεταβλητών όπως επίσης και ο τρόπος σύνδεσης των συντελεστών. Εδώ μπορούν να συμπεριληφθούν και τα μοντέλα εκείνα τα οποία δεν είναι γραμμικά, όπως πχ. εκθετικά μοντέλα, τα οποία με διάφορες διαδικασίες μπορούν να μετασχηματιστούν σε γραμμικά.

Η γενική μορφή του τυπικού μοντέλου που επιλύεται με την εφαρμογή της πολλαπλής παλινδρόμησης, είναι:

$$y_i = \hat{b}_0 + \sum_{k=1}^{p-1} \hat{b}_k \cdot x_{ik} + \varepsilon_i \quad (9.33)$$

όπου  $p$  είναι ο αριθμός των συντελεστών παλινδρόμησης του μοντέλου.

Η τιμή του συντελεστή  $\hat{b}_0$ , δίνει η τιμή που μπορεί να πάρει η μεταβλητή  $Y$ , όταν γίνει ταυτόχρονη μηδένιση των τιμών όλων των ανεξάρτητων μεταβλητών.

Ο συντελεστής  $\hat{b}_k$  δείχνει τη μέση μεταβολή (αύξηση ή μείωση) της εξαρτημένης μεταβλητής, όταν η μεταβλητή  $x_{ik}$  μεταβληθεί κατά μία μονάδα, με την προϋπόθεση ότι όλες οι υπόλοιπες  $(p-2)$  ανεξάρτητες μεταβλητές παραμένουν στα-

θερές, ανεξάρτητα από την τιμή που έχουν. Κατ' επέκταση, το ίδιο ισχύει και για τους υπόλοιπους συντελεστές παλινδρόμησης.

Οι ίδιες βασικές προϋποθέσεις που αναφέρθηκαν στην περίπτωση της απλής γραμμικής παλινδρόμησης επεκτείνονται και ισχύουν και για την περίπτωση της εφαρμογής της πολλαπλής γραμμικής παλινδρόμησης. Επιπρόσθετα, επειδή στην εξίσωση υπάρχουν περισσότερες από μία ανεξάρτητες μεταβλητές, αυτές δεν θα πρέπει να εμφανίζουν υψηλή συσχέτιση μεταξύ τους, δηλαδή δεν θα πρέπει να εμφανίζεται **πολυσυγγραμμότητα** (multicollinearity).

Το αντιπροσωπευτικό μοντέλο της πολλαπλής παλινδρόμησης μπορεί να δοθεί και με μορφή πινάκων ως εξής:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1,p-1} \\ 1 & x_{21} & \cdots & x_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{n,p-1} \end{bmatrix} \cdot \begin{bmatrix} \hat{b}_0 \\ \hat{b}_1 \\ \vdots \\ \hat{b}_{p-1} \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} \Leftrightarrow Y = X \cdot \hat{b} + \varepsilon \quad (9.34)$$

Όπως φαίνεται στη σχέση (9.34), η εξαρτημένη μεταβλητή  $Y$  είναι ένας πίνακας διαστάσεων  $n$  γραμμών επί μία στήλη:  $(n \times 1)$ , οι ανεξάρτητες μεταβλητές είναι ένας πίνακας διαστάσεων  $(n \times p)$ , οι συντελεστές παλινδρόμησης αποτελούν έναν πίνακα διαστάσεων  $(p \times 1)$  και τέλος τα σφάλματα είναι ένας πίνακας διαστάσεων  $(n \times 1)$ .

Στην πολλαπλή γραμμική παλινδρόμηση, οι παράμετροι του πληθυσμού εκτιμώνται μέσω της μεθόδου των ελαχίστων τετραγώνων κατά τρόπο ανάλογο μ' αυτόν που προαναφέρθηκε. Η μόνη διαφορά έγκειται στο γεγονός ότι εδώ το σύστημα εξισώσεων που προκύπτει δεν είναι  $2 \times 2$  όπως στην περίπτωση της απλής γραμμικής παλινδρόμησης, αλλά  $p \times p$ , οπότε γίνεται ευρεία χρήση της θεωρίας των πινάκων βάσει της οποίας λύνονται οι περισσότερες σχέσεις που ακολουθούν. Δηλαδή, επιδιώκεται η ελαχιστοποίηση του αθροίσματος των τετραγώνων των σφαλμάτων της παραγόμενης εξίσωσης παλινδρόμησης μέσω του μηδενισμού των μερικών παραγώγων του αθροίσματος αυτού ως προς τον καθένα από τους συντελεστές παλινδρόμησης. Οι εκτιμητές που προκύπτουν κατά αυτόν τον τρόπο, ενσωματώνουν όλες τις επιθυμητές ιδιότητες των εκτιμητών, δηλαδή είναι αμερόληπτοι, αποτελεσματικοί, επαρκείς και συνεπείς. Υπολογίζονται με βάση τη σχέση:

$$\hat{b} = (X'X)^{-1} \cdot (X'Y) \quad (9.35)$$

όπου  $X'$  είναι ο ανάστροφος πίνακας του  $X$  και έχει διαστάσεις  $(p \times n)$  και  $(X'X)^{-1}$

### Παράδειγμα 9.7

Σε δείγμα 94 δέντρων, μετρήθηκαν οι διάμετροι των κορμών τους, ανά ένα μέτρο και το ολικό τους ύψος. Με εφαρμογή μεθόδου τμηματικής ογκομέτρησης υπολογίστηκε ο ολικός έμφλοιος κορμικός όγκος τους.

Να καταρτιστεί εξίσωση μη-γραμμικής παλινδρόμησης με εξαρτημένη μεταβλητή τον όγκο των δέντρων και ανεξάρτητη μεταβλητή το ολικό ύψος τους.

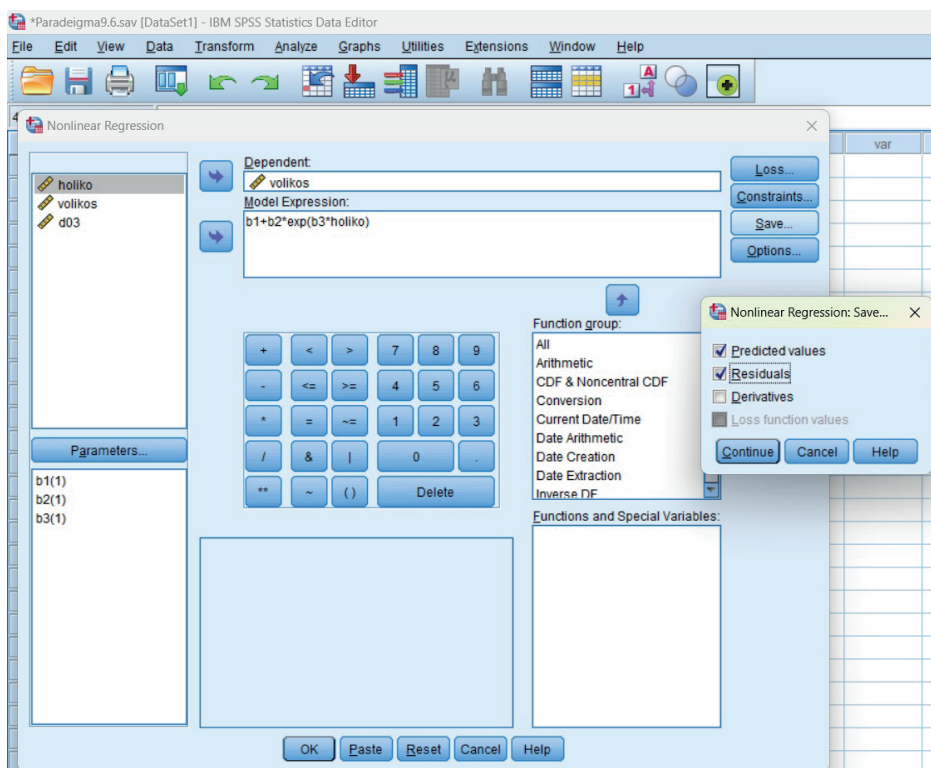
#### Απάντηση

Πριν την εφαρμογή μη-γραμμικής παλινδρόμησης, θα πρέπει να εξαντληθεί κάθε πιθανότητα, ένα γραμμικό μοντέλο να έχει τη δυνατότητα να περιγράψει με ικανοποιητικό τρόπο, τα δεδομένα. Μετά τη διερεύνηση αυτή και εφόσον η απλή γραμμική παλινδρόμηση και η πολλαπλή γραμμική παλινδρόμηση απέτυχαν να παράγουν ικανοποιητικά αποτελέσματα, τότε εφαρμόζεται η μη-γραμμική παλινδρόμηση. Η εφαρμογή της στο SPSS περιγράφεται στον Πίνακα 9.13.

**Πίνακας 9.13.** Επίλυση προβλήματος μη-γραμμικής παλινδρόμησης, με το SPSS.

**Ένδειξη οθόνης**

	holiko	volikos
1	3.30	.0081
2	4.30	.0225
3	4.30	.0363
4	4.30	.0096
5	4.40	.0260
6	4.40	.0284
7	4.50	.0166
8	4.70	.0149
9	4.80	.0478
10	4.80	.0469
11	4.80	.0176
12	4.80	.0169
13	4.90	.0173
14	4.90	.0293
15	4.90	.0199
16	4.90	.0193
17	4.90	.0230
18	5.00	.0277
19	5.00	.0206
20	5.00	.0144
21	5.00	.0188
22	5.10	.0752
23	5.30	.0509
24	5.30	.0325



## Περιγραφή

1. Εισάγουμε τα δεδομένα
2. Επιλέγουμε Regression → Nonlinear
3. Στη θυρίδα διαλόγου που αναδύεται τοποθετούμε την εξαρτημένη μεταβλητή στο «Dependent» γράφουμε με κωδικοποίηση του SPSS την εξίσωση την οποία επιθυμούμε να προσαρμόσει το SPSS στα δεδομένα
4. Από τις διαθέσιμες εντολές επιλέγουμε «Save» και από το παράθυρο που αναδύεται επιλέγουμε «Predicted values», «Residuals»
5. Στη συνέχεια επιλέγουμε «Continue» και OK
6. Τέλος, μέσω της επιλογής «Graphs» από το κύριο menu των επιλογών επιλέγουμε «Regression Variable Plots» και δημιουργούμε το στικτόδιάγραμμα των δεδομένων με τα αντίστοιχα θηκογράμματα και την γραμμή προσαρμογής του μοντέλου

## Εντολές σε γλώσσα SPSS:

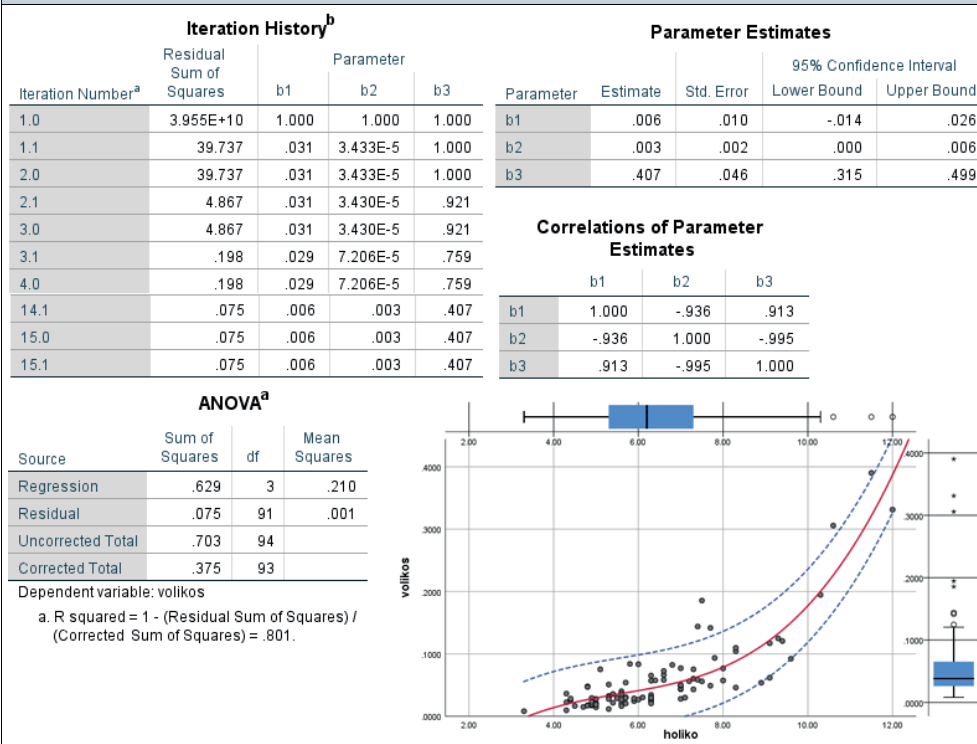
```
DATASET ACTIVATE DataSet1.
* NonLinear Regression.
MODEL PROGRAM b1=1 b2=1 b3=1.
COMPUTE PRED_ =b1+b2*exp(b3*holiko).
```



```
NLR volikos
/OUTFILE='C:\Users\maria\spss22540\SPSSFNLR.TMP'
/PRED PRED_
/SAVE PRED RESID
/CRITERIA SCONVERGENCE 1E-8 PCON 1E-8.
```

```
STATS REGRESS PLOT YVARS=volikos XVARS=holiko
/OPTIONS CATEGORICAL=BAR GROUP=1 BOXPLOTS
INDENT=15 YSCALE=75
/FITLINES APPLYTO=TOTAL.
```

**Παραγόμενο αποτέλεσμα:**



Όπως φαίνεται στο συγκεκριμένο Παράδειγμα 9.7, ζητήθηκε η προσαρμογή του ασυμπτωτικού μη-γραμμικού μοντέλου:  $\hat{y}_i = b_1 + b_2 \cdot e^{(b_3 \cdot x_i)}$ , το οποίο είναι εσωτερικά μη-γραμμικό, δηλαδή είναι αδύνατο να μετασχηματιστεί σε γραμμική μορφή, εκτός αν τεθεί η προϋπόθεση ότι  $b_1 = 0$ . Το μοντέλο επιλύθηκε με την εφαρμογή του αλγορίθμου βελτιστοποίησης των Levenberg-Marquardt, με εφαρμογή επαναληπτικής διαδικασίας, ενώ χρησιμοποιήθηκε το μέσο τετραγωνικό σφάλμα, ως το μέτρο αξιολόγησης του μοντέλου.

Πίνακας 9.16. Εφαρμογή της παλινδρόμησης Ridge, στην R.

**Κώδικας R (R-script)**

```

library(glmnet)
library(car)
# Load the data
data <- read.csv("C:/Users/MD/ Ridge.csv")
colnames(data) <- make.names(colnames(data))
colnames(data)
colnames(data) <- c("h", "d1", "d2", "d3", "d4", "y")
# Prepare the predictor and response variables
x <- as.matrix(data[, c('h', 'd1', 'd2', 'd3', 'd4')])
y <- data$y
# Fit the ridge regression model with a sequence of lambda values
ridge_model <- glmnet(x, y, alpha = 0)
# Plot the Ridge trace
plot(ridge_model, xvar = "lambda", label = TRUE)
title(main = "Ridge Trace")
# Perform cross-validation to find the best lambda
cv_ridge <- cv.glmnet(x, y, alpha = 0)
best_lambda <- cv_ridge$lambda.min
print(best_lambda)
# Fit the final model with the best lambda
final_ridge_model <- glmnet(x, y, alpha = 0, lambda = best_lambda)
ridge_coefficients <- coef(final_ridge_model)
print(ridge_coefficients)

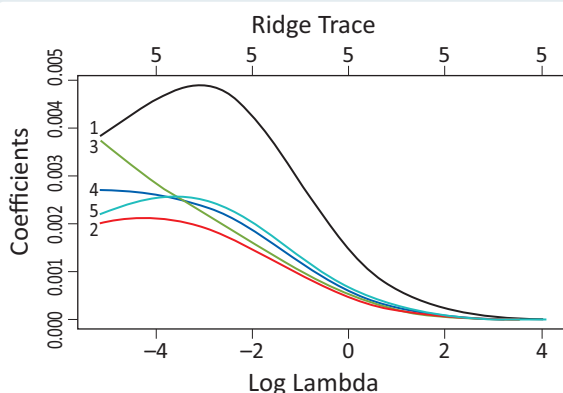
```

**Παραγόμενο αποτέλεσμα**

```

> print(best_lambda)
[1] 0.00602468
> ridge_coefficients <-
coef(final_ridge_model)
> print(ridge_coefficients)
(Intercept)  -0.101289505
h             0.003833729
d1           0.002030172
d2           0.003749633
d3           0.002708144
d4           0.002194844

```



Σε περίπτωση που πρέπει να αναλυθούν ταυτόχρονα ποσοτικές και κατηγορικές μεταβλητές, τότε η κοινή ανάλυσή τους θα πρέπει να χρησιμοποιεί τη μετρική απόστασης Gower, η οποία λειτουργεί με αποδεκτό τρόπο για συνδυασμό ποσοτικών και ποιοτικών μεταβλητών.

**Πίνακας 10.15.** Επιλογή κατάλληλων τεχνικών στην Ιεραρχική ταξινόμηση, ανάλο-γα με τον τύπο των μεταβλητών που αναλύονται.

Βασικές διαθέσιμες τεχνικές στην Ιεραρχική ταξινόμηση				
Μεταβλητές	Μετρική απόστασης	Μετασημα-τισμός	Μέθοδος σύνδεσης	Δενδρόγραμμα
<b>Ποσοτικές</b>	Ευκλείδεια, L1 (Manhattan), Cosine, Chebyshev, Block, and Minkowski, Pearson correlation	Μπορεί να χρησιμοποιηθεί κανονικοποίηση	Simple (Nearest-neighbor), Complete (Furthest-Neighbor), Ward's, Centroid clustering, Median clustering	Η συγχώνευση βασίζεται σε αριθμητικές αποστάσεις
<b>Κατηγορι-κές</b>	Gower, συντελεστής απλός ομοιότητας, δείκτης του Jaccard	Απαραίτητη η κωδικοποίηση	Simple (Nearest-neighbor), Complete (Furthest-neighbor), Ward's	Η συγχώνευση βασίζεται σε κατηγορική ομοιότητα

### Παράδειγμα 10.10

Ερωτήθηκαν εννιά άτομα τυχαία από μια γειτονιά, στην οποία πρόκειται να κατασκευαστεί πάρκο τσέπης, σχετικά με τις προτιμήσεις τους σε είδος δέντρου (μεταξύ τριών διαφορετικών ειδών), σχήμα (μεταξύ τριών διαφορετικών σχημάτων) και μέγεθος πάρκου (μεταξύ τριών διαφορετικών μεγεθών). Οι απαντήσεις που δόθηκαν καταγράφονται στον παρακάτω πίνακα:

Ερωτώμενος	Απαντήσεις σχετικά με		
	Είδος δέντρου	Σχήμα	Μέγεθος
1	είδος1	κυκλικό	μικρό
2	είδος2	τετραγωνικό	μεγάλο
3	είδος3	τριγωνικό	μέσο
4	είδος2	κυκλικό	μεγάλο
5	είδος1	τετραγωνικό	μικρό

Ερωτώμενος	Απαντήσεις σχετικά με		
	Είδος δέντρου	Σχήμα	Μέγεθος
6	είδος2	τετραγωνικό	μεγάλο
7	είδος3	τριγωνικό	μέσο
8	είδος2	κυκλικό	μεγάλο
9	είδος1	τετραγωνικό	μικρό

Μπορεί να ομαδοποιηθεί η συμπεριφορά των ερωτώμενων ως προς τις προτιμήσεις τους;

### Απάντηση

Πρόκειται για τη διερεύνηση των προτιμήσεων 9 ατόμων, των οποίων οι απαντήσεις δημιούργησαν τρεις κατηγορικές μεταβλητές, τριών επιπέδων η καθεμία. Θα χρησιμοποιηθεί η R, για την εφαρμογή της ιεραρχικής ταξινόμησης στα κατηγορικά δεδομένα. Στον Πίνακα 10.16, δίνεται ο βασικός κώδικας και τα παραγόμενα αποτελέσματα.

**Πίνακας 10.16.** Εφαρμογή της τεχνικής της ιεραρχικής ταξινόμησης, με την R.

#### Κώδικας R (R-script)

```
install.packages("cluster")
library(cluster)
# Sample data with categorical variables
data <- data.frame(
  Color = as.factor(c("1", "2", "3", "2", "1", "2", "3", "2", "1")),
  Shape = as.factor(c("Circle", "Square", "Triangle", "Circle", "Square", "Square",
"Triangle", "Circle", "Square")),
  Size = as.factor(c("Small", "Large", "Medium", "Large", "Small", "Large", "Me-
dium", "Large", "Small"))
)
print(data)
# Compute Ward distance matrix
Wards_dist <- daisy(data, metric = "Ward")
# Print distance matrix
print(as.matrix(Wards_dist))
# Perform hierarchical clustering
hc <- agnes(Wards_dist, method = "average")
```

```
# Print clustering result
```

```
print(hc)
```

```
# Plot dendrogram
```

```
plot(hc, which.plots = 2, main = "Dendrogram of Hierarchical Clustering")
```

### Παραγόμενο αποτέλεσμα

```
# Print distance matrix
```

```
> print(as.matrix(Wards_dist))
```

	1	2	3	4	5	6	7	8	9
1	0.000000	1.000000	1	0.6666667	0.3333333	1.000000	1	0.6666667	0.3333333
2	1.000000	0.000000	1	0.3333333	0.6666667	0.000000	1	0.3333333	0.6666667
3	1.000000	1.000000	0	1.000000	1.000000	1.000000	0	1.000000	1.000000
4	0.6666667	0.3333333	1	0.000000	1.000000	0.3333333	1	0.000000	1.000000
5	0.3333333	0.6666667	1	1.000000	0.000000	0.6666667	1	1.000000	0.000000
6	1.000000	0.000000	1	0.3333333	0.6666667	0.000000	1	0.3333333	0.6666667
7	1.000000	1.000000	0	1.000000	1.000000	1.000000	0	1.000000	1.000000
8	0.6666667	0.3333333	1	0.000000	1.000000	0.3333333	1	0.000000	1.000000
9	0.3333333	0.6666667	1	1.000000	0.000000	0.6666667	1	1.000000	0.000000

```
>
```

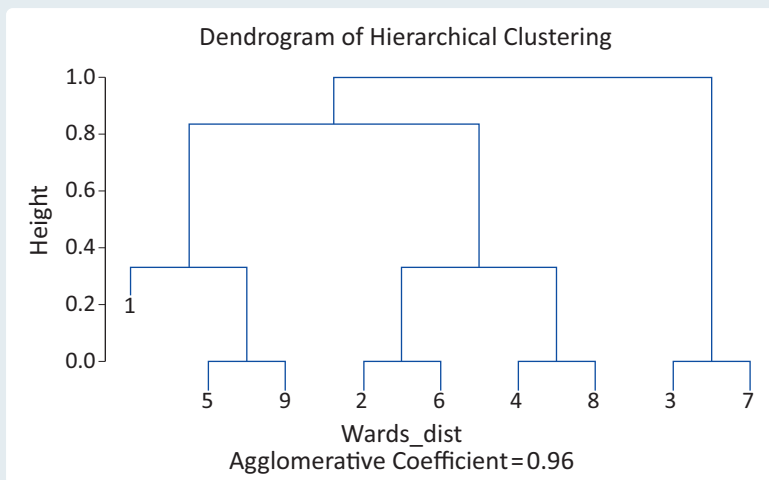
Agglomerative coefficient: **0.962963**

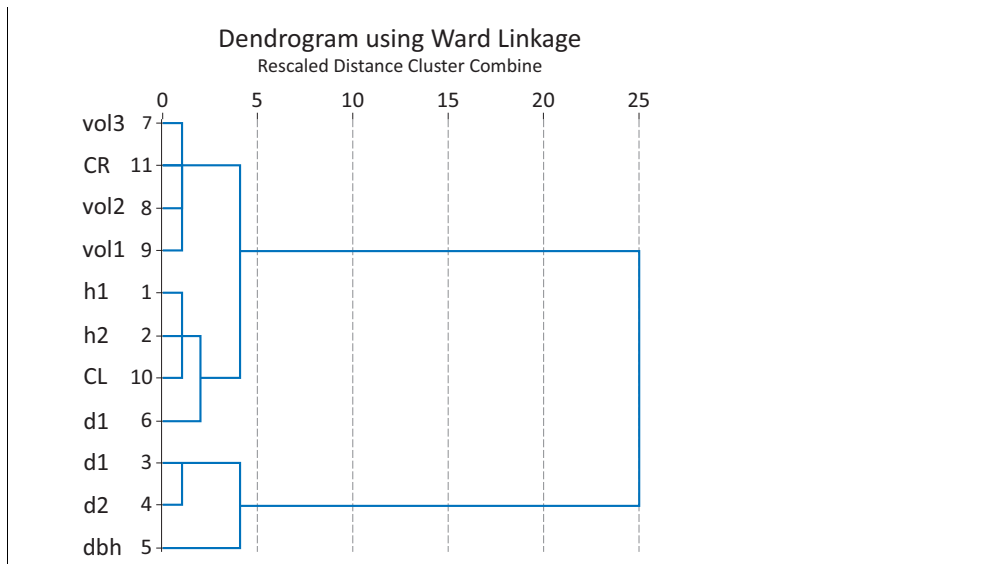
Order of objects:

**[1] 1 5 9 2 6 4 8 3 7**

Height (summary):

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.0000	0.0000	0.1667	0.3125	0.4583	1.0000





Όπως φαίνεται στον Πίνακα 10.17, ο πρώτος πίνακας δίνει τον αριθμό των στοιχείων ανά στήλη τα οποία αναλύθηκαν και την πληροφορία της ύπαρξης ελλειπουσών τιμών η οποία είναι πολύ σημαντική, γιατί μεγάλο ποσοστό ελλειπουσών τιμών μπορεί να οδηγήσει σε παραπλανητικά αποτελέσματα.

Ο επόμενος πίνακας (Proximity Matrix) δίνει όπως αναφέρει και το όνομά του πληροφορίες σχετικά με την εγγύτητα μεταξύ των μεταβλητών. Δηλαδή οι μεταβλητές που έχουν μικρότερη απόσταση μεταξύ τους είναι περισσότερο όμοιες. Όπως μπορεί να παρατηρήσει κανείς, η διαγώνιος του πίνακα είναι μηδέν, γιατί η απόσταση της κάθε μεταβλητής από τον εαυτό της είναι μηδενική. Οι αποστάσεις αυτές που δίνονται στον πίνακα εγγύτητας, είναι οι τετραγωνικές Ευκλείδειες αποστάσεις γιατί αυτή η επιλογή έγινε κατά την κατάστρωση της διαδικασίας.

Ο πίνακας της συσσώρευσης (Agglomeration Schedule) δίνει πληροφόρηση σχετικά με τη διαδικασία με την οποία εξελίσσεται η ομαδοποίηση. Στον πίνακα που έχει παραχθεί αξίζει να προσεχθούν οι ακόλουθες πληροφορίες. Φαίνονται 10 βήματα (Stage), σε καθένα από τα οποία μια μεταβλητή (δίνεται ο αύξοντας αριθμός της) μπορεί να αποτελέσει ομάδα με κάποια άλλη κοντινή, από άποψη ομοιότητας, μεταβλητή. Στη συνέχεια, οι συντελεστές ποσοτικοποιούν την διαφορετικότητα μεταξύ των μεταβλητών. Μεγάλη τιμή του συντελεστή αντιπροσωπεύει μεγαλύτερη διαφορετικότητα μεταξύ των μεταβλητών. Αυτή η στήλη των τιμών των συντελεστών μπορεί να αξιοποιηθεί για τη βέλτιστη επιλογή του αριθμού των ομάδων. Δηλαδή, στο σημείο που εμφανίζεται μια απότομη και μεγάλη σε μέγεθος αλλαγή τιμών, σηματοδοτείται η πιθανή βέλτιστη τιμή του αριθμού των ομάδων.

## Ευρετήριο όρων

- αθροιστικές συχνότητες ( $\Phi$ ), 42
- ακρίβεια (accuracy), 423
- αμερόληπτος εκτιμητής, 202
- αναλογία διασποράς (variation ratio, Vr), 357
- αναλογίας της εγκυρότητας του ερωτηματολογίου (Content Validity Ratio, CVR), 423
- ανάλυση διακύμανσης (ANalysis Of Variance, ANOVA), 245
- ανάλυση συστάδων (cluster analysis), 396
- ανεξάρτητες μεταβλητές, 299
- ανεξάρτητη μεταβλητή, 279
- αξιοπιστία (reliability), 422
- αξιοπιστία (reliability) ερωτηματολογίου, 425
- αξιοπιστία εσωτερικής συνέπειας (Internal Consistency Reliability), 425
- απλή γραμμική παλινδρόμηση (simple linear regression, SLR), 299
- απλό γεγονός ή ενδεχόμενο, 134
- απογραφές, 21
- αποκομμένος μέσος όρος (trimmed mean), 61
- αποκοπή (intercept), 302
- απόλυτες συχνότητες (fi), 42
- αποτελεσματικός εκτιμητής, 203
- αριθμητικός μέσος όρος ή μέση τιμή (arithmetic mean or average), 51
- αρμονικός μέσος (harmonic mean), 68
- ασυμμετρία (skewness), 117
- ασυνεχής ή διακριτή μεταβλητή, 157
- βαθμοί ελευθερίας, 306
- βασικά στοιχεία συνδυαστικής ανάλυσης (combinatorial analysis), 151
- βασικές ιδιότητες της διακύμανσης, 102
- βασικές μέθοδοι προσδιορισμού του δειγματοχώρου, 155
- βασικοί ορισμοί και έννοιες της θεωρίας των πιθανοτήτων – βασικοί συμβολισμοί, 134
- βέννεια διαγράμματα (venn diagrams), 137
- βιβλιογραφική αναφορά πίνακα, 33
- γεγονός, 134
- γενικός τίτλος πίνακα, 33
- γεωμετρικός μέσος (geometric mean), 64
- γραμμική παλινδρόμηση, 299
- γραφήματα, 33
- Δασική Στατιστική, 4
- δείγμα, 18
- δειγματοληπτικά σχέδια με πιθανότητα, 18
- δειγματοληπτικά σχέδια χωρίς πιθανότητα, 19
- δειγματοληπτικές μονάδες, 18
- δειγματοληψία, 22
- δειγματοχώρος ή δειγματικός χώρος, 134
- διαδικασία «post hoc tests», 464
- διακύμανση της μεταβλητής
  - » ανάμεσα στις ομάδες (between group variance), 245
  - » μέσα στις ομάδες (within group variance), 245
- διάμεσος (median), 74
- διασπορά ή διακύμανση (variance), 98
- διασταυρωμένη επικύρωση (cross-validation), 349

- διάστημα εμπιστοσύνης αριθμητικού μέσου όρου, 204
- διάστημα εμπιστοσύνης εκτιμητή, 204
- διάστημα εμπιστοσύνης της διαφοράς δύο αριθμητικών μέσων όρων, 215
- διάφορες εκφράσεις συχνότητας, 357
- διαχείριση ποιοτικών δεδομένων, 353
- διερευνητική ανάλυση δεδομένων (exploratory data analysis, eda), 25
- δίπλευρος έλεγχος, 229, 240
- διωνυμική κατανομή (binomial distribution), 159
- διωνυμική λογιστική παλινδρόμηση (binomial logistic regression), 388
- δοκιμασία Friedman (Friedman test), 468
- » Kolmogorov-Smirnov (K-S) για ένα δείγμα, 454
  - » Kruskal-Wallis H (Kruskal-Wallis H Test), 459
  - » Mann-Whitney U (Wilcoxon Rank-Sum Test), 444
  - » Wilcoxon (Wilcoxon Signed-Rank Test), 436
- δύναμη ελέγχου της δοκιμασίας, 276
- εγκυρότητα (validity), 422
- » (validity) ερωτηματολογίου, 423
  - » κατασκευής του ερωτηματολογίου, 425
  - » του περιεχομένου, 423
- εκτιμητική, 202
- ελάχιστα τετράγωνα (Least squares method), 304
- έλεγχος ανάλυσης διακύμανσης για τη διαπίστωση της γραμμικότητας της εξίσωσης, 308
- έλεγχος ανεξαρτησίας (Chi-square test for independence/of association), 375
- έλεγχος καλής προσαρμογής (goodness of fit test), 364
- έλεγχος ομοιογένειας (Test of Homogeneity), 387
- έλεγχος υποθέσεων, 227
- » υποθέσεων (hypothesis testing), 201
- έλεγχος υποθέσεων αριθμητικού μέσου όρου, 229
- » υποθέσεων δύο πληθυσμών, 240
  - » υποθέσεων ενός πληθυσμού, 226
  - » υποθέσεων περισσότερων του ενός πληθυσμών, 240
  - » υποθέσεων περισσότερων των δύο πληθυσμών (ANOVA), 244
- έλεγχος υπόθεσης ANOVA κατά ένα παράγοντα, 245
- έλεγχος υπόθεσης κατά δύο παράγοντες (two-way ANOVA), 255
- έλεγχος υπόθεσης κατά έναν παράγοντα (one-way ANOVA), 245
- έλεγχος υπόθεσης των συντελεστών παλινδρόμησης, 306
- εμπειρική κατανομή συχνότητας, 41
- » κατανομή συχνοτήτων συνεχούς μεταβλητής, 46
- εναλλακτική υπόθεση, 227
- ενδοτεταρτημοριακό εύρος (inter-quartile range), 94
- έννοια και ορισμός της πιθανότητας, 138
- εξέταση τακτικών κατηγορικών μεταβλητών, 358
- εξαντλητικές μέθοδοι, 21
- εξαρτημένη μεταβλητή, 299
- εξωτερική εγκυρότητα, 425
- επακρίβεια (precision), 423
- επαρκής, 203
- επικεφαλίδες των στηλών πίνακα, 33
- επικρατούσα τιμή, 357
- επίπεδο εμπιστοσύνης, 204
- » σημαντικότητας ( $\alpha$ ), 204, 228
- εσωτερική εγκυρότητα, 425
- εύρος (range), 91
- εφαρμοσμένη στατιστική, 3
- θεώρημα του Bayes (Bayes' theorem or Bayes' law or Bayes' rule), 149
- θηκόγραμμα (box-plot), 29
- ιδιότητες αριθμητικού μέσου όρου, 55
- » εκτιμητή, 202
  - » τυπικής απόκλισης, 107



- ιεραρχική ταξινόμηση (Hierarchical clustering), 400  
 ιστόγραμμα (histogram), 32  
 καμπύλη του Lorenz (Lorenz curve), 113  
 κανονική και τυπική κανονική κατανομή (normal and standard normal distribution), 178  
 κατανομή Poisson (Poisson distribution), 167  
 κατανομή δειγματοληψίας, 202  
 κατηγορικές μεταβλητές, 18, 353  
 κατηγορικές ονομαστικές μεταβλητές, 357  
 κλίση (slope), 302  
 κρίσιμη περιοχή, 228  
 κριτήριο του Akaike (Akaike Information Criterion, AIC), 327  
 κυρίως σώμα πίνακα, 33  
 κύρτωση (kurtosis), 124  
 Μ-εκτιμητές (M-estimators), 74  
 μερική συσχέτιση, 292  
 μέση απόλυτη απόκλιση (mean absolute deviation), 95  
 μέσοι όροι, 51  
 μεταβλητές, 17  
 μετάθεση - διάταξη (permutation - ordered arrangements), 151  
 μέτρα αξιολόγησης της προσαρμοσμένης εξίσωσης, 311  
 μέτρα ασυμμετρίας και κύρτωσης, 117, 357  
 μέτρα διασποράς, 91, 357  
 μέτρα κεντρικής τάσης, 51  
   » κεντρικής τάσης (θέσης), 357  
   » κεντρικής τάσης ή μέτρα θέσης, 51  
 μη-γραμμική παλινδρόμηση (nonlinear regression, NLR), 299, 335  
 μηδενική υπόθεση, 227  
 μη-παραμετρικοί έλεγχοι, 435  
 μονογραφική μέθοδος, 23  
 μονόπλευρος έλεγχος, 232  
 μπεϋζιανό κριτήριο (Bayesian Information Criterion, BIC), 327  
 ομοιογένεια διακύμανσης, 303  
 οπτική αναπαράσταση των συνόλων και πράξεων μεταξύ τους, 137  
 παλινδρόμηση (regression), 279, 298  
 παλινδρόμηση Lasso (Lasso regression), 349  
 παλινδρόμηση Ridge (Ridge regression analysis, RRA), 342  
 παραγοντική ανάλυση μικτών μεταβλητών (Factorial Analysis of Mixed Data, FAMD), 410  
 παραγοντικός έλεγχος υπόθεσης (factorial ANOVA), 275  
 πείραμα, 134  
 περιγραφικά στατιστικά κατάλληλα για κατηγορικές, 355  
 περιγραφική στατιστική, 201  
 πιθανότητα υπό συνθήκη, 146  
 πιθανότητα υπό συνθήκη (δεσμευμένη πιθανότητα) (conditional probability), 146  
 πιθανότητα υπό συνθήκη (δεσμευμένη πιθανότητα), 146  
 πίνακας συνδιακύμανσης, 281  
 πολλαπλασιαστική αρχή (multiplicative principle or product rule), 155  
 πολλαπλή γραμμική παλινδρόμηση (multiple linear regression, MLR), 299, 324  
 πολυσυγγραμμικότητα (multicollinearity), 325  
 ποσοσημόρια (quantiles), 84  
 ποσοτικές μεταβλητές, 18, 353  
 προ-επεξεργασία (pre-processing process) δεδομένων, 25  
 προσαρμοσμένος (διορθωμένος) συντελεστής προσδιορισμού ( $R_{adj}^2$ ), 326  
 ρίζα του μέσου τετραγωνικού σφάλματος (root mean square error, RMSE), 311, 326  
 σταθμισμένος αριθμητικός μέσος όρος (weighted arithmetic mean), 59  
 στατιστικά ανεξάρτητα, 144  
 στατιστικά στοιχεία – προ-επεξεργασία, 21  
 στατιστικές μονάδες, 16  
 στατιστικός πληθυσμός, 16  
 στικτό διάγραμμα, 300  
 στοιχεία επαγωγικής στατιστικής, 201  
 στοιχεία πιθανοτήτων, 133  
 στοχαστική ή στατιστική ανεξαρτησία γεγονότων (statistically independent events), 144

- συνάρτηση πιθανότητας (probability mass function, PMF), 158
- συνάρτηση πυκνότητας - πιθανότητας (probability density function, PDF), 158
- συνδιακύμανση (ή συνδιασπορά), (covariance), 280
- συνδυασμοί με επανάθεση, 154
- » χωρίς επανάθεση, 154
- συνδυασμός (combination – unordered arrangements), 154
- συνδυαστική ανάλυση, 151
- συνεπής, 203
- συνεχής καταγραφές, 21, 22
- συνεχής μεταβλητή, 157
- συντελεστές παλινδρόμησης, 300
- συντελεστής  $\tau$  του Kendall (Kendall's Tau), 415
- συντελεστής Cramer's  $V$ , 412
- συντελεστής κύμανσης τεταρτημορίου (coefficient of quartile variation), 111
- » μεταβλητότητας ή κύμανσης (coefficient of variation), 109
  - » προσδιορισμού ( $R^2$ ), 312
  - » συσχέτισης του Pearson ( $r$ ), 285
  - » του Spearman,  $\rho$  (Spearman's Rank Correlation,  $\rho$ ), 415
  - »  $\phi$  (Phi coefficient,  $\phi$ ), 412
- συσχέτιση (correlation), 279, 285
- σχέση μεταξύ αριθμητικού μέσου – διαμέσου - επικρατούσας τιμής, 83
- σχέση μεταξύ αριθμητικού, γεωμετρικού, αρμονικού και τετραγωνικού μέσου, 73
- σχετικές αθροιστικές συχνότητες ( $F_i$ ), 42
- σχετικές συχνότητες ( $p_i$ ), 42
- τεταρτημόρια (quartiles), 85
- τετραγωνικός μέσος (quadratic mean), 71
- τεχνική K-modes (K-modes), 397
- τιμές ( $x_i$ ), 41
- τυπική απόκλιση (standard deviation), 104
- » κανονική κατανομή, 180
- τυπικό σφάλμα εκτίμησης θεωρητικών τιμών (SEE), 326
- τυπικό σφάλμα εκτίμησης θεωρητικών τιμών (standard error of estimate, SEE), 312
- τύπος ή επικρατούσα τιμή (mode), 80
- τυχαία ή στοχαστική μεταβλητή (random or stochastic variable), 157
- τυχαία πειράματα ή πειράματα τύχης, 134
- τυχαίες ή στοχαστικές μεταβλητές, 17
- υπολογισμός ασυμμετρίας και κύρτωσης με τη γλώσσα προγραμματισμού R, 129
- » διαμέσου σε μη ομαδοποιημένα στοιχεία, 75
  - » διαμέσου σε ομαδοποιημένα στοιχεία, 78
  - » επικρατούσας τιμής σε ομαδοποιημένα στοιχεία, 82
  - » επικρατούσας τιμής σε μη ομαδοποιημένα στοιχεία, 80
  - » του δεύτερου συντελεστή ασυμμετρίας του Pearson ( $Sk_2$ ), 121
  - » του προσαρμοσμένου συντελεστή ασυμμετρίας του Fisher-Pearson ( $G_1$ ), 121
  - » του πρώτου συντελεστή ασυμμετρίας του Pearson ( $Sk_1$ ), 120
  - » του συντελεστή ασυμμετρίας με το στατιστικό πακέτο SPSS, 122
- υποσημειώσεις, 33
- φυλλογράφημα (Stem-and-Leaf Plot), 31
- F-κατανομή (F-distribution or Snedecor's F-distribution or Fisher-Snedecor distribution), 198
- pairwise comparisons, 473
- t-κατανομή ή Student – κατανομή (t-distribution or Student-distribution), 194
- $\chi^2$ -κατανομή (chi-squared distribution), 192