

**ΑΡΙΘΜΗΤΙΚΗ ΑΝΑΛΥΣΗ**  
**(ΣΗΜΕΙΩΣΕΙΣ ΜΕΤΑΠΤΥΧΙΑΚΟΥ ΜΑΘΗΜΑΤΟΣ)**

**ΒΑΣΙΛΕΙΟΣ Α. ΔΟΥΓΑΛΗΣ**

**Μαθηματικό Τμήμα Πανεπιστημίου Αθηνών**

**ΑΘΗΝΑ 1995**

## ΠΕΡΙΕΧΟΜΕΝΑ

Πρόλογος

Συνοπτική βιβλιογραφία

### 1. Αριθμητική Γραμμική Άλγεβρα

- 1.1 Απαλοιφή Gauss
- 1.2 Δείκτης κατάστασης πίνακα
- 1.3 Σφάλματα ετροχύλευσης στην απαλοιφή Gauss
- 1.4 Η ανάλυση Cholesky για συμμετρικούς, θετικά ορισμένους πίνακες
- 1.5 Μέθοδοι ελαχιστοποίησης
- 1.6 Η μέθοδος των ευζυγών κλίσεων

### 2. Αριθμητική λύση μη γραμμικών ευστημάτων

- 2.1 Παραγωγίσιμες συναρτήσεις στον  $\mathbb{R}^n$
- 2.2 Τοπικά θεώρηματα εύγκλισης. Το θεώρημα της ευστολής
- 2.3 Μέθοδος του Νεύτωνα: Τοπική εύγκλιση και ταχύτητα εύγκλισης
- 2.4 Το θεώρημα του Kantorovich για την εύγκλιση της μεθόδου του Νεύτωνα

### 3. Αριθμητική λύση συνήθων διαφορικών εξισώσεων

- 3.1 Πρόβλημα αρχικών τιμών. Η μέθοδος του Euler
- 3.2 Μέθοδοι Runge-Kutta
- 3.3 Πολυβηματικές μέθοδοι
- 3.4 Άκαμπα προβλήματα. Απόλυτη ευστάθεια και γενικεύσεις της

#### 4. Παρεμβολή και προσέγγιση

- 4.1 Παρεμβολή με πολυώνυμο Lagrange
- 4.2 Παρεμβολή και προσέγγιση με τμηματικά γραμμικές συναρτήσεις
- 4.3 Παρεμβολή με τμηματικά κυβικές συναρτήσεις Hermite
- 4.4 Παρεμβολή με κυβικές splines

## ΠΡΟΛΟΓΟΣ

Οι σημειώσεις αυτές γράφτηκαν και χρησιμοποιήθηκαν ως κύριο διδακτικό βοήθημα για το μεταπτυχιακό μάθημα 350 "Αριθμητική Ανάλυση" που δίδαξα κατά το χειμερινό εξάμηνο του ακαδημαϊκού έτους 1986-7 στο Μαθηματικό Τμήμα του Πανεπιστημίου Κρήτης σε μεταπτυχιακούς και προχωρημένους προπτυχιακούς φοιτητές. Σκοπός του μαθήματος ήταν να εκθέσει το ακροατήριο, μέσα στη χρονική διάρκεια ενός διδακτικού εξαμήνου, σε αποδείξεις ορισμένων θεωρημάτων κεντρικής σημασίας σε μερικές βασικές περιοχές της κλασικής Αριθμητικής Ανάλυσης, όπως αριθμητική γραμμική και μη γραμμική άλγεβρα, αριθμητική λύση συστημάτων διαφορικών εξισώσεων και θεωρία παρεμβολής και προσέγγισης συναρτήσεων. (Μια εισαγωγή σε θέματα αριθμητικής λύσης μερικών διαφορικών εξισώσεων γίνεται σε άλλο μεταπτυχιακό μάθημα του Μαθηματικού Τμήματος).

Η επιλογή των θεμάτων από την ευρύτερη περιοχή της Αριθμητικής Ανάλυσης, που είναι δυνατόν να διδαχθούν σε ένα τέτοιο μάθημα μέσα σε ένα εξάμηνο, είναι προφανές ότι είναι σε ένα βαθμό αυθαίρετη και αυτανακλά πιθανότατα τις προσωπικές προτιμήσεις και τα ενδιαφέροντα του διδάσκοντος. Δεν έχουν λοιπόν οι σημειώσεις αυτές αξιώσεις πληρότητας.

Στο κεφάλαιο 1 (Αριθμητική Γραμμική Άλγεβρα) εξετάζεται η αριθμητική λύση συστημάτων γραμμικών εξισώσεων. Στο πρώτο μέρος του το κεντρικό αποτέλεσμα είναι η αντίστροφη ανάλυση του Wilkinson για την μελέτη της επίρραξης των εσφαλμάτων ετροχύλωσης στην απαλοιφή Gauss. Ακολουθούν την απόδειξη που δίνεται στο κλασικό βιβλίο των Forsythe και Moler [1.2] (βλ. Συνοπτική Βιβλιογραφία). Το δεύτερο μέρος του κεφαλαίου αφορά την αριθμητική λύση γραμμικών συστημάτων με πραγματικούς, συμμετρικούς και θετικά ορισμένους πίνακες με ανάλυση Cholesky και, κατά κύριο λόγο, με μεθόδους ελαχιστοποίησης (μεθόδους καθόδου μεγίστης κλίσεως και ευζυχών κλίσεων). Ακολουθούν την ανάπτυξη του θέματος που κάνουν οι Golub και Van Loan, [1.4], συμπληρώνοντας τις αποδείξεις εύκλειους και φράγματος του εσφαλματος για τις μεθόδους ελαχιστοποίησης. Λόγω ελλείψεως χρόνου δεν εξετάστηκαν οι κλασικές επαναληπτικές μέθοδοι, καθώς και μέθοδοι για



το γραμμικό πρόβλημα ελαχίστων τετραγώνων και για το πρόβλημα ιδιοτιμών. (Μερικά θέματα που παρέλειψα είτε καλύπτονται σε άλλα προπτυχιακά μαθήματα του προγράμματος σπουδών του Μαθηματικού Τμήματος είτε παρουσιάσθηκαν από τους φοιτητές σε μορφή σεμιναρίου στο τέλος του εξαμήνου).

Στο Κεφάλαιο 2 (Αριθμητική Λύση μη Γραμμικών Συστημάτων) μετά από μία εισαγωγή στον διαφορικό λογισμό πολλών μεταβλητών και στα τοπικά θεωρήματα εύχλισης, αποδεικνύεται το θεώρημα του Kantorovich για την εύχλιση της μεθόδου του Νεύτωνα. Ακολουθήσα το κλασικό βιβλίο [2.3] των Ortega και Rheinboldt καθώς και το βιβλίο [0.5] του Ortega. Δεν υπήρχε δυστυχώς χρόνος ε' αυτό το μάθημα για να καλυφθούν σπουδαιότερα θέματα όπως η θεωρία των μεθόδων "του τύπου του Νεύτωνα" ή το μη γραμμικό πρόβλημα ελαχίστων τετραγώνων καθώς και γενικές μέθοδοι βελτιστοποίησης.

Το Κεφάλαιο 3 είναι αφιερωμένο στην αριθμητική λύση του προβλήματος αρχικών τιμών για συστήματα συνήθων διαφορικών εξισώσεων πρώτης τάξης με μεθόδους Runge-Kutta και πολυβηματικές μεθόδους. Για τις μεθόδους Runge-Kutta, με βάση την διατριβή του Crouzeix (βλ. παρ. 3.2), αποδεικνύεται ένα αποτέλεσμα των Butcher-Crouzeix που δίνει ικανές συνθήκες για εύχλιση με ορισμένη τάξη ακρίβειας στην περίπτωση μη ακάμπτων συστημάτων. Στις πολυβηματικές μεθόδους, για το βασικό αποτέλεσμα του Dahlquist, ότι δηλ. η εύχλιση είναι ισοδύναμη με συνέπεια και ευστάθεια, ακολουθήσα το κλασικό βιβλίο [3.4] του Henrici. Στο τέλος του κεφαλαίου γίνεται μία εισαγωγή στο πρόβλημα της αριθμητικής λύσης ακάμπτων συστημάτων και σε θέματα απόλυτης ευστάθειας και των γενικεύσεών της σε μη γραμμικά συστήματα. Δεν υπήρχε καιρός για να εξετασθούν θέματα όπως συνοριακά προβλήματα δύο ημερών, ειδικές μέθοδοι για διαφορικές εξισώσεις ανώτερης τάξης κ.ά.

Από την τεράστια και σημαντική περιοχή της θεωρίας Προσέγγισης περιορίστηκα, λόγω ελλείψεως χρόνου, σε ένα μόνο πρόβλημα δηλ., στην προσέγγιση συναρτήσεων μίας μεταβλητής με παρεμβολή με πολυώνυμα Lagrange και με τμηματικά πολυωνυμικές συναρτήσεις όπως

τμηματικά γραμμικές συναρτήσεις, τμηματικά κυβικές συναρτήσεις Hermite και κυβικές splines. Ακολούθησα, κατά κύριο λόγο, τα βιβλία των De Boor [4.2] και Schultz [4.10]. Από τα πολλά και σημαντικά θέματα της θεωρίας προσέγγισης που δεν έγινε δυνατόν να μελετηθούν, η αριθμητική ολοκλήρωση και η ομοιόμορφη προσέγγιση καλύπτονται σε προπτυχιακό μάθημα του Μαθηματικού Τμήματος.

Κατά το γράψιμο των σημειώσεων αυτών αντιμετώπισα το γνωστό πρόβλημα της μεταγλώττισης στα ελληνικά της ξένης ορολογίας. Ζητώ την επιείκεια του αναγνώστη αν οι μεταφράσεις μου ορισμένων όρων δεν ταιριάζουν στο γλωσσικό του αισθητήριο ή αν από άγνοια δεν υιοθέτησα ήδη καθιερωμένη στα ελληνικά ορολογία.

Θα ήθελα να ευχαριστήσω τους φοιτητές που παρακολούθησαν το μάθημα δίνοντάς μου έτσι την ευκαιρία να το διδάξω και να γράψω αυτές τις σημειώσεις. Ιδιαίτερα ευχαριστώ τον κ. Γεώργιο Ζουράρη που εντόπισε και διόρθωσε πολλά ουσιαστικά και τυπογραφικά λάθη στο κείμενο. Για ό,τι εφέλαμα απομένουν είμαι φυσικά ο ίδιος υπεύθυνος. Εκφράζω επίσης τις ευχαριστίες μου στο συνάδελφο κ. Γ. Ακριβή με τον οποίο συζητούσα διαρκώς πολλά θέματα σχετικά με το μάθημα και τις σημειώσεις.

Ιδιαίτερα επίσης ευχαριστώ την δ. Μαρία Σταυρακάκη που με ακρίβεια, ταχύτητα και επιμέλεια έγραψε το μεγαλύτερο μέρος αυτών των σημειώσεων στον μικροϋπολογιστή της. Μέρος του Κεφαλαίου 4 χράφηκε από την κ. Αίλινα Ζαχαριουδάκη την οποία επίσης ευχαριστώ. Τέλος θα ήθελα να ευχαριστήσω το Ινστιτούτο Υπολογιστικών Μαθηματικών του Ερευνητικού Κέντρου Κρήτης για οικονομική υποστήριξη στην δακτυλογράφηση, φωτοτύπηση και έκδοση των σημειώσεων στην παρούσα τους μορφή.

Ηράκλειο, Ιούνιος 1987

Β.Α. Δουχαλής

## ΣΥΝΟΠΤΙΚΗ ΒΙΒΛΙΟΓΡΑΦΙΑ

### 0. Γενικά βιβλία

- 0.1 R. Bulirsch and J. Stoer, "An introduction to numerical analysis", Springer-Verlag, Berlin-Heidelberg-New York 1980.
- 0.2 S.D. Conte and C. de Boor, "Elementary numerical analysis: an algorithmic approach", 3<sup>d</sup> ed, McGraw Hill, New York 1980.
- 0.3 G.E. Forsythe, M.A. Malcolm and C.B. Moler, "Computer methods for mathematical computation", Prentice-Hall, Englewood Cliffs N.J. 1977.
- 0.4 E. Isaacson and H.B. Keller, "Analysis of numerical methods", Wiley, New York 1966.
- 0.5 J.M. Ortega, "Numerical analysis: a second course", Academic Press, New York 1972.
- 0.6 B. Wendroff, "Theoretical numerical analysis", Academic Press, New York 1966.
- 0.7 Ν. Αποστολάτος, "Αριθμητική Ανάλυση I, II", Αθήνα 1971.
- 0.8 Α. Μπακόπουλος, "Αριθμητική Ανάλυση", Αθήνα 1981.
- 0.9 Α. Χατζηδόμος, "Αριθμητική Ανάλυση I, II", Ιωάννινα 1978, 1979.
- 0.10 Η. Χαύτης, "Αριθμητικές μέθοδοι, προγραμματισμός και ανάλυση, μέρος 1", Θεσσαλονίκη 1985.

### 1. Αριθμητική Γραμμική Άλγεβρα

- 1.1 D.K. Faddeev and V.N. Faddeeva, "Computational methods of linear algebra", Freeman, San Francisco 1963.

- 1.2 G.E. Forsythe and C.B. Moler, "Computer solution of linear algebraic systems", Prentice-Hall, Englewood Cliffs N.J. 1967.
- 1.3 A. George and J.W. Liu, "Computer solution of large, sparse positive definite systems", Prentice-Hall, Englewood Cliffs N.J. 1981.
- 1.4 G.H. Golub and C.F. Van Loan, "Matrix computations", Johns Hopkins U. Press, Baltimore 1983.
- 1.5 C.L. Lawson and R.J. Hanson, "Solving least squares problems", Prentice-Hall, Englewood Cliffs N.J. 1974.
- 1.6 B.N. Parlett, "The symmetric eigenvalue problem", Prentice-Hall, Englewood Cliffs N.J. 1980.
- 1.7 G.W. Stewart, "Introduction to matrix computations", Academic Press, New York 1973.
- 1.8 R.S. Varga, "Matrix iterative analysis", Prentice-Hall, Englewood Cliffs N.J. 1962.
- 1.9 J.H. Wilkinson, "Rounding errors in algebraic processes", Prentice-Hall, Englewood Cliffs N.J. 1963.
- 1.10 J.H. Wilkinson, "The algebraic eigenvalue problem", Clarendon Press, Oxford 1965.
- 1.11 D.M. Young, "Iterative solution of large linear systems", Academic Press, New York 1971.

## 2. Αριθμητική λύση μη γραμμικών συστημάτων

- 2.1 J.E. Dennis and R.B. Schnabel, "Numerical methods for unconstrained optimization and nonlinear equations", Prentice-Hall, Englewood Cliffs N.J. 1983.
- 2.2 P.E. Gill, W. Murray and M.H. Wright, "Practical optimization", Academic Press, London 1981.
- 2.3 J.M. Ortega and W.C. Rheinboldt, "Iterative solution of nonlinear equations in several variables", Academic Press, New York 1970.
- 2.4 W.C. Rheinboldt, "Methods for solving systems of nonlinear equations", SIAM, Philadelphia 1974.

## 3. Αριθμητική λύση συνήθων διαφορικών εξισώσεων

- 3.1 K. Dekker and J.G. Verwer, "Stability of Runge-Kutta methods for stiff nonlinear differential equations", North-Holland, Amsterdam 1984.
- 3.2 C.W. Gear, "Numerical initial value problems in ordinary differential equations", Prentice-Hall, Englewood Cliffs N.J. 1971.
- 3.3 R.D. Grigorieff, "Numerik gewöhnlicher Differentialgleichungen", Bd.I,II, B.G. Teubner, Stuttgart, Bd.I: 1972, Bd.II: 1977.
- 3.4 P.Henrici, "Discrete variable methods in ordinary differential equations", Wiley, New York 1962.
- 3.5 H.B. Keller, "Numerical methods for two-point boundary value problems", Blaisdell, Waltham Mass. 1968.

- 3.6 J.D. Lambert, "Computational methods in ordinary differential equations", Wiley, London 1973.
- 3.7 L.F. Shampine and M.K. Gordon, "Computer solution of ordinary differential equations", Freeman, San Francisco 1975.
- 3.8 M. Crouzeix et A.L. Mignot, "Analyse numerique des equations differentielles", Masson, Paris 1984.

#### 4. Παρεμβολή, προέχχιση και αριθμητική ολοκλήρωση

- 4.1 N.I. Achieser, "Theory of approximation", English Translation, Ungar, New York 1956.
- 4.2 C. de Boor, "A practical guide to splines", Springer-Verlag, New York 1978.
- 4.3 E.W. Cheney, "Introduction to approximation theory", McGraw-Hill, New York 1966.
- 4.4 P.J. Davis, "Interpolation and approximation", Blaisdell, Waltham Mass. 1963.
- 4.5 P.J. Davis and P. Rabinowitz, "Methods of numerical integration", Academic Press, New York 1975.
- 4.6 I.P. Natanson, "Constructive function theory", vols 1-3, (Transl. from Russian), Ungar, New York 1964-5.
- 4.7 M.J.D. Powell, "Approximation theory and methods", Cambridge U.P., Cambridge 1981.
- 4.8 J.R. Rice, "The approximation of functions", (2 vols), Addison-Wesley, Reading Mass., vol.1 1964, vol.2 1969.

- 4.9 T.J. Rivlin, "An introduction to the approximation of functions", Blaisdell, Waltham Mass. 1969.
- 4.10 M.H. Schultz, "Spline analysis", Prentice-Hall, Englewood Cliffs N.J. 1973.

#### 5. Σημειώσεις μαθημάτων Πανεπιστημίου Κρήτης

- 5.1 Γ.Δ. Ακρίβης, "Θεωρία προεσχίσεως και αριθμητική ολοκλήρωση", χειρόγραφες σημειώσεις, Χειμ. Εξ. 1985-6.
- 5.2 Γ.Δ. Ακρίβης, "Εισαγωγή στην Αριθμητική Ανάλυση", δακτυλογραφημένες σημειώσεις, Ηράκλειο 1986.
- 5.3 Β.Α. Δουχαλής, "Διδακτικές σημειώσεις για το μάθημα Αριθμητική Ανάλυση II", χειρόγραφες σημειώσεις, Χειμ. Εξ. 1983-4.
- 5.4 Β.Α. Δουχαλής, "Αριθμητική λύση γραμμικών συστημάτων στον υπολογιστή", δακτυλογραφημένες σημειώσεις, Ηράκλειο 1986.
- 5.5 Β.Α. Δουχαλής, "Διδακτικές σημειώσεις για το Μάθημα Αριθμητική Ανάλυση I", χειρόγραφες σημειώσεις, Εαρ. Εξ. 1983-4.
- 5.6 Β.Α. Δουχαλής, Διδακτικές σημειώσεις για το μεταπτυχιακό μάθημα 351 "Αριθμητική λύση Περιόδων Διαφορικών Εξισώσεων" (Μέθοδοι Galerkin/πεπερασμένων στοιχείων), χειρόγραφες σημειώσεις, ετά Αγγλικά, Εαρ. εξ. 1984-5.

1. ΑΡΙΘΗΤΙΚΗ ΓΡΑΜΜΙΚΗ ΑΛΓΕΒΡΑ

## 1.1 ΑΠΑΛΟΙΦΗ GAUSS

Έστω  $A=(a_{ij}), 1 \leq i, j \leq n$ , πραγματικός  $n \times n$  πίνακας (θα συμβολίζουμε  $A \in \mathbb{R}^{n \times n}$ ),  $b=(b_1, \dots, b_n)^T \in \mathbb{R}^n$ . Υποθέτουμε ότι ο  $A$  είναι αντιστρέψιμος' έστω  $x=(x_1, \dots, x_n)^T$  η λύση του γραμμικού συστήματος

$$(1) \quad Ax=b.$$

θα αναλύσουμε μία βασική μέθοδο για την επίλυση του (1), την απαλοιφή Gauss. Έστω  $A^{(1)}=A, b^{(1)}=b$ . Στο πρώτο βήμα ας υποθέσουμε

(1)

ότι  $a_{11} \neq 0$ . (Αλλιώς, με εναλλαγή δύο γραμμών μπορούμε να φέρουμε ένα μη μηδενικό στοιχείο - φερ' ειπείν εκείνο με τον μικρότερο δείκτη

(1)

γραμμής - της πρώτης στήλης στην θέση (1,1)). Το στοιχείο  $a_{11} \neq 0$  λέγεται αδηγός του πρώτου βήματος. Ορίζοντας τώρα τους πολλαπλασιαστές

(1) (1)

$$m_{i1} = a_{i1} / a_{11}, \quad 2 \leq i \leq n,$$

πολλαπλασιάζοντας την πρώτη εξίσωση του (1) επί  $m_{i1}$  και αφαιρώντας από την  $i$ -ετή, παίρνουμε το ισοδύναμο σύστημα

$$A^{(2)}x = b^{(2)},$$

όπου

(2)

$$a_{ij} =$$

(1)

$$a_{ij}$$

αν  $i=1, 1 \leq j \leq n$

$$0$$

αν  $j=1, 2 \leq i \leq n$

(1) (1)

$$a_{ij} - m_{i1} a_{1j}$$

αν  $2 \leq i, j \leq n$



## 1.1.2

$$(2) \quad b_i = \begin{cases} (1) & \text{αν } i=1 \\ (1) & \text{αν } 2 \leq i \leq n. \\ b_i - m_{i1} b_1 \end{cases}$$

Προχωρούμε ανάλογα μετατρέποντας σε μηδενικά τα στοιχεία της δεύτερης στήλης του  $A^{(2)}$  κάτω από το διαγώνιο στοιχείο. Στο  $k$ -στό βήμα της μεθόδου έχουμε το ισοδύναμο προς το (1) σύστημα

$$A^{(k)} x = b^{(k)},$$

όπου τα  $A^{(k)}$ ,  $b^{(k)}$  είναι της μορφής

$$A^{(k)} = \begin{bmatrix} (k) & & & (k) \\ a_{11} & \dots & & a_{1n} \\ & \ddots & & \vdots \\ & & (k) & (k) \\ & 0 & a_{kk} & \dots & a_{kn} \\ & & \vdots & & \vdots \\ & & (k) & & (k) \\ & & a_{nk} & \dots & a_{nn} \end{bmatrix}, \quad b^{(k)} = \begin{bmatrix} (k) \\ b_1 \\ \vdots \\ (k) \\ b_k \\ \vdots \\ (k) \\ b_n \end{bmatrix}$$

Έστω τώρα ότι ο οδηγός  $a_{kk}$  είναι διάφορος του μηδενός. (Αν είναι ίσος με το μηδέν, με εναλλαγή της  $k$ -στής γραμμής του  $A^{(k)}$  με κάποια  $i$ -στή γραμμή ( $i > k$ ) μπορούμε να φέρουμε ένα μη μηδενικό -π.χ. εκείνο με το μικρότερο δυνατό  $i$  - στοιχείο στην θέση του οδηγού). Ορίζουμε τους πολλαπλασιαστές:

$$m_{ik} = a_{ik} / a_{kk}, \quad k+1 \leq i \leq n,$$

πολλαπλασιάζουμε την  $k$ -στή γραμμή επί  $m_{ik}$  και αφαιρούμε από την

$i$ -ετή. Προκύπτει έτσι το σύστημα

$$A^{(k+1)}x = b^{(k+1)},$$

όπου

$$a_{ij}^{(k+1)} = \begin{cases} a_{ij}^{(k)} & \text{αν } i \leq k \\ 0 & \text{αν } i > k \text{ και } j \leq k \\ a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} & \text{αν } k+1 \leq i, j \leq n, \end{cases}$$

$$b_i^{(k+1)} = \begin{cases} b_i^{(k)} & \text{αν } i \leq k \\ b_i^{(k)} - m_{ik} b_k^{(k)} & \text{αν } k+1 \leq i \leq n. \end{cases}$$

Στό τέλος αυτής της διαδικασίας μετά από  $n-1$  βήματα του αλγορίθμου, προκύπτει το σύστημα

$$(2) \quad A^{(n)}x = b^{(n)},$$

ισοδύναμο με το (1)· ο πίνακας  $A^{(n)}$  είναι ο άνω τριγωνικός αντίστρο-

$(n)$

ψιμος πίνακας με στοιχεία  $a_{ij}^{(n)}$ , όπου βέβαια

$$a_{ij}^{(n)} = \begin{cases} 0 & i > j \\ a_{ij}^{(i)} & i \leq j \end{cases}$$

$(n)$   $(i)$

Γιά το διάνυσμα  $b^{(n)}$  έχουμε  $b_i = b_i^{(i)}$ ,  $1 \leq i \leq n$ . Με τον υπολογισμό των

1.1.4

$A^{(n)}, b^{(n)}$ , ολοκληρώνεται η πρώτη φάση της απαλοιφής, η τριγωνοποίηση του (1). Κατά την δεύτερη φάση, την οπισθοδρόμηση, υπολογίζουμε τους αγνώστους με την σειρά  $x_n, x_{n-1}, \dots, x_1$  λύνοντας το άνω τριγωνικό σύστημα (2) με τις αναδρομικές σχέσεις:

$$\begin{cases} x_n = b_n^{(n)} / a_{nn}^{(n)} \\ x_i = (b_i^{(n)} - \sum_{j=i+1}^n a_{ij}^{(n)} x_j^{(n)}) / a_{ii}^{(n)}, \quad i=n-1, n-2, \dots, 1. \end{cases}$$

(Θα σημειωθεί ότι αν ο πίνακας  $A$  δεν είναι αντιστρέψιμος, τότε για

(k)

κάποιο  $k$ , όλα τα στοιχεία  $a_{ik}$ ,  $k \leq i \leq n$  θα είναι μηδέν. Παραλείπουμε τότε το  $k$ -στό βήμα της απαλοιφής και προχωρούμε στο  $k+1$ -στό. Δηλ. η τριγωνοποίηση ενός μη αντιστρέψιμου πίνακα είναι πάντα δυνατή' βέβαια

(i)

στον  $A^{(n)}$  τουλάχιστον ένα διαγώνιο στοιχείο  $a_{ii}$  θα είναι μηδέν. Το αν το σύστημα είναι συμβιβαστό ή αδύνατο εξαρτάται φυσικά και από το  $b$ ).

Η φάση της τριγωνοποίησης του  $A$  κατά την απαλοιφή είναι δυνατόν να ερμηνευθεί σε γλώσσα πινάκων με την λεγόμενη "ανάλυση LU", πολύ σημαντική για την θεωρία και τις εφαρμογές:

**ΘΕΩΡΗΜΑ 1.** Για κάθε  $A \in \mathbb{R}^{n \times n}$  υπάρχει  $n \times n$  πίνακας μεταθέσεως  $P$ , τέτοιος ώστε

$$(3) \quad PA = LU,$$

όπου ο  $n \times n$  πίνακας  $L$  είναι κάτω τριγωνικός με μονάδες στην διαγώνιο του και ο πίνακας  $U$  είναι άνω τριγωνικός.

Απόδειξη: Χρησιμοποιούμε τον συμβολισμό και την διαδικασία της απαλοιφής. Υποθέτουμε πρώτα ότι κατά την απαλοιφή δεν έγιναν καθόλου εναλλαγές γραμμών, δηλ. ότι για κάθε  $k$  στον  $A^{(k)}$  έχουμε είτε

(k)  $a_{kk} \neq 0$  είτε  $a_{ik} = 0, -k \leq i \leq n$ . Θεωρούμε τον πίνακα

$$M_1 = \begin{pmatrix} 1 & & & & \\ -m_{21} & 1 & & & 0 \\ \vdots & & \ddots & & \\ -m_{n1} & & & & 1 \end{pmatrix}$$

Από τον ορισμό του  $A^{(2)}$  προκύπτει (με  $A^{(1)}=A$ ) ότι

$$A^{(2)} = M_1 A^{(1)}$$

(1)

(Αν  $a_{i1} = 0, 1 \leq i \leq n$ , θέτουμε  $M_1=I$ ). Γενικά, στο  $k$ -στό βήμα έχουμε

$$A^{(k+1)} = M_k A^{(k)},$$

όπου ο  $n \times n$  πίνακας  $M_k$  έχει στοιχεία

$$(M_k)_{ij} = \begin{cases} 1 & \text{αν } i=j \\ -m_{ik} & \text{αν } k+1 \leq i \leq n, j=k \\ 0 & \text{αλλιώς} \end{cases}$$

(k)

(k)

αν  $a_{kk} \neq 0$ , ή  $M_k=I$  αν  $a_{ik} = 0, k \leq i \leq n$ . Συμπεραίνουμε ότι

$$A^{(n)} = M_{n-1} \dots M_1 A,$$



Έστω τώρα ότι ο  $A$  είναι αντιστρέψιμος και ότι υπολογίσθηκαν οι πίνακες  $P, L, U$  με την διαδικασία της απαλοιφής. Για να λύσουμε το σύστημα (1) πολλαπλασιάζουμε και τα δύο μέλη με  $P$  οπότε προκύπτει το ισοδύναμο σύστημα

$$(6) \quad LUx = Pb,$$

που λύνεται σε δύο βήματα: πρώτα υπολογίζουμε το ευδιάμεσο διάνυσμα  $y$  ως λύση του κάτω τριγωνικού συστήματος

$$(7) \quad Ly = Pb,$$

με τον προφανή αλγόριθμο που υπολογίζει πρώτα το  $y_1$  και γενικά το  $y_k$  συναρτήσει των  $y_i$ ,  $i < k$ . Κατόπιν υπολογίζουμε το  $x$  λύνοντας το άνω τριγωνικό σύστημα

$$(8) \quad Ux = y$$

με τον αλγόριθμο της οπισθοδρόμησης.

#### Παρατηρήσεις

1. Στην πράξη, η λύση του συστήματος (1) με απαλοιφή Gauss γίνεται σε δύο φάσεις, όπως περιγράψαμε πιο πάνω:

(α) Στην φάση της ανάλυσης  $PA=LU$ , κατά την οποία εργαζόμαστε μόνο με τον πίνακα  $A$ : υπολογίζουμε τα (ευδιάφερα) στοιχεία των πινάκων  $L$  και  $U$  με την διαδικασία της απαλοιφής και καταγράφουμε τις πληροφορίες εναλλαγής γραμμών που πιθανόν θα γίνουν σε ορισμένα βήματα, δηλ. την δράση του  $P$ . Η φάση αυτή απαιτεί  $n^3/3 + O(n)$  πράξεις (πράξη=πολλαπλασιασμός ή διαίρεση). Οι πολλαπλασιαστές  $m_{ij}$  (δηλ. τα στοιχεία του  $L$ ) καταχωρούνται κατά στήλες στις θέσεις μνήμης όπου ήταν αποθηκευμένα τα στοιχεία  $a_{ij}$ ,  $i > j$  του  $A$  που μετατρέπονται σε μηδενικά. Τα στοιχεία  $u_{ij}$  του  $U (= A^{(n)})$  γράφονται κατά γραμμές πάνω στις θέσεις των  $a_{ij}$ ,  $i \leq j$ . Οι πληροφορίες εναλλαγής γραμμών καταχω-



2. (α) Υπολογίστε το πλήθος των πολλαπλασιασμών και διαιρέσεων που απαιτούνται για την φάση της ανάλυσης  $PA=LU$  καθώς και το πλήθος των προεξαφαιρέσεων και του απαιτούμενου χώρου μνήμης.

(β) Υπολογίστε το πλήθος των πράξεων (πράξη=πολλαπλασιασμός ή διαίρεση) που απαιτούνται για την επίλυση ενός  $n \times n$  άνω ή κάτω τριγωνικού ευστήματος. Υπολογίστε επίσης το πλήθος των προεξαφαιρέσεων και της απαιτούμενης μνήμης.

(γ) Δείξτε ότι ο υπολογισμός του αντιστρόφου του  $A$  (βλ. παρατήρηση 2) απαιτεί  $n^3+O(n^2)$  πράξεις και  $2n^2+O(n)$  θέσεις μνήμης.

(δ) Στις εφαρμογές επάνω ενδιαφερόμαστε για τον  $A^{-1}$  καθεαυτόν. Συνήθως ενδιαφερόμαστε για την λύση  $x$  ενός ευστήματος  $Ax=b$  που γράφεται και σαν  $x=A^{-1}b$ . Για τον υπολογισμό της λύσης με απαλοιφή Gauss χρειαζόμαστε όπως είδαμε  $n^3/3+O(n^2)$  πράξεις και  $n^2+O(n)$  θέσεις μνήμης. Για τον υπολογισμό όμως του  $A^{-1}b$  χρειαζόμαστε περίπου τον τριπλάσιο αριθμό πράξεων και τον διπλάσιο χώρο μνήμης. Συνεπώς ουδέποτε υπολογίζουμε διανύσματα της μορφής  $A^{-1}y$ . Βέτομε  $x=A^{-1}y$  και λύσουμε το σύστημα  $Ax=y$ . Π.χ.: πώς υπολογίζουμε, όσο το δυνατό φθηνότερα από άποψη πλήθους πράξεων και χώρου μνήμης, τα διανύσματα  $A^{-5}y$ ,  $A^{-1}BA^{-1}y$ ,  $ABA^{-1}y$ , όπου  $A, B \in \mathbb{R}^{n \times n}$  ( $A$  αντιστρέψιμος),  $y \in \mathbb{R}^n$ .

(ε) Δείξτε ότι ο υπολογισμός της ορίζουσας ενός  $n \times n$  πίνακα απαιτεί  $n^3/3+O(n)$  πράξεις και μπορεί να γίνει εύκολα ως παραπροϊόν της ανάλυσης  $PA=LU$  (Συγκρίστε με τον αριθμό των πράξεων που απαιτούνται από τον ευθέτη ορισμό της ορίζουσας).

3. (α) Έστω η  $n$  μετάθεση  $i_k \mapsto k$ ,  $k=1, 2, \dots, n$ , των  $n$  πρώτων φυσικών και έστω  $P$  ο πίνακας μεταθέσεως που αντιστοιχεί στην  $\pi$ . Τότε η  $k$ -στή γραμμή του  $P$  είναι το διάνυσμα  $(e^{(i_k)})^T$  όπου  $e^{(i)} = \delta_{ij}$ ,  $1 \leq i, j \leq n$ .

(β) Δείξτε ότι  $PP^T=I$ , δηλ. ότι, για κάθε πίνακα μεταθέσεως,  $P^{-1}=P^T$ . Συνεπώς  $\det(P)=\pm 1$ . Βεβαιωθείτε ότι ο  $P^T$  παριστάνει την μετάθεση  $\pi^{-1}$ , δηλ. την  $k \mapsto i_k$ ,  $1 \leq k \leq n$ .

(γ) Έστω  $A \in \mathbb{R}^{n \times n}$ ,  $b \in \mathbb{R}^n$ . Βρείτε συναρτήσεις των  $a_{ij}, b_i$  και της  $n$  τα στοιχεία των  $PA, AP, P^T A, AP^T, PAP^T, Pb, P^T b$ .



(δ) Έστω  $P_i$  οι πίνακες μεταθέσεων που αντιστοιχούν στις μεταθέσεις  $\pi_i$ ,  $i=1,2$ . Δείξτε ότι το γινόμενο  $P_1 P_2$  είναι πίνακας μεταθέσεων. Σε ποιά μετάθεση αντιστοιχεί;

4. (α) Υποθέστε ότι ο  $A \in \mathbb{R}^{n \times n}$  είναι αντιστρέψιμος και ότι η ανάλυση LU του μπορεί να γίνει χωρίς εναλλαγές γραμμών, (δηλ. ότι  $P=I$  στην (3)). Τότε το ζευγάρι  $L, U$  (με τις γνωστές ιδιότητες) είναι μοναδικό. Η' άλλα λόγια, έστω ότι  $A$  αντιστρέψιμος και υποθέστε ότι

$$(i) A=LU$$

$$(ii) L \text{ κάτω τριγωνικός με } L_{ii}=1$$

$$(iii) U \text{ άνω τριγωνικός.}$$

Δείξτε τότε ότι οι  $L, U$  είναι μοναδικοί.

(β) Η μοναδικότητα των  $L, U$  (με τις παραπάνω προϋποθέσεις) μας επιτρέπει να τους υπολογίσουμε με κατασκευές διαφορετικές από εκείνη της απόδειξης του θεωρήματος 1. Π.χ. χρησιμοποιώντας την ιδιότητα  $A=LU$  στοιχείο προς στοιχείο, δείξτε ότι τα στοιχεία των  $L, U$  δίνονται από τους τύπους:

$$L_{ij} = (a_{ij} - \sum_{k=1}^{j-1} L_{ik} U_{kj}) / U_{ii}, \quad j < i$$

$$L_{ii} = 1$$

$$U_{ij} = a_{ij} - \sum_{k=1}^{i-1} L_{ik} U_{kj}, \quad j > i$$

όπου συμβολίζουμε  $\sum_{k=p}^q = 0$  αν  $q < p$ . Δώστε ένα ακριβή αλγόριθμο για τον υπολογισμό και αποθήκευση των  $L, U$ , γραμμή προς γραμμή. (Η κατασκευή αυτή των  $L$  και  $U$  λέγεται μέθοδος του Crout).

(γ) Υποθέστε ότι ισχύουν οι υποθέσεις του ερωτήματος (α) και ότι επιπλέον ο  $A$  είναι συμμετρικός. Δείξτε ότι  $U=BL^T$  όπου  $B$  ο διαγώνιος πίνακας με  $B_{ii}=U_{ii}$ . Δώστε αλγόριθμο ανάλογο με του του (β) για την κατασκευή των  $B, L$ .

## 1.1.11

5. Για  $A \in \mathbb{R}^{n \times n}$  ορίζουμε τις κύριες ορίζουσες  $\delta_i$  ως εξής:

$$\delta_1 = a_{11}, \quad \delta_2 = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix}, \quad \dots, \quad \delta_i = \begin{vmatrix} a_{11} & \dots & a_{1i} \\ \vdots & \ddots & \vdots \\ a_{i1} & \dots & a_{ii} \end{vmatrix}$$

Έστω ότι  $\delta_i \neq 0, 1 \leq i \leq n$ . Τότε δείξτε ότι το θεώρημα 1 ισχύει με  $P=1$ ,

δηλ. ότι υπάρχουν πίνακες  $L, U$  με τις ιδιότητες (i), (ii), (iii) της άσκησης 4(α). Επιπλέον το ζευγάρι  $L, U$  είναι μοναδικό. (Υπόδειξη: για

την ύπαρξη δείξτε με επαγωγή ότι  $a_{kk}^{(k)} \neq 0, 1 \leq k \leq n$ , χωρίς εναλλαγές

γραμμών.

6. Έστω ότι ο πίνακας  $A \in \mathbb{R}^{n \times n}$  είναι αντιστρέψιμος και έχει κυριαρχική διαγώνιο, κατά γραμμές, δηλ. ότι

$$|a_{ii}| \geq \sum_{j=1, j \neq i}^n |a_{ij}|, \quad 1 \leq i \leq n.$$

(α) Δείξτε ότι υπάρχει (μοναδικό) ζευγάρι  $L, U$  με τις ιδιότητες (i), (ii), (iii) της άσκησης 4(α). (Υπόδειξη: από τις υποθέσεις μας

(1)

$a_{11} = a_{11} \neq 0$ . θεωρούμε τον  $(n-1) \times (n-1)$  υποπίνακα  $\tilde{A}^{(2)}$  του  $A^{(2)}$  που προκύπτει αν αφαιρέσουμε την πρώτη γραμμή και την πρώτη ετήλη του  $A^{(2)}$ . Δείξτε ότι ο  $\tilde{A}^{(2)}$  είναι αντιστρέψιμος και έχει κυριαρχική

διαγώνιο— κατά γραμμές. Συμπέρασμα:  $a_{11}^{(2)} \neq 0$ , δηλ. δεν χρειάζεται

22

εναλλαγή γραμμών· για να υπολογίσουμε τον  $A^{(3)}$ . Η συνέχεια προφανής, με επαγωγή).

(β) Δείξτε ότι τα στοιχεία του  $U$  που προκύπτουν από την ανάλυση

$LU$  του  $A$  του μέρους (α) ικανοποιούν την ανισότητα  $|u_{ij}| \leq 2 \max_{i,j} |a_{ij}|$ .

(γ) Αν ο  $A$  είναι αντιστρέψιμος και έχει κυριαρχική διαγώνιο κατά ετήλες (προφανής ο ανάλογος ορισμός), τότε ισχύει πάλι η (α) (με ανάλογη υπόδειξη) και η εκτίμηση  $|L_{ij}| \leq 1$ .

## 1.2 ΔΕΙΚΤΗΣ ΚΑΤΑΣΤΑΣΗΣ ΠΙΝΑΚΑ

Αν ο  $A \in \mathbb{R}^{n \times n}$  είναι αντιστρέψιμος ο αλγόριθμος της απαλοιφής Gauss της παρ. 1.1 υπολογίζει την λύση του συστήματος  $Ax=b$  υπό την προϋπόθεση βέβαια ότι οι αριθμητικές πράξεις γίνονται ακριβώς. Στην πραγματικότητα όμως αυτό δεν είναι δυνατόν όπως ξέρουμε οι αριθμητικές πράξεις υπόκεινται σε εσφάλματα ετροχχύλευσης λόγω της πεπερασμένης ακρίβειας της αριθμητικής κάθε υπολογιστή. Τα εσφάλματα ετροχχύλευσης ενσωματώνονται και αποβαίνουν συχνά καταστρεφικά έτσι ώστε η υπολογιστική λύση  $\tilde{x}$  να απέχει πολύ από την θεωρητική λύση  $x$ . Ένας από τους κύριους σκοπούς μας στη συνέχεια θα είναι να αναλύσουμε την επίρροή αυτών των εσφαλμάτων στην απαλοιφή.

Θα αρχίσουμε εξετάζοντας ε' αυτήν την παράγραφο το θεωρητικό πρόβλημα της ευαισθησίας της θεωρητικής λύσης  $x$  του γραμμικού συστήματος  $Ax=b$  σε διαταραχές (μεταβολές) στα δεδομένα  $A$  και  $b$  του προβλήματος. Λέμε ότι ένα σύστημα έχει κακή κατάσταση (ill-conditioned) - ή ότι είναι "αεστάθης" - αν είναι πολύ ευαίσθητο σε διαταραχές, δηλ. αν "μικρές" διαταραχές των στοιχείων των  $A, b$  είναι δυνατόν να επιφέρουν "μεγάλες" μεταβολές στη λύση του. Το πρόβλημα αυτό είναι ανεξάρτητο από την μελέτη οποιουδήποτε αλγορίθμου για την αριθμητική λύση του συστήματος. Είναι προφανές όμως ότι και μόνο η (γενικά μη ακριβής) παράσταση των στοιχείων  $a_{ij}$  και  $b_i$  του υπολογιστή ευνοεί μία "διαταραχή" στα δεδομένα.

Επιπλέον, όπως θα δούμε, η υπολογιστική λύση  $\tilde{x}$  μπορεί να θεωρηθεί ως ακριβής λύση κάποιου "παραπλήσιου" προς το  $Ax=b$  συστήματος. Η κατάσταση (ευαισθησία σε διαταραχές) ενός γραμμικού συστήματος - που εξαρτάται από το μέγεθος του δείκτη κατάστασης του πίνακα  $A$  - είναι λοιπόν ένας σημαντικός παράγοντας που επηρεάζει την ακρίβεια της υπολογιστικής λύσης  $\tilde{x}$ . Ένας άλλος παράγοντας είναι η ευστάθεια (ή η αστάθεια) του συγκεκριμένου αλγορίθμου που υπολογίζει το  $\tilde{x}$ , πρόβλημα που θα εξετάσουμε στην επόμενη παράγραφο.

Αρχίζουμε με μία εύστοχη επανάληψη πάνω στις νόρμες (στάθμες) διανυσμάτων και πινάκων για περισσότερες λεπτομέρειες και αποδείξεις βλ. π.χ. [5.4, παρ. 3]. Μία (διανυσματική) νόρμα στον  $\mathbb{C}^n$  είναι μία

απεικόνιση  $\|\cdot\|: \mathbb{C}^n \rightarrow [0, \infty)$  με τις ιδιότητες

$$(i) \quad \|x\| \geq 0 \quad \forall x \in \mathbb{C}^n, \quad \|x\|=0 \Leftrightarrow x=0.$$

$$(ii) \quad \|\lambda x\| = |\lambda| \|x\| \quad \forall x \in \mathbb{C}^n, \lambda \in \mathbb{C}.$$

$$(iii) \quad \|x+y\| \leq \|x\| + \|y\| \quad \forall x, y \in \mathbb{C}^n.$$

Μερικές χρήσιμες νόρμες στην αριθμητική ανάλυση είναι οι

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i| \quad (\text{η λεγόμενη } l_\infty \text{ ή max(imum) νόρμα})$$

$$\|x\|_1 = \sum_{i=1}^n |x_i| \quad (l_1 \text{ νόρμα})$$

$$\|x\|_2 = \left( \sum_{i=1}^n |x_i|^2 \right)^{1/2} \quad (l_2 \text{ ή Ευκλείδεια νόρμα}).$$

Για την  $l_2$  νόρμα έχουμε ότι  $\|\cdot\|_2 = (\cdot, \cdot)_2$ , όπου  $(x, y)_2 = \sum_{i=1}^n x_i \bar{y}_i$ ,  $x, y \in \mathbb{C}^n$ , είναι το Ευκλείδειο εσωτερικό γινόμενο στον  $\mathbb{C}^n$ . Ένα θεμελιώδες αποτέλεσμα, απόρροια του ότι ο  $\mathbb{C}^n$  έχει πεπερασμένη διάσταση, είναι ότι οποιοδήποτε δύο νόρμες στον  $\mathbb{C}^n$  είναι ισοδύναμες (ή εσυγκρίσιμες). Δηλ. ότι, δεδομένων δύο νορμών  $\|\cdot\|_\alpha$  και  $\|\cdot\|_\beta$  στον  $\mathbb{C}^n$ , υπάρχουν θετικές σταθερές  $c_1, c_2$  (που εξαρτώνται γενικά από το  $n$ ), τέτοιες ώστε να ισχύει

$$(1) \quad c_1 \|x\|_\alpha \leq \|x\|_\beta \leq c_2 \|x\|_\alpha, \quad \forall x \in \mathbb{C}^n.$$

Λέμε ότι η ακολουθία  $\{x^k\}_{k=1}^\infty$  διανυσμάτων του  $\mathbb{C}^n$  συγκλίνει στο διάνυσμα  $x \in \mathbb{C}^n$  (γράφουμε  $x^k \rightarrow x, k \rightarrow \infty$ ) αν για κάποια νόρμα  $\|\cdot\|$  του  $\mathbb{C}^n$  ισχύει

$$\lim_{k \rightarrow \infty} \|x^k - x\| = 0$$

Από την (1) έπεται ότι η σύγκλιση είναι ανεξάρτητη της χρησιμοποιού-

μειξης νόρμας. Επίσης από την ισοδυναμία των νορμών  $\|\cdot\|$  και  $\|\cdot\|_\infty$

έπεται ότι  $x^k \rightarrow x, k \rightarrow \infty \Leftrightarrow \forall i, 1 \leq i \leq n, x_i^k \rightarrow x_i, k \rightarrow \infty$ .

Κάθε νόρμα  $\|\cdot\|$  του  $\mathbb{C}^n$  παράγει μία αντίστοιχη νόρμα ("φυσική" νόρμα, νόρμα "τελεστού") στον χώρο  $\mathbb{C}^{n \times n}$  των τετραγωνικών μιγαδικών πινάκων: για  $A \in \mathbb{C}^{n \times n}$  ορίζουμε την νόρμα  $\|A\|$  ως

$$(2) \quad \|A\| = \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \|Ax\| / \|x\| \quad (= \sup_{\substack{x \in \mathbb{C}^n \\ \|x\| \leq 1}} \|Ax\| = \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|=1}} \|Ax\|).$$

Εύκολα προκύπτει ότι αν  $A, B \in \mathbb{C}^{n \times n}$  τότε

- (α)  $\|A\| \geq 0, \|A\|=0 \Leftrightarrow A=0,$
- (β)  $\|\lambda A\| = |\lambda| \|A\| \quad \forall \lambda \in \mathbb{C},$
- (γ)  $\|A+B\| \leq \|A\| + \|B\|,$
- (δ)  $\|AB\| \leq \|A\| \|B\|.$

Οι νόρμες πινάκων που παράγονται από τις διανυσματικές νόρμες  $\|\cdot\|_\infty, \|\cdot\|_1, \|\cdot\|_2$  δίνονται, αντίστοιχα, συνάρτησει του  $A=(a_{ij}) \in \mathbb{C}^{n \times n}$  από τους τύπους:

$$\|A\|_\infty = \max_{1 \leq i \leq n} \left( \sum_{j=1}^n |a_{ij}| \right),$$

$$\|A\|_1 = \max_{1 \leq j \leq n} \left( \sum_{i=1}^n |a_{ij}| \right),$$

$$\|A\|_2 = \max_{1 \leq i \leq n} [\lambda_i(A^*A)]^{1/2},$$

όπου  $A^*$  είναι ο ανάστροφος συζυγής του  $A$ , δηλ. όπου  $(A^*)_{ij} = \bar{a}_{ji}$ , και όπου με  $\lambda_i(B)$  συμβολίζουμε τις ιδιοτιμές του  $B \in \mathbb{C}^{n \times n}$ . Αν ο πίνακας  $A$  είναι αυτοσυζυγής (ερμιτιανός), δηλ. αν  $A=A^*$ , τότε ο τελευταίος

τύπος απλουστεύεται σε  $\|A\|_2 = \max_{1 \leq i \leq n} |\lambda_i(A)|$ .

Από τον ορισμό (2) έπεται ότι οποιεσδήποτε δύο νόρμες πινάκων στον  $\mathbb{C}^{n \times n}$  είναι ισοδύναμες. Σύγκλιση ακολουθίας πινάκων ορίζουμε εντελώς ανάλογα με την σύγκλιση ακολουθίας διανυσμάτων. Αξιοσημείωτη είναι τέλος η ιδιότητα (Neumann) ότι αν  $\|A\| < 1$ , τότε ο πίνακας  $I-A$  είναι αντιστρέψιμος και ικανοποιεί τις ανισότητες

$$(3) \quad (1+\|A\|)^{-1} \leq \|(I-A)^{-1}\| \leq (1-\|A\|)^{-1}.$$

Προχωρούμε τώρα στην διερεύνηση της ευαισθησίας της λύσης του εστήματος  $Ax=b$ , όπου  $A$   $n \times n$  αντιστρέψιμος πίνακας σε διαταραχές των  $A$  και  $b$ . Ας μεταβάλλουμε κατ' αρχήν μόνο το δεύτερο μέλος  $b$  σε  $b+\delta b$ . Έστω  $x+\delta x$  η λύση του νέου εστήματος

$$(4) \quad A(x+\delta x) = b+\delta b.$$

Συμπεραίνουμε ότι  $\delta x = A^{-1}\delta b$ , δηλ. ότι για οποιαδήποτε νόρμα  $\|\cdot\|$ ,  $\|\delta x\| \leq \|A^{-1}\| \|\delta b\|$ , από την οποία έπεται για  $b \neq 0$  η

$$(5) \quad \|\delta x\|/\|x\| \leq \|A\| \|A^{-1}\| (\|\delta b\|/\|b\|).$$

Μπορούμε εύκολα να κατασκευάσουμε παραδείγματα (δεδομένων των  $A, b, \|\cdot\|$ ) διαταραχών  $\delta b$  για τα οποία η (5) ισχύει ως ισότητα (βλ. Ασκ. 3).

Συμπεραίνουμε ότι η ποσότητα  $\|A\| \|A^{-1}\|$  είναι ένας συντελεστής που προσδιορίζει πόσο μεγάλη μπορεί να γίνει η σχετική μεταβολή  $\|\delta x\|/\|x\|$  της λύσης του εστήματος όταν η σχετική μεταβολή (ή σχετικό "εφάλμα" αν θεωρήσουμε π.χ. ότι το  $\delta b$  παριστάνει το εφάλμα της προσεγγιστικής παράστασης στοιχείου του  $b$  στον υπολογιστή) του δεύτερου μέλους είναι  $\|\delta b\|/\|b\|$ . Ο αριθμός

$$(6) \quad \kappa(A) = \|A\| \|A^{-1}\|$$

λέγεται δείκτης κατάστασης του  $A$  ως προς την νόρμα  $\|\cdot\|$ . Αν ο  $\kappa(A)$

είναι πολύ μεγάλος τότε λέμε ότι ο πίνακας έχει κακή κατάσταση. (Αν ο  $\kappa(A)$  είναι μικρός - πάντα  $\kappa(A) \geq 1$  - λέμε ότι ο  $A$  έχει "καλή κατάσταση" - σημειώστε ότι η ευκρισιμότητα δύο οποιωνδήποτε νόρμων πινάκων οδηγεί σε ευκρισιμότητα των αντιστοίχων δεικτών κατάστασης: Για δύο οποιαδήποτε νόρμες  $\|\cdot\|_\alpha, \|\cdot\|_\beta$  υπάρχουν σταθερές  $c_1, c_2$  τέτοιες ώστε  $c_1 \kappa_\alpha(A) \leq \kappa_\beta(A) \leq c_2 \kappa_\alpha(A)$  για κάθε  $A$  αντιστρέψιμο). Αν η κατάσταση ενός πίνακα είναι κακή, είναι δυνατόν μία μικρή μεταβολή στο  $b$  να προκαλέσει μεγάλη μεταβολή στην λύση.

Ανάλογες παρατηρήσεις ισχύουν αν αντί του  $b$  μεταβάλλουμε τώρα τα στοιχεία του  $A$  μόνο, δηλ. θεωρήσουμε το σύστημα

$$(7) \quad (A + \delta A)(x + \delta x) = b,$$

όπου  $x = A^{-1}b$  και όπου υποθέτουμε ότι η διαταραχή  $\delta A$  είναι αρκετά μικρή έτσι ώστε και ο πίνακας  $A + \delta A$  να είναι και αυτός αντιστρέψιμος. Υποθέτουμε ευχεκριμένα ότι

$$(8) \quad \|\delta A\| \|A^{-1}\| < 1$$

από την οποία έπεται ότι  $\|( \delta A ) A^{-1} \| < 1 \Rightarrow I + ( \delta A ) A^{-1}$  αντιστρέψιμος  $\Rightarrow ( I + ( \delta A ) A^{-1} ) A = A + \delta A$  αντιστρέψιμος. Η (7) δίνει τώρα ότι  $\delta x = -(A + \delta A)^{-1} (\delta A)x = -A^{-1} (I + (\delta A)A^{-1})^{-1} (\delta A)x$ . Παίρνοντας νόρμες και χρησιμοποιώντας τις (3), (8) έχουμε  $\|\delta x\| \leq \|A^{-1}\| \|\delta A\| \|x\| / (1 - \|\delta A\| \|A^{-1}\|)$ , απ' την οποία για  $x \neq 0$  προκύπτει (με  $\kappa(A) = \|A\| \|A^{-1}\|$ ) η σχέση

$$(9) \quad \|\delta x\| / \|x\| \leq [\kappa(A) / (1 - \|\delta A\| \|A^{-1}\|)] (\|\delta A\| / \|A\|).$$

Βλέπουμε δηλ. ξανά τον ρόλο του  $\kappa(A)$  ως δείκτη επιρροής των σχετικών εφαλμάτων (μεταβολών)  $\|\delta A\| / \|A\|$  πάνω στην λύση του  $Ax = b$ . Στην γενική περίπτωση, όπου και ο  $A$  και το  $b$  μεταβάλλονται, υποθέτουμε ότι

$$(10) \quad (A + \delta A)(x + \delta x) = b + \delta b$$

και επίσης ότι ισχύει η (8), μπορούμε να δείξουμε για  $b \neq 0$  την

ανισότητα

$$(11) \quad \|δx\|/\|x\| \leq [\kappa(A)/(1-\|A^{-1}\| \|δA\|)] [\|δA\|/\|A\| + \|δb\|/\|b\|],$$

της οποίας οι (5), (9) αποτελούν ειδικές περιπτώσεις.

### Παρατηρήσεις

1. Ο δείκτης κατάστασης ενός πίνακα  $A$  τείνει στο  $\infty$  αν ο  $A$  τείνει να γίνει μη αντιστρέψιμος. Μάλιστα μπορεί να αποδειχθεί (Kahan)

$$(12) \quad (\kappa(A))^{-1} = \inf\{\|A-B\|/\|A\|, B \text{ μη αντιστρέψιμος}\}$$

δηλ. ότι ο  $(\kappa(A))^{-1}$  μετράει την (αχετική ως προς το μέγεθος του  $A$ ) απόσταση του  $A$  από το εύρολο των μη αντιστρεψίμων πινάκων. (Η απόδειξη του  $\leq$  στην (12) είναι εύκολη: Έστω  $B$   $n \times n$  μη αντιστρέψιμος. Τότε  $\exists x \neq 0$  τ.ώ.  $Bx=0$ . Συνεπώς  $\|A-B\| \|x\| \geq \|(A-B)x\| = \|Ax\| \geq \|x\|/\|A^{-1}\|$ . Η απόδειξη του  $\geq$  είναι δυσκολότερη).

2. Το μέγεθος του δείκτη κατάστασης δεν έχει όμως σχέση (για  $n > 2$ ) με το μέγεθος της ορίζουσας  $\det A$ , αν ο  $A$  είναι αντιστρέψιμος. Π.χ. ο  $100 \times 100$  διαγώνιος πίνακας  $D$  με στοιχεία  $d_{ii}=0.1$  έχει  $\det D=10^{-100}$  ενώ η κατάσταση του είναι φυσικά η καλύτερη δυνατή ( $\kappa(D)=1$  ως προς οποιαδήποτε νόρμα). Αν' την άλλη μεριά, για τον  $n \times n$  άνω τριγωνικό πίνακα  $A$  με στοιχεία  $a_{ii}=1$ ,  $1 \leq i \leq n$ ,  $a_{ij}=-1$ ,  $i < j$ , ισχύει ότι  $\det A=1$  αλλά  $\kappa_{\infty}(A)=n2^{n-1}$ .

3. Από την σχέση  $\|A\|_2 = \max_i (\lambda_i(A^*A))^{1/2}$  εύκολα βλέπουμε ότι  $\kappa_2(A) = (\mu_{\max}/\mu_{\min})^{1/2}$ , όπου  $\mu_{\max}$  ( $\mu_{\min}$ ) είναι η μέγιστη (ελάχιστη) ιδιοτιμή του πίνακα  $A^*A$ . Αν ο πίνακας  $A$  είναι αυτοσυζυγής, δηλ. αν  $A^*=A$ , τότε η παραπάνω σχέση απλοποιείται στην  $\kappa_2(A) = |\lambda_{\max}/\lambda_{\min}|$  όπου  $\lambda_{\max}$  ( $\lambda_{\min}$ ) είναι η μέγιστη (ελάχιστη) ιδιοτιμή του  $A$ .

4. Πολύ γνωστοί για την κακή κατάστασή τους, ακόμα και για μικρό  $n$ ,



## 1.2.7

είναι οι λεγόμενοι πίνακες Hilbert  $H_n$ ,  $n=1,2,3,\dots$  που ορίζονται ως  $(H_n)_{ij} = (i+j-1)^{-1}$ ,  $1 \leq i, j \leq n$ . Είναι πίνακες συμμετρικοί και θετικά ορισμένοι (γιατί;) των οποίων οι δείκτες κατάστασής τους για  $2 \leq n \leq 10$  δίδονται από τον πίνακα

$$\begin{array}{cccccccc} n & : & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ \kappa_2(H_n) & : & 1.9 \cdot 10^1 & 5.2 \cdot 10^2 & 1.6 \cdot 10^4 & 4.8 \cdot 10^5 & 1.5 \cdot 10^7 & 4.8 \cdot 10^8 & 1.5 \cdot 10^{10} \end{array}$$

$$\begin{array}{cc} n & : & 9 & 10 \\ \kappa_2(H_n) & : & 4.9 \cdot 10^{11} & 1.6 \cdot 10^{13} \end{array}$$

5. Έστω  $x \neq 0$  η ακριβής λύση του συστήματος  $Ax=b$  και έστω  $\tilde{x}$  προσέγγιση της  $x$ . Το υπόλοιπο της  $\tilde{x}$  ορίζεται ως  $r=A\tilde{x}-b$ . Εύκολα βλέπουμε ότι  $r=A(\tilde{x}-x)$ , δηλ. ότι  $\|\tilde{x}-x\| = \|A^{-1}r\| \leq \|A^{-1}\| \|r\| \leq \kappa(A) \|r\| \|x\| / \|b\|$ . Συμπεραίνουμε ότι

$$\|\tilde{x}-x\| / \|x\| \leq \kappa(A) \|r\| / \|b\|,$$

δηλ. ότι αν το υπόλοιπο  $r$  μίας προσεγγιστικής λύσης είναι μικρό αυτό δεν σημαίνει αναγκαστικά ότι το (σχετικό) σφάλμα της  $\|\tilde{x}-x\| / \|x\|$  θα είναι μικρό αν  $\kappa(A) \gg 1$ .

### Ασκήσεις 1.2

1. Στη βιβλιογραφία της αριθμητικής γραμμικής άλγεβρας συνηθίζεται να ονομάζεται "νόρμα πινάκων" μία απεικόνιση  $\|\cdot\|: \mathbb{C}^{n \times n} \rightarrow [0, \infty)$  που ικανοποιεί τις ιδιότητες (α)-(δ) στην σελίδα 1.2.3 (Ειδική περίπτωση λοιπόν "νόρμας πινάκων" είναι η "νόρμα τελεστών" - "φυσική" νόρμα όπως λέγεται - που παράγεται από μία διανυσματική νόρμα και που ορίζεται από τον τύπο (2)).

(α) Δείξτε ότι η παράσταση  $\|A\|_F = \left( \sum_{i,j=1}^n |a_{ij}|^2 \right)^{1/2}$  - η λεγόμενη νόρμα Frobenius του  $A$  - ορίζει μία "νόρμα πινάκων" που δεν

είναι "φυσική", δηλ. δεν παράγεται από διανυσματική νόρμα. Δείξτε ότι ισχύει  $\|A\|_2 \leq \|A\|_F \leq n^{1/2} \|A\|_2$  για κάθε  $A$ .

(β) Δείξτε ότι η παράσταση  $\max_{i,j} |a_{ij}|$  δεν είναι νόρμα πινάκων.

2. (α) Αν  $U^*U=I$  δείξτε ότι  $\|U^*AU\|_2 = \|A\|_2 \quad \forall A \in \mathbb{C}^{n \times n}$ .

(β) Δείξτε ότι  $\max_i |\lambda_i(A)| \leq \|A\|$  για κάθε  $A \in \mathbb{C}^{n \times n}$  και κάθε νόρμα  $\|\cdot\|$ .

(γ) Βρείτε σταθερές σύγκρισης για τα ζευγάρια νορμών από τις  $\|A\|_1, \|A\|_\infty, \|A\|_2$ .

3. (α) Έστω  $A \in \mathbb{R}^{n \times n}$  συμμετρικός αντιστρέψιμος πίνακας. θεωρείστε το πρόβλημα (4). Βρείτε ένα δεύτερο μέλος  $b$  και μία διαταραχή  $\delta b$  στον  $\mathbb{R}^n$  έτσι ώστε να ισχύουν οι ισότητες  $\|\delta x\|_2 = \|A^{-1}\|_2 \|\delta b\|_2$  και  $\|b\|_2 = \|A\|_2 \|x\|_2$  και κατά συνέπεια και η (5) εαν ισότητα. (Υπόδειξη: ο  $A$  ως συμμετρικός πίνακας έχει  $n$  ορθοκανονικά ιδιοδιανύσματα και  $n$  πραγματικές ιδιοτιμές. Δηλ. υπάρχει ορθογώνιος  $n \times n$  πίνακας  $Q (Q^T Q = I)$  τέτοιος ώστε  $Q^T A Q = D = \text{diag}(\lambda_1, \dots, \lambda_n)$  όπου  $\lambda_i$  οι ιδιοτιμές του  $A$  τις οποίες διατάσσουμε κατά φθίνουσα απόλυτη τιμή ως  $|\lambda_1| \geq |\lambda_2| \dots \geq |\lambda_n| > 0$ . Διαλέξτε  $b = Q(1, 0, 0, \dots, 0)^T$ ,  $\delta b = Q(0, 0, \dots, 0, 1)^T$ ).

(β) θεωρείστε το σύστημα  $Ax=b$  όπου

$$A = \begin{pmatrix} 4.1 & 2.8 \\ 9.7 & 6.6 \end{pmatrix} \quad b = \begin{pmatrix} 4.1 \\ 9.7 \end{pmatrix}$$

Δείξτε ότι αν  $\|\cdot\| = \|\cdot\|_1$ ,  $\delta b = (0.01, 0)^T$ , η (5) ισχύει εαν ισότητα.

4. (Η άσκηση αυτή αναφέρεται σε εφαρμογές της (12) της παρατήρησης 1.)

(α) Χωρίς να βρεθεί ο  $A^{-1}$  δείξτε ότι για τον πίνακα

$$A = \begin{pmatrix} 1.01 & .99 \\ .99 & 1.01 \end{pmatrix}$$

ισχύει ότι  $\kappa_{\infty}(A) = \|A\|_{\infty} \|A^{-1}\|_{\infty} \geq 100$ . (Υπόδειξη: ο πίνακας  $B: b_{ij}=1$  δεν είναι αντιστρέψιμος).

(β) Έστω  $A$  αντιστρέψιμος άνω ή κάτω τριγωνικός πίνακας. Δείξτε ότι  $\kappa_{\infty}(A) \geq \|A\|_{\infty} / \min_i |a_{ii}|$

(γ) Αν ο  $A$  αντιστρέψιμος και ο  $B$  ικανοποιεί την σχέση  $\|A-B\| < 1/\|A^{-1}\|$ , τότε ο  $B$  είναι αντιστρέψιμος.

5. (α) Να αποδειχθεί η (11).

(β) Αν οι πίνακες  $A, B$  είναι αντιστρέψιμοι, τότε

$$\|B^{-1}-A^{-1}\|/\|B^{-1}\| \leq \kappa(A) \|A-B\|/\|A\|.$$

(γ) Να αποδειχθούν οι ισχυρισμοί της παρατήρησης 2.

(δ) Να αποδειχθούν οι ισχυρισμοί της παρατήρησης 3.

6. Έχουμε πάντοτε ότι  $\kappa(A) = \|A\| \|A^{-1}\| \geq \|AA^{-1}\| = 1$ . Τι είδους πίνακες έχουν  $\kappa(A)=1$ ; Δείξτε ότι

(α) Αν ο  $A$  είναι διαγώνιος με ίσα στοιχεία, τότε  $\kappa(A)=1$ .

(β) Αν ο  $A$  είναι ισομετρία ως προς  $\|\cdot\|$ , δηλ. αν  $\|Ax\|=\|x\| \quad \forall x$ , τότε  $\kappa(A)=1$ .

(γ) Αν ο  $A$  είναι τέτοιος ώστε  $A^*A=I$ , τότε  $\kappa_2(A)=1$ .

7. (α) Αν ο πραγματικός πίνακας  $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$  είναι αντιστρέψιμος

και  $\epsilon = (a^2+b^2+c^2+d^2)/2|ad-bc|$ , δείξτε ότι  $\kappa_2(A) = \epsilon + (\epsilon^2 - 1)^{1/2}$ .

(β) Δείξτε ότι ο πίνακας

$$\begin{pmatrix} 100 & 99 \\ 99 & 98 \end{pmatrix}$$

(καθώς και οι πίνακες που προέρχονται από αυτόν με μεταθέσεις στοιχείων του) έχει τον μεγαλύτερο δείκτη κατάταξης μέσα στο εύρος των  $2 \times 2$  αντιστρέψιμων πινάκων που τα στοιχεία τους είναι θετικοί ακέραιοι  $\leq 100$ .

δ. θεωρείστε για  $\epsilon \in \mathbb{R}$  την λύση  $x(\epsilon)$  του συστήματος

$$\begin{cases} (A + \epsilon F)x(\epsilon) = b + \epsilon f, & \epsilon \neq 0 \\ x(0) = x \end{cases}$$

όπου  $A, F \in \mathbb{R}^{n \times n}$  δεδομένοι πίνακες ( $A$  αντιστρέψιμος) και  $b, f$  δεδομένα διανύσματα στον  $\mathbb{R}^n$ . Δείξτε ότι η  $\epsilon \mapsto x(\epsilon)$  είναι διαφορίσιμη σε μία περιοχή του μηδενός. Χρησιμοποιώντας το ανάπτυγμα Taylor της  $x(\epsilon)$  γύρω από το  $\epsilon=0$  δείξτε ότι

$$\|x(\epsilon) - x\| / \|\epsilon\| \leq \kappa(A) [(\|F\|/\|A\|) + (\|f\|/\|b\|)] + O(\epsilon^2).$$

## 1.3 ΣΦΡΑΝΑΤΑ ΣΤΡΟΓΓΥΛΕΥΣΗΣ ΣΤΗΝ ΑΠΛΑΙΟΤΗ ΓΑΥΣΣ

Προχωρούμε τώρα στην μελέτη του δευτέρου παράγοντα που, όπως αναφέραμε στις εισαγωγικές παρατηρήσεις της προηγούμενης παραγράφου, επηρεάζει την υπολογιστική λύση  $\tilde{x}$  του  $Ax=b$ ,  $A \in \mathbb{R}^{n \times n}$ ,  $b \in \mathbb{R}^n$ . Ο παράγοντας αυτός είναι η ευστάθεια (ή αστάθεια) του αλγορίθμου της απαλοιφής σε (μικρές) διαταραχές που θα προέλθουν βέβαια από τα σφάλματα ετρογχύλευσης στις πράξεις. Η ευστάθεια του αλγορίθμου της απαλοιφής είναι ένα θέμα ανεξάρτητο από την κατάσταση του ευστήματος. Στο τέλος όμως και οι δύο αυτοί παράγοντες θα συμβάλουν, όπως θα δούμε, στην διαμόρφωση του σφάλματος της  $\tilde{x}$ .

Κατ' αρχήν πρέπει να παρατηρήσουμε ότι στην πράξη η πρώτη φάση της απαλοιφής, δηλ. η ανάλυση  $PA=LU$ , δεν γίνεται όπως στην παρ. 1.1, όπου στο γενικό βήμα  $k$  φέρναμε στη θέση του οδηγού απλώς ένα μη

(k)

μηδενικό στοιχείο από τα  $a_{ik}$ ,  $i \geq k$ , π.χ. εκείνο με το μικρότερο δείκτη  $i$ . Αν ο  $A$  είναι αντιετρήσιμος (πράγμα που θα υποθέσουμε από εδώ κι εμπρός) η εύρεση ενός μη μηδενικού οδηγού είναι εξασφαλισμένη σε κάθε βήμα. Είναι προφανές όμως ότι θα πρέπει να περιμένουμε

(k)

προβλήματα αν κάποιος οδηγός  $a_{kk}$  είναι π.χ. πολύ μικρός σε απόλυτη τιμή. Ο λόγος είναι φυσικά η κατά προέχχιση παράσταση των αριθμών και η πεπερασμένη ακρίβεια με την οποία γίνονται οι πράξεις στο

(k)

υπολογιστή. Ένας (απόλυτα) μικρός οδηγός  $a_{kk}$  π.χ. μπορεί να οδηγήσει σε μεγάλους (σε απόλυτη τιμή) πολλαπλασιαστές στο βήμα  $k$

(k) (k)

$m_{ik} = a_{ik} / a_{kk}$  και συνεπώς σε μεγάλη απόλυτη ακρίβεια κατά τις

(k) (k)

αφαιρέσεις  $a_{ij} - m_{ik} a_{kj}$  στις οποίες είναι δυνατόν - υποθέτοντας ότι

(k) (k)

(k)

$a_{ij}, a_{kj} = O(1)$  - ο "μεγάλος" όρος  $-m_{ik} a_{kj}$  να εξαφανίσει του "μικρό"

(κ)

$a_{ij}$  μετά από την ετρογχύλευση στην ακρίβεια του υπολογιστή. Τέτοιου είδους εφάλματα ετρογχύλευσης διαδίδονται στους υπολογισμούς και καταστρέφουν την ακρίβεια της  $\tilde{x}$ . Για αριθμητικά παραδείγματα τέτοιων φαινομένων βλ. π.χ. [5.4, παρ. 2.2]. Η απλή απαλοιφή λοιπόν, όπως περιγράφηκε στη παρ. 1.1, δεν είναι ευεταθής αλγόριθμος. Προεπαθούμε λοιπόν να φέρουμε σε κάθε βήμα στην θέση του οδηγού ~~σε~~ το δυνατόν μεγαλύτερα σε απόλυτη τιμή στοιχεία. Μιά τέτοια στρατηγική λέγεται οδήγηση (pivoting). Μιά προφανής επιλογή στο

(κ)

k-ετό βήμα είναι να βρούμε εκείνο το στοιχείο από τα  $a_{ik}$ ,  $i \leq k$ , με την μεγαλύτερη απόλυτη τιμή και με εναλλαγή γραμμών να το φέρουμε στην θέση του οδηγού. Αυτή είναι η λεγόμενη μερική οδήγηση (partial pivoting) ή οδήγηση κατά γραμμές. (Η μερική οδήγηση ευχυνά βελτιώνεται με μία παραλλαγή της, την μερική οδήγηση με στάθμιση, βλ. [5.4, σελ. 21]). Το κόστος των εσυκρίσεων των στοιχείων κατά την μερική οδήγηση είναι  $O(n^2)$  και ευνενπώς εσυμπτωτικά πολύ μικρότερο από το κόστος της ανάλυσης LU.

Η μερική οδήγηση είναι στην πράξη εχεδών πάντα ασφαλής. Υπάρχουν όμως παραδείγματα "παθολογικών" ευστημάτων όπου δεν βελτιώνεται και πολύ την ακρίβεια στις πράξεις χιιά μεγάλο n (βλ. παρακάτω). Σε τέτοιες περιπτώσεις καταφεύχουμε στην λεγόμενη ολική οδήγηση (total pivoting) ή οδήγηση κατά γραμμές και ετήλες κατά την οποία με εναλλαγές γραμμών και ετηλών φέρουμε στην θέση του

(κ)

(κ)

οδηγού  $a_{kk}$  το στοιχείο εκείνο του υποπίνακα  $a_{ij}$ ,  $k \leq i, j \leq n$  με την μέγιστη απόλυτη τιμή. Μπορεί να αποδειχθεί (βλ. παρακάτω) ότι η ολική οδήγηση είναι πάντα ασφαλής, δηλ. ότι ο "ευντελεστής μεγέθυνσης" των εφαλμάτων ετρογχύλευσης αυξάνεται πολύ αρχά με το n (στη θεωρία' πρακτικά δεν αυξάνεται). Πάντως επανιπώτα χρησιμοποιούμε ολική οδήγηση στην πράξη χιιά και το κόστος της διπλασιάζει το κόστος των  $n^3/3$  περίπου πράξεων της ανάλυσης LU. Θα αναλύσουμε λοιπόν την ευστάθεια του αλγόριθμου της απαλοιφής Gauss (ανάλυση PA=LU και λύση

των δύο τριγωνικών ευστημάτων (1.1.7) και (1.1.8)\* με μερική οδήγηση, τις εναλλαγές γραμμών της οποίας καταγράφει ο πίνακας P. Αρχίζουμε με μία εύτομη επανάληψη περί παράστασης πραγματικών αριθμών στον υπολογιστή' για περισσότερες λεπτομέρειες βλ. [5.2] ή [5.4, παρ. 2.2].

Όπως είναι γνωστό σε κάθε υπολογιστή οι πραγματικοί αριθμοί παριστάνονται από ένα πεπερασμένο σύνολο ρητών, τους λεγόμενους αριθμούς της μηχανής (ή αριθμούς κινητής υποδιαστολής, πεπερασμένης ακρίβειας). Οι ρητοί αυτοί εκφράζονται με t ψηφία ("ακρίβεια" της μηχανής) ε' ένα αριθμητικό σύστημα με βάση β. Κάθε αριθμός της μηχανής είναι της μορφής  $y = \pm d_1 d_2 \dots d_t \cdot \beta^e$  (ή  $y=0$ ), όπου οι ακέραιοι  $d_i$ ,  $1 \leq i \leq t$  είναι ψηφία στο σύστημα με βάση β, δηλ.  $0 \leq d_i \leq \beta - 1$  με  $d_t \neq 0$ . Ο εκθέτης e είναι ακέραιος θετικός ή αρνητικός αλλά περιορίζεται ε' ένα πεπερασμένο διάστημα. Είναι φανερό ότι κάθε τέτοιο σύνολο αριθμών είναι πεπερασμένο. Κάθε πραγματικός x που βρίσκεται μέσα στο εύρος των αριθμών της μηχανής προεχχίζεται είτε με "αποκοπή" είτε με "ετροχχύλευση" από ένα κουτιό του αριθμό της μηχανής που συμβολίζουμε με  $f(x)$ . Κατά τα γνωστά ισχύει

$$(1) \quad f(x) = x(1 + \delta), \quad \text{όπου } |\delta| \leq u = \begin{cases} \beta^{1-t} & \text{για αποκοπή} \\ (\beta^{1-t})/2 & \text{για ετροχχύλευση} \end{cases}$$

Ο (μικρός) αριθμός u λέγεται "μοναδιαίο εφάλμα ετροχχύλευσης". Η (1) λέει απλούστατα ότι το σχετικό εφάλμα  $|(f(x)-x)/x|$  είναι μικρότερο ή ίσο του u.

Δεδομένων δύο αριθμών της μηχανής x και y θα συμβολίζουμε με  $f(x @ y)$  όπου @ = +, -, ' ή / το αποτέλεσμα της κατά προέχχειν πράξεως, αντίστοιχα αφαίρεσης, πολ/εμού, διαίρεσής τους ετην αριθμητική μονάδα της μηχανής με απλή ακρίβεια.

\* Η αναφορά σε τύπους, θεωρήματα, παρατηρήσεις, ασκήσεις κλπ. (z) μιάς διαφορετικής παραγράφου x.y θα συμβολίζεται με (x.y.z). Η αναφορά σε τύπους, κλπ. (z) της ίδιας παραγράφου θα συμβολίζεται με (z).

## 1.3.4

Υποθέτουμε (απλουστευτική αλλά αρκετά ρεαλιστική παραδοχή) ότι πρώτα υπολογίζεται ακριβώς ο πραγματικός αριθμός  $x@y$  - που βρίσκεται μέσα στο εύρος των αριθμών της μηχανής - και κατόπιν ετροχχυλεύεται (ή αποκόπτεται) στον αριθμό μηχανής  $fl(x@y)$ . Συνεπώς εάν ευνέπεια της (1) έχουμε ότι

$$(2) \quad fl(x@y) = (x@y)(1+\delta), \quad |\delta| \leq u.$$

θα μάς είναι χρήσιμη και η εξής εναλλακτική μορφή της (ψευδο)ισότητας (2) - βλ. Άσκηση 1:

$$(2') \quad fl(x@y) = (x@y)/(1+\delta'), \quad |\delta'| \leq u.$$

Χρησιμοποιώντας τώρα την (2) ή την (2') μπορούμε να βρούμε ανάλογες εκφράσεις για πιο πολύπλοκες πράξεις. Π.χ. ορίζοντας, για  $x, y, z$  αριθμούς μηχανής, την παράσταση του αθροίσματός τους  $fl(x+y+z)$  ως  $fl(fl(x+y)+z)$  έχουμε από κατ' επανάληψιν χρήση της (2) ότι  $fl(x+y+z) = fl((x+y)(1+\delta_1)+z) = ((x+y)(1+\delta_1)+z)(1+\delta_2)$   
 $= (x+y)(1+\delta_1)(1+\delta_2)+z(1+\delta_2)$ ,  $|\delta_i| \leq u$ . Είναι προφανές ότι το  $fl(x+y+z)$  είναι διάφορο γενικά του  $fl(x+z+y)$  κλπ. (Η προερχιστική πρόθεση

δεν είναι προεταίριετική!) Γιαυτό γράφοντας  $fl(\sum_{i=1}^m x_i)$  θα εννοούμε ότι πρώτα προσθέτουμε τα  $x_1$  και  $x_2$ , στο αποτέλεσμα το  $x_3$  κ.ο.κ, δηλ. ότι υλοποιούμε τον αλγόριθμο

$$\begin{aligned} s &\leftarrow x_1 \\ \left[ \begin{array}{l} \text{Γιά } i=2, \dots, m \\ s \leftarrow s+x_i \end{array} \right. \end{aligned}$$

όπου  $s \leftarrow$  συμβολίζουμε την εκχώρηση της τιμής της δεξιάς μεταβλητής στη θέση που βρίσκεται αποθηκευμένη η αριστερή. (Υποθέτουμε ειωηηρά ότι όλες οι πράξεις δίνουν αποτελέσματα μετά στο εύρος των αριθμών της μηχανής).



## 1.3.5

θα εφαρμόσουμε τώρα τα παραπάνω εν είδει παραδείγματος ε' ένα πρόβλημα που εμφανίζεται επανειλημμένα στην αριθμητική γραμμική

άλγεβρα. θεωρούμε την πράξη  $f(\sum_{i=1}^n x_i y_i)$ , δηλ. τον υπολογισμό του Ευκλείδειου εσωτερικού γινομένου δύο πραγματικών διανυσμάτων  $x, y \in \mathbb{R}^n$  με στοιχεία αριθμούς μηχανής. θα επιχειρήσουμε να βρούμε μία

εύχρηστη έκφραση του  $f(\sum_{i=1}^n x_i y_i)$  συναρτήσει των  $x_i, y_i, n$  και  $u$ . Γιαυτόν τον σκοπό αποδεικνύουμε πρώτα τα εξής βοηθητικά αποτελέσματα:

**Λήμμα 1.** Έστω  $0 \leq u \leq 1$  και  $n \in \mathbb{N}$ . Τότε

$$(\alpha) \quad 1 - nu \leq (1 - u)^n.$$

$$(\beta) \quad \text{Αν } 0 < nu \leq 0.01, \text{ τότε } (1 + u)^n \leq 1 + 1.01nu.$$

$$(\gamma) \quad \text{Αν } |\delta_i| \leq u, \quad i = 1, 2, \dots, n \text{ και } 0 < nu \leq 0.01, \text{ τότε}$$

$$(3) \quad 1 - nu \leq \prod_{i=1}^n (1 + \delta_i) \leq 1 + 1.01nu.$$

Απόδειξη. Το (α) προκύπτει εύκολα είτε με επαγωγή είτε χρησιμοποιώντας το θεώρημα του Taylor για την συνάρτηση  $f(u) = (1 - u)^n$  γύρω από το 0 με δύο όρους. Για το (β) χρησιμοποιούμε το γεγονός ότι  $1 + x \leq e^x$  για  $x \geq 0$  (προφανές) και το ότι  $e^x \leq 1 + 1.01x$  για  $x \in [0, 0.01]$ , που προκύπτει από την ανισότητα  $e^x \leq 1 + x + x^2$  που εύκολα ισχύει για  $0 \leq x \leq 0.01$ . Συνεπώς  $(1 + u)^n \leq e^{nu} \leq 1 + 1.01nu$ . Τέλος το (γ) είναι συνέπεια των ανισοτήτων (α) και (β) και της  $1 - u \leq 1 + \delta_i \leq 1 + u$ . @

Σημειώστε ότι η (3), γράφεται και στην μορφή

$$(3') \quad \prod_{i=1}^n (1 + \delta_i) = 1 + 1.01n\theta \text{ για κάποιο } \theta: |\theta| \leq 1.$$

Προχωρούμε τώρα στην ζητούμενη έκφραση του  $f(\sum_{i=1}^n x_i y_i)$ .  
 Τυπίζουμε ξανά ότι έχει σημασία η σειρά με την οποία γίνονται οι  
 πράξεις στο εσωτερικό γινόμενο. Θα υποθέσουμε ότι ο αριθμός

$s = f(\sum_{i=1}^n x_i y_i)$  υπολογίζεται από τον αλγόριθμο

$$(4) \quad \begin{cases} s \leftarrow x_1 y_1 \\ \text{Γι } i=2, \dots, n \\ s \leftarrow s + (x_i y_i) \end{cases}$$

**Λήμμα 2.** Έστω  $u$  το μοναδιαίο εφάλμα ετροχύλευσης και  $x_i, y_i$ ,  
 $1 \leq i \leq n$  αριθμοί μηχανής. Τότε αν  $n u \leq 0.01$ , έχουμε

$$(5) \quad f(\sum_{i=1}^n x_i y_i) = \sum_{i=1}^n x_i y_i [1 + 1.01(n+2-i)\theta_i u],$$

όπου  $|\theta_i| \leq 1$ ,  $1 \leq i \leq n$ .

Απόδειξη: Λόγω του αλγορίθμου (4) έχουμε

$$(6) \quad f(\sum_{i=1}^n x_i y_i) = \underbrace{f(f(\dots(f(x_1 y_1) + f(x_2 y_2)) + f(x_3 y_3)) + \dots + f(x_4 y_4)) + \dots + f(x_n y_n))}_{n-1 \text{ "f"}} = x_1 y_1 \prod_{i=1}^n (1 + \delta_i^{(1)}) + x_2 y_2 \prod_{i=1}^{n-1} (1 + \delta_i^{(2)}) + x_3 y_3 \prod_{i=1}^{n-2} (1 + \delta_i^{(3)}) + \dots + x_n y_n \prod_{i=1}^2 (1 + \delta_i^{(n)})$$

(j)

όπου  $|\delta_i| \leq u$ . Για  $k > 1$  έχουμε από την (3') και την υπόθεση  $nu \leq 0.01$  ότι

$$(7) \quad x_k y_k \prod_{i=1}^{n-k+2} (1+\delta_i) = x_k y_k [1+1.01(n+2-k)\theta_k u]$$

για κάποιο  $\theta_k : |\theta_k| \leq 1$ . Για  $k=1$  έχουμε από την (3') ότι

$$x_1 y_1 \prod_{i=1}^n (1+\delta_i) = x_1 y_1 [1+1.01n\theta'_1 u], \quad |\theta'_1| \leq 1. \quad \text{Θέτουμε } \theta_1 = n\theta'_1/(n+1)$$

οπότε  $|\theta_1| \leq 1$ . Γιαυτό το  $\theta_1$  λοιπόν ισχύει η (7) και για  $k=1$ . Οι (6) και (7) δίνουν τώρα την (5). @

Ας αρχίσουμε τώρα την μελέτη της επίρροής των εφαλμάτων ετρογχύλευσης στον αλγόριθμο της απαλοιφής με μερική οδήγηση. Χρησιμοποιώντας τις (1), (2) σε κάθε πράξη της απαλοιφής είναι δυνατόν να παρακολουθήσουμε την ευσέρευση των εφαλμάτων ετρογχύλευσης. Η λογιστική όμως ενός τέτοιου εγχειρήματος (που λέγεται και "απ' ευθείας" (forward) ανάλυση του εφάλματος) είναι αρκετά πολύπλοκη. Επιπλέον μιά τέτοια στρατηγική συνήθως δίνει πολύ μεγάλα και μη ρεαλιστικά φράγματα εφαλμάτων. Μιά άλλη τεχνική, η λεγόμενη "αντίστροφη" (inverse, backward) ανάλυση του εφάλματος οφείλεται στον άγγλο μαθηματικό J.H. Wilkinson ο οποίος απέδειξε (δεκαετία του '60) ότι η υπολογιστική λύση  $\tilde{x}$  που δίνει η απαλοιφή μπορεί να θεωρηθεί ως ακριβής λύση (δηλ. λύση με αριθμητική απεριόριστη ακρίβεια) όχι του ευστήματος  $Ax=b$  αλλά ενός "παραπλήσιου" ευστήματος  $(A+\delta A)\tilde{x}=b$ . Η ευστάθεια ή μη του αλγορίθμου εκφράζεται ποσοτικά από το αν η σχετική μεταβολή  $\|\delta A\|/\|A\|$  είναι μικρή ή όχι. Στη συνέχεια θα εκτιμήσουμε αυτήν την μεταβολή στην  $\|\cdot\|_\infty$  νόρμα.

Υποθέτουμε, για να απλουτεύουμε το πρόβλημα ελαφρώς, ότι τα στοιχεία των  $A, b$  είναι αριθμοί της μηχανής (βλέπε παρατήρηση 1). Επίσης υποθέτουμε ότι έχουν γίνει ήδη εκ των προτέρων (αν και αυτό φυσικά δεν μπορεί να γίνει στην πράξη) όλες οι εναλλαγές γραμμών

(k) που υπαγορεύει η μερική οδήγηση έτσι ώστε για κάθε  $k$  τα στοιχεία  $a_{kk}$

που θα προκύπτουν από τους υπολογισμούς στις θέσεις των οδηγών να

(k)

είναι μεγαλύτερα ή ίσα απολύτως απ' όλα τα στοιχεία  $a_{ik}$   $i > k$ .

Προφανώς η εναλλαγή γραμμών δεν επηρεάζει την αριθμητική για την οποία ενδιαφεράμαστε εδώ. Υποθέτουμε επίσης ότι όλες οι πράξεις οδηγούν σε αριθμούς μέσα στο εύρος των αριθμών της μηχανής έτσι ώστε να μην σταματά ο υπολογισμός λόγω over- ή underflow. (Για

(k)

απλούστευση του συμβολισμού δεν συμβολίζουμε με  $\tilde{L}, \tilde{U}, \tilde{m}_{ik}, \tilde{a}_{ij}$  κλπ. τα μεγέθη που υπολογίζουμε με αριθμητική πεπερασμένης ακρίβειας

(k)

κατά την απαλοιφή, αλλά χρησιμοποιούμε τα σύμβολα  $L, U, m_{ik}, a_{ij}$  κλπ. έχοντας βέβαια υπ' όψιν μας ότι τα μεγέθη αυτά θα είναι ευ γένει διαφορετικά από τα αντίστοιχα μεγέθη που εμφανίζονται στην παρ. 1.1 ως αποτελέσματα ακριβών υπολογισμών). Στο k-στό βήμα της απαλοιφής υπολογίζουμε τους πολλαπλασιαστές ως

$$(8) \quad m_{ik} = fl(a_{ik} / a_{kk}), \quad k+1 \leq i \leq n.$$

(k+1)

Τα στοιχεία  $a_{ij}$  του πίνακα  $A^{(k)}$  υπολογίζονται κατόπιν βάσει του τύπου

$$(9) \quad a_{ij}^{(k+1)} = \begin{cases} 0 & , \text{ αν } j=k, \quad k+1 \leq i \leq n \\ \begin{matrix} (k) & (k) \\ fl(a_{ij} - m_{ik} a_{kj}) \end{matrix} & , \text{ αν } k+1 \leq i, j \leq n \\ (k) \\ a_{ij} & , \text{ αλλιώς} \end{cases}$$

Στο τέλος αυτής της διαδικασίας ορίζουμε τους nxn πίνακες

$$(10) \quad U = A^{(n)}$$

$$(11) \quad L : L_{ij} = \begin{cases} 0, & \text{αν } i < j \\ 1, & \text{αν } i = j \\ m_{ij}, & \text{αν } i > j \end{cases}$$

Παύζουμε ότι  $A \neq LU$  γιατί τα στοιχεία των  $L, U$  έχουν υπολογισθεί με αριθμητική πεπερασμένης ακρίβειας αναδρομικά από τις (8) & (9). Έχουμε όμως

**Λήμμα 3.** Αν οι πίνακες  $L, U$  ορίζονται από τις (10), (11) έχουμε

$$(12) \quad LU = A + E,$$

όπου

$$(13) \quad E = \sum_{k=1}^{n-1} E^{(k)}, \quad E^{(k)} \in \mathbb{R}^{n \times n}.$$

Για  $1 \leq k \leq n-1$ , τα στοιχεία  $\varepsilon_{ij}^{(k)}$  του  $E^{(k)}$  δίδονται από

$$(14) \quad \varepsilon_{ij}^{(k)} = \begin{cases} a_{ik}^{(k)} \delta_{ik}^{(k)} & k+1 \leq i \leq n, j=k \\ -m_{ik}^{(k)} a_{kj}^{(k)} \delta'_{ij} - a_{ij}^{(k+1)} \delta''_{ij}, & k+1 \leq i, j \leq n \\ 0 & \text{αλλιώς} \end{cases}$$

όπου τα  $\delta_{ij}, \delta'_{ij}, \delta''_{ij}$  εξαρτώνται και από το  $k$  και είναι τέτοια ώστε  $|\delta_{ij}|, |\delta'_{ij}|, |\delta''_{ij}| \leq u$ , όπου  $u$  το μοναδιαίο εφάλμα ετρογχύλευσης

Απόδειξη: Η (8) δίνει, λόγω της (2),  $m_{ik}^{(k)} = (a_{ik}^{(k)} / a_{kk}^{(k)})(1 + \delta_{ik}^{(k)})$ ,  $i \geq k+1$ , όπου  $|\delta_{ik}^{(k)}| \leq u$ . Γράφουμε την σχέση αυτή ως

$$(15) \quad a_{kk}^{(k)} m_{ik} - a_{ik}^{(k)} - \varepsilon_{ik} = 0, \quad i \geq k+1,$$

όπου ορίσαμε

$$(16) \quad \varepsilon_{ik} = a_{ik}^{(k)} \delta_{ik}, \quad i \geq k+1.$$

Η (9) δίνει, με χρήση των (2) και (2')

$$\begin{aligned} a_{ij}^{(k+1)} &= fl(a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)}) = fl(a_{ij}^{(k)} - fl(m_{ik} a_{kj}^{(k)})) = \\ &= (a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} (1 + \delta'_{ij})) / (1 + \delta''_{ij}), \quad k+1 \leq i, j \leq n, \end{aligned}$$

όπου  $|\delta'_{ij}|, |\delta''_{ij}| \leq u$ . Ξαναγράφουμε την σχέση αυτή ως

$$(17) \quad a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} + \varepsilon_{ij}, \quad k+1 \leq i, j \leq n,$$

όπου ορίσαμε

$$(18) \quad \varepsilon_{ij} = -m_{ik} a_{kj}^{(k)} \delta'_{ij} - a_{ij}^{(k)} \delta''_{ij}, \quad k+1 \leq i, j \leq n.$$

(Σημειώστε ότι τα  $\varepsilon_{ij}$  της (16) και της (18) θα ήταν μηδέν αν δεν γίνονταν εφάλματα ετρογχύλευσης κατά το  $k$ -στό βήμα). Ορίζουμε τώρα

$\varepsilon_{ij}^{(k)} = 0$  για τα υπόλοιπα  $i, j$  (δηλ. συνολικά ορίζουμε τα  $\varepsilon_{ij}^{(k)}$  όπως έτην

(14)) και θεωρούμε τον  $n \times n$  πίνακα  $E^{(k)}$  με στοιχεία  $\varepsilon_{ij}^{(k)}$ ,  $1 \leq i, j \leq n$ .

Έστω ο  $n \times n$  πίνακας  $L^{(k)}$  όπου

$$L_{ij}^{(k)} = \begin{cases} m_{ik}, & \text{αν } k+1 \leq i \leq n, j=k \\ 0, & \text{αλλιώς} \end{cases}$$

Πορούμε τώρα εύκολα να δείξουμε, χρησιμοποιώντας τις σχέσεις (15), (17) και (9) ότι

$$A^{(k+1)} = A^{(k)} - L^{(k)}A^{(k)} + E^{(k)}.$$

Αθροίζοντας και τα δύο μέλη αυτής της ισότητας ως προς  $k$  από  $k=1$  έως  $n-1$  παίρνουμε.

$$(19) \quad L^{(1)}A^{(1)} + L^{(2)}A^{(2)} + \dots + L^{(n-1)}A^{(n-1)} + A^{(n)} = A^{(1)} + \sum_{k=1}^{n-1} E^{(k)}.$$

Υπολογίζοντας τα στοιχεία του πίνακα  $L^{(1)}A^{(1)}$  βλέπουμε ότι για  $1 \leq i, j \leq n$

$$(L^{(1)}A^{(1)})_{ij} = \begin{cases} 0, & \text{αν } i=1 \\ m_{i1}a_{1j}, & \text{αλλιώς} \end{cases} = (L^{(1)}A^{(n)})_{ij}.$$

Συμπεώς  $L^{(1)}A^{(1)} = L^{(1)}A^{(n)}$ . Κατά τον ίδιο τρόπο, επειδή η  $k$ -στή γραμμή του  $A^{(k)}$  συμπίπτει με την  $k$ -στή γραμμή του  $A^{(n)}$ , συμπεραίνουμε ότι ο πίνακας  $L^{(k)}A^{(k)}$  (που εξαρτάται, εκτός από τα  $m_{ik}$ , μόνο από την  $k$ -στή γραμμή του  $A^{(k)}$ ) συμπίπτει με τον  $L^{(k)}A^{(n)}$ .

Συμπεώς η (19) δίνει, λόγω της (13)

$$(L^{(1)} + L^{(2)} + \dots + L^{(n-1)} + I)A^{(n)} = A + E,$$

από την οποία προκύπτει η (12) λόγω της (10) και του ότι (16) και

$$\sum_{k=1}^{n-1} L^{(k)} + I = L, \text{ βλ. (11) @}$$

Στο επόμενο βήμα φράσουμε τα στοιχεία του πίνακα  $E$  καθώς και

την νόρμα  $\|E\|_\infty$ .

**Λήμμα 4.** Για τον πίνακα  $E$  που κατασκευάσαμε στο Λήμμα 3. ισχύει

$$(20) \quad \|E\|_\infty < n^2 \rho \cdot \|A\|_\infty \cdot \mu,$$

όπου  $\mu$  το μοναδιαίο εφάλμα ετροχγύλευσης και όπου

$$(21) \quad \rho = \max_{i,j,k} |a_{ij}| / \|A\|_\infty,$$

όπου τα  $a_{ij}$  είναι τα ενδιάμεσα προϊόντα των υπολογισμών της απαλοιφής που ορίζονται αναδρομικά από την (9).

Απόδειξη. Για να να φράξουμε τα στοιχεία  $E_{ij} = \sum_{k=1}^{n-1} \varepsilon_{ij}$  του  $E$

πρέπει να φράξουμε, λόγω της (14), τις ποσότητες  $m_{ik}$  και  $a_{ij}$ .  
 θυμόμαστε την υπόθεση ότι οι γραμμές του  $A$  έχουν διαταχθεί κατάλληλα  
 έτσι ώστε οι οδηγοί που προκύπτουν από τους υπολογισμούς με

αριθμητική πεπερασμένης ακρίβειας, δηλ. τα  $a_{kk}$ , να είναι τα απολύτως

μεγαλύτερα, για κάθε  $k$ , από τα στοιχεία  $a_{ik}$ ,  $k \leq i \leq n$ . Επειδή

$m_{ik} = f_1(a_{ik} / a_{kk})$ , ανακαλύπτας την παραδοχή μας ότι πρώτα γίνεται η

πράξη  $a_{ik} / a_{kk}$  ακριβώς και μετά ετροχγυλεύεται ή αποκόπτεται το  
 αποτέλεσμα, βλέπουμε ότι αυτό σημαίνει ότι

$$(22) \quad |m_{ik}| \leq 1 \text{ για όλα τα } i, k, (i > k).$$



Ορίζοντας την ποσότητα  $\rho$  από την (21) - θα επανέλθουμε βέβαια γιά να εσολιάσουμε την σημασία της και τις θεωρητικές και πρακτικές εκτιμήσεις της - πρὸς το παρόν αρκούμαστε να παρατηρήσουμε ὅτι μπορεί να υπολογισθεῖ εύκολα εαν παραπροϊόν της απαλοιφής - "λύσουμε" το

(k)

πρόβλημα της εκτίμησης των  $a_{ij}$  γιατί βέβαια εἶναι ὁριεμό

(k)

$$(23) \quad |a_{ij}| \leq \rho \|A\|_{\infty}.$$

Από τις (14), (22) και (23) συμπεραίνουμε ὅτι

$$(24) \quad |e_{ij}| \leq \rho \|A\|_{\infty} \cdot \begin{cases} u & \text{αν } k+1 \leq i \leq n, j=k \\ 2u & \text{αν } k+1 \leq i, j \leq n \\ 0 & \text{αλλιῶς.} \end{cases}$$

Ἐετω  $|B|$  ο πίνακας με στοιχεία τις απόλυτες τιμές  $|b_{ij}|$  των στοιχείων ενός πίνακα  $B$ . Επίσης αν  $B, C \in \mathbb{R}^{n \times n}$  θα γράψουμε  $B \leq C$  αν και μόνο αν  $b_{ij} \leq c_{ij}$ ,  $1 \leq i, j \leq n$ . Ανακαλύπτας τον ὁριεμό του πίνακα  $E$  ἔχουμε από την (24) ὅτι

$$(25) \quad |E| \leq \rho \|A\|_{\infty} u C,$$

ὅπου ο  $C$  εἶναι ο  $n \times n$  πίνακας που δίνεται από το ἄθροισμα των  $n$   $n \times n$  πινάκων

$$\begin{pmatrix} 0 & 0 & \dots & 0 \\ 1 & 2 & \dots & 2 \\ 1 & 2 & \dots & 2 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 2 & \dots & 2 \end{pmatrix} + \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 1 & 2 & \dots & 2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & 2 & \dots & 2 \end{pmatrix} + \dots + \begin{pmatrix} 0 & \dots & \dots & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & \dots & \dots & 0 \\ 0 & 0 & \dots & 0 & 1 & 2 \end{pmatrix}$$

Προφανῶς ο  $C$  εἶναι ο πίνακας

$$C = \begin{pmatrix} 0 & 0 & 0 & 0 & \dots & 0 & 0 \\ 1 & 2 & 2 & 2 & \dots & 2 & 2 \\ 1 & 3 & 4 & 4 & \dots & 4 & 4 \\ 1 & 3 & 5 & 6 & \dots & 6 & 6 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 3 & 5 & \dots & \dots & 2n-4 & 2n-4 \\ 1 & 3 & 5 & \dots & \dots & 2n-3 & 2n-2 \end{pmatrix}$$

Άρα η (25) δίνει

$$\begin{aligned} \|E\|_{\infty} &\leq \rho \|A\|_{\infty} \text{ u } \|C\|_{\infty} = \rho \|A\|_{\infty} \text{ u } \sum_{j=1}^n C_{nj} = \\ &= \rho \|A\|_{\infty} \text{ u } \left( \sum_{j=1}^{n-1} (2j-1) + 2n-2 \right) = (n^2-1) \rho \|A\|_{\infty} \text{ u } \text{ ó.é.δ } @ \end{aligned}$$

Συχνωνεύουμε τώρα τα λήμματα 3 και 4 στο παρακάτω

**ΘΕΩΡΗΜΑ 1.** Οι πίνακες  $L, U$  που υπολογίζουμε κατά την απαλοιφή Gauss με μερική οδήγηση χρησιμοποιώντας αριθμητική πεπερασμένης ακρίβειας με μοναδιαίο εφάλμα ετρογχύλευσης  $\text{u}$  ικανοποιούν ακριβώς την ιδιότητα

$$(26) \quad LU = A + E,$$

όπου

$$(27) \quad \|E\|_{\infty} < n^2 \rho \|A\|_{\infty} \text{ u},$$

και όπου το  $\rho$  ορίζεται από την (21). @

Το πρώτο αυτό θεώρημα μας λέει ότι τα προϊόντα  $L, U$  των υπολογισμών της απαλοιφής είναι η (ακριβής) ανάλυση  $LU$  ενός πίνακα  $A+E$ , κοντινού στον  $A$ , εφ' όσον το  $n$  δεν είναι πολύ μεγάλο και εφ'

(k)

όπου η ποσότητα  $\max_{i,j,k} |a_{ij}|$  - περί της οποίας περιεσσότερα παρακάτω - παραμένει φραγμένη. Επαναλαμβάνουμε ότι έχουμε υποθέσει ότι οι εναλλαγές γραμμών έγιναν όλες εκ των προτέρων.

Μετά την κατασκευή των  $L$  και  $U$ , η κατά προσέγγιση λύση  $\tilde{x}$  του συστήματος  $Ax=b$  θα βρεθεί με την επίλυση (με αριθμητική πεπερασμένης ακρίβειας) των τριγωνικών συστημάτων  $Ly=b$ ,  $U\tilde{x}=y$ . Σφάλματα ετρόχυλεύσης θα υπερεέλθουν φυσικά και ε' αυτές τις πράξεις. Ας εξετάσουμε λεπτομερώς π.χ. την αριθμητική στο πρώτο σύστημα  $Ly=b$ . Τα  $y_i$  θα βρεθούν από τον αλγόριθμο

$$(28) \quad \begin{cases} y_1 = f_1(b_1/L_{11}) \\ y_i = f_i((-L_{i1}y_1 - L_{i2}y_2 - \dots - L_{i,i-1}y_{i-1} + b_i)/L_{ii}), \quad i=2, \dots, n \end{cases}$$

Στον (28) αγνοούμε το γεγονός ότι  $L_{ii}=1$  και συνεπώς ότι η διαίρεση διά  $L_{ii}$  είναι τετριμμένη - και ότι δεν εκτελείται! - έτσι ώστε τα αποτελέσματα να ισχύουν για οποιοδήποτε κάτω (και κατ' επέκταση και άνω) τριγωνικό σύστημα. Σημαντική είναι η σειρά των πράξεων κατά του υπολογισμού του  $y_i$ . Όπως χράφτηκε ο (28), υπονοεί ότι ο αλγόριθμος της κατασκευής του  $y_i$  είναι ο

$$(29) \quad \begin{aligned} & y_i \leftarrow -\sum_{j=1}^{i-1} L_{ij} y_j \\ & y_i \leftarrow y_i + b_i \\ & y_i \leftarrow y_i / L_{ii} \end{aligned}$$

όπου ο υπολογισμός του εσωτερικού γινομένου  $\sum_{j=1}^{i-1} L_{ij} y_j$  γίνεται όπως στην (4). Θα ήταν π.χ. διαφορετικά τα αποτελέσματα αν ο αλγόριθμος ήταν ο



το οποίο ικανοποιεί ακριβώς η  $y$  (αντί του  $Ly=b$ ) όταν υπάρχουν εφάλματα ετροχγύλευσης. Εύκολα βλέπουμε ότι

$$\|\delta L\| \leq 1.01 \text{ u } \begin{pmatrix} |L_{11}| & & & & 0 \\ |L_{21}| & 2|L_{22}| & & & \\ 2|L_{31}| & 2|L_{32}| & 2|L_{33}| & & \\ \vdots & \vdots & \vdots & \ddots & \\ (n-1)|L_{n1}| & (n-1)|L_{n2}| & (n-2)|L_{n3}| & \dots & 2|L_{nn}| \end{pmatrix}$$

Συνοψώς

$$(32') \quad \|\delta L\|_{\infty} \leq 1.01 \text{ u } \max_{i,j} |L_{ij}| n(n+1)/2 \leq 1.01 \text{ u } n(n+1)/2,$$

(ανακαλώντας ότι  $\max_{i \neq j} |L_{ij}| \leq 1$ , βλ. (22),  $L_{ii}=1$ ).

Ευτελώς παρόμοια ισχύουν για το άνω τριγωνικό εύστημα  $U\tilde{x}=y$ : Η λύση  $\tilde{x}$  ικανοποιεί ακριβώς την εξίσωση

$$(33) \quad (U+\delta U)\tilde{x}=y,$$

όπου, όπως στην (32'),

$$(34) \quad \|\delta U\|_{\infty} \leq 0.01 \text{ u } \max_{i,j} |U_{ij}| n(n+1)/2 \leq 0.01 \text{ u } \|U\|_{\infty} n(n+1)/2,$$

επειδή  $U=A^{(n)}$  και λόγω της (21). Παίρνουμε λοιπόν το εξής

**Λήμμα 5.** Η προσεχχιστική λύση  $\tilde{x}$  του ευστήματος  $Ax=b$  που παίρνουμε με απαλοιφή Gauss με μερική οδήγηση είναι ακριβής λύση του ευστήματος

$$(35) \quad (L+\delta L)(U+\delta U)\tilde{x}=b,$$

όπου οι ποσότητες  $\|\delta L\|_{\infty}$ ,  $\|\delta U\|_{\infty}$  φράσσονται όπως στις (32'), (34). @

Τέλος φθάνουμε στο κεντρικό αποτέλεσμα αυτής της παραγράφου:

**ΘΕΩΡΗΜΑ 2.** Η προεχχιστική λύση  $\tilde{x}$  του  $Ax=b$  που δίνει η απαλοιφή Gauss με μερική οδήγηση είναι ακριβής λύση του συστήματος

$$(36) \quad (A+\delta A)\tilde{x}=b,$$

όπου

$$(37) \quad \|\delta A\|_{\infty} \leq 1.01 (n^3+3n^2) \rho \|A\|_{\infty} u,$$

όπου το  $\rho$  ορίζεται από την (21),  $u$  είναι το μοναδιαίο εφάλμα ετροχύλευσης και όπου υποθέτουμε ότι  $n^2 u \leq 1$ .

Απόδειξη: Η (35) δίνει ότι το  $\tilde{x}$  ικανοποιεί την

$$(LU+\delta LU+L\delta U+\delta L\delta U)\tilde{x}=b.$$

Αλλά, από το θεώρημα 1 έχουμε ότι  $LU=A+E$ . Συνεπώς η (36) ικανοποιείται με

$$\delta A=E+\delta LU+L\delta U+\delta L\delta U.$$

$$\text{Άρα } \|\delta A\|_{\infty} \leq \|E\|_{\infty} + \|\delta L\|_{\infty} \|U\|_{\infty} + \|L\|_{\infty} \|\delta U\|_{\infty} + \|\delta L\|_{\infty} \|\delta U\|_{\infty}.$$

Από την (20)  $\|E\|_{\infty} < n^2 \rho \|A\|_{\infty} u$ . Επίσης έχουμε ότι  $\|U\|_{\infty} = \max_i \sum_j |U_{ij}| \leq n \rho \|A\|_{\infty}$ ,  $\|L\|_{\infty} = \max_i \sum_j |L_{ij}| \leq n$ . Συνεπώς η (37) έπεται από τα παραπάνω, τις εκτιμήσεις (32), (34) και την υπόθεσή μας ότι  $n^2 u \leq 1$ . @

Συμπεραίνουμε ότι το αν η μέθοδος είναι ευσταθής, δηλ. το αν το  $\|\delta A\|_{\infty} / \|A\|_{\infty}$  είναι μικρό, θα εξαρτηθεί κυρίως από τον λεγόμενο "δυντε-

(κ)

λεστή μεγέθυνσης"  $\rho = \max_{i,j,k} |a_{ij}| / \|A\|_{\infty}$ . (Ο παράγοντας  $n^3+3n^2$

αφείλεται στην χρήση νορμώ πινάκων και σε χονδραιοδείξ εκτιμήσεις του φράγματος του  $\| \delta A \|_\infty$  στην απόδειξη του θεωρήματος 2 και δεν πρέπει να λαμβάνεται εσφαλρά υπ' όψιν - βλ. και Παρατήρηση 2 -). Δεν είναι δύσκολο να δούμε, επειδή  $|m_{ik}| \leq 1$ , ότι

$$\max_{i,j}^{(k+1)} |a_{ij}| = \max_{i,j}^{(k)} |a_{ij} - m_{ik} a_{ki}| \leq 2 \max_{i,j}^{(k)} |a_{ij}|$$

και συνεπώς ότι  $\max_{i,j,k}^{(k)} |a_{ij}| \leq 2^{n-1} \max_{i,j} |a_{ij}|$ . (Μάλιστα ο Wilkinson έχει δώσει παράδειγμα πίνακα για τον οποίο ισχύει η ισότητα (1)). Γενικά λοιπόν  $\rho \leq 2^{n-1}$ . Στην πράξη όμως τέτοια μελέθυση είναι εξαιρετικά επάνια. (Υπάρχει γενική ευμφωνία ότι εχεδόν πάντα για μερική οδήγηση (εμπειρικά)

$$\max_{i,j,k}^{(k)} |a_{ij}| \lesssim 10 \text{ αν } \max_{i,j} |a_{ij}| \leq 1).$$

Για ολική οδήγηση (για την οποία ισχύει ένα ανάλογο θεώρημα 2 και για την οποία εμφανίζεται πάλι ο παράγων  $\rho$  στην εκτίμηση του  $\| \delta A \|_\infty$ ), ο Wilkinson έδειξε ότι το φράγμα του  $\rho$  είναι της τάξης του  $n^{1 \ln n / 4}$  (δηλ. αυξάνεται πολύ αρχότερα από το  $2^n$ ) αν και δεν είναι χυωστά παραδείγματα όπου  $\rho \gg n$ .

Με βάση πολλά εμπειρικά δεδομένα και αριθμητικά πειράματα με μερική οδήγηση - παρά την ύπαρξη πινάκων για τους οποίους ο συντελεστής μελέθυσης αυξάνει εκθετικά με  $n$  - ευνήθως υποθέτουμε στην πράξη ότι

$$(38) \quad \frac{\| \delta A \|_\infty}{\| A \|_\infty} \approx \epsilon \beta^{-t},$$

όπου ευνήθως το  $\epsilon$  είναι της τάξης του  $\beta$ . Επί το ευτηρητικότερον, για μεγάλους πίνακες, υποθέτουμε ότι

$$(38') \quad \frac{\| \delta A \|_\infty}{\| A \|_\infty} \approx n \epsilon,$$

που επαυτίτητα δεν συμβαίνει. Στην πράξη λοιπόν δεχόμαστε ότι η απαλοιφή Gauss με μερική οδήγηση είναι συνήθως ((38)) ευσταθής ή το πολύ-πολύ ότι επιτρέπει μία ασθενή (ανάλογη του  $n$ ) μεγέθυνση των εφαλμάτων ((38')). Επαναλαμβάνουμε ότι αυτές οι εκτιμήσεις είναι εμπειρικές αλλά συνήθως ρεαλιστικές.

Τέλος, ας κλείσουμε την παράγραφο αυτή συνδιάζοντας τα της ευστάθειας του αλγορίθμου της απαλοιφής με μερική οδήγηση (θεώρημα 2) και τα αποτελέσματα της παρ. 1.2 για την ευαισθησία της λύσης του  $Ax=b$  σε διαταραχές για να εκτιμήσουμε το εφάλμα της προεχχιστικής λύσης  $\tilde{x}$ . Ας υποθέσουμε λοιπόν ότι η  $\tilde{x}$  ικανοποιεί την (36) και ας θέσουμε

$$(39) \quad \frac{\|\delta A\|}{\|A\|} = \mu,$$

για κάποια νόρμα  $\|\cdot\|$ . Οι (36), (39), (1.2.8), (1.2.9) δίνουν ( $\tilde{x}=x+\delta x$ ) ότι αν  $\kappa=\kappa(A)=\|A\| \|A^{-1}\|$ ,  $\|\delta A\| \|A^{-1}\| = \mu\kappa < 1$ , τότε

$$(40) \quad \|\tilde{x}-x\| / \|x\| \leq \mu\kappa / (1-\mu\kappa),$$

δηλ. ότι (αν π.χ.  $\mu$  μικρό) τότε το εχετικό εφάλμα της  $\tilde{x}$  (ως προς  $\|x\|$ ) μπορεί να είναι μεγάλο αν ο δείκτης κατάστασης  $\kappa$  είναι μεγάλος. Σημειώστε ότι για το υπόλοιπο  $r=A\tilde{x}-b$  έχουμε ότι  $r=A\tilde{x}-b = -\delta A\tilde{x}$  από την οποία έπεται ότι

$$(41) \quad \|r\| / (\|A\| \|\tilde{x}\|) \leq \mu,$$

δηλ. ότι το "εχετικό" υπόλοιπο (ως προς την ποσότητα  $\|A\| \|\tilde{x}\|$  που μετράει την "κλίμακα" του προβλήματος) είναι πάντα μικρότερο ή ίσο του  $\mu$ , όπου το  $\mu$  ορίστηκε από την (39). Ανακαλύπτας ότι για  $\|\cdot\|=\|\cdot\|_\infty$ , συνήθως  $\mu=O(u)$ , βλέπουμε ότι το υπόλοιπο της προεχχιστικής λύσης που δίνει η απαλοιφή Gauss με μερική οδήγηση είναι εχεδόν πάντα πολύ μικρό, ανεξάρτητα από την κατάσταση του συστήματος. Μάλιστα αυτό ισχύει και για μη αντιστρέψιμο  $A$  (!) - αρκεί βέβαια να υπάρχει λύση  $\tilde{x}$  του (36) -.



Τέλος, παρατηρώντας ότι  $\tilde{x}-x = A^{-1}r$ , βλέπουμε ότι  $\|\tilde{x}-x\| \leq \|A^{-1}\| \|r\|$  και ευνεπώς από την (41) ότι

$$(42) \quad \|x-\tilde{x}\|/\|\tilde{x}\| \leq \kappa\mu,$$

που πάλι δείχνει την σημασία του δείκτη κατάστασης στο σχετικό εφάλμα της  $\tilde{x}$ . Η (42) μάλιστα μας οδηγεί, αν ξέρουμε το  $\kappa$  και το  $\mu$ , στην εκτίμηση του απολύτου εφάλματος  $\|\tilde{x}-x\| \leq \|\tilde{x}\|\kappa\mu$  ευναρτήσει της υπολογιστικής λύσης  $\tilde{x}$  - ένα παράδειγμα εκτίμησης εκ των υετέρων (a posteriori), δηλ. με γνώση του αλγοριθμικού αποτελέσματος  $\tilde{x}$ , σε αντίθεση με την εκτίμηση του  $\|x-\tilde{x}\|$  που δίνει η (40) που είναι εκτίμηση εκ των προτέρων (a priori), δηλ. διατυπώνεται ευναρτήσει της ακριβούς λύσης  $x$ .

Εκτιμήσεις εφαλμάτων όπως οι ανισότητες (40) και (42) δείχνουν ότι, όσο αφορά εφάλματα της προεσχειτικής λύσης με απαλοιφή Gauss με μερική οδήγηση, ο δείκτης κατάστασης  $\kappa$  εμφανίζεται σαν όρος του γινομένου  $\kappa\mu$  όπου το  $\mu$  είναι της μορφής  $c_n u$  και όπου στην πράξη η σταθερά  $c_n$  δεν αυξάνεται γρήγορα με το  $n$  (βλ. (38), (38')). Έτσι π.χ. για ένα εύστημα με "μεγάλο" δείκτη κατάστασης ο αλγόριθμος είναι δυνατόν να δώσει λογικά αποτελέσματα αν κάνουμε τις πράξεις με αριθμητική διπλής ακρίβειας, οπότε  $u=O(\beta^{1-2t})$ . βεβαίως το κόστος των αριθμητικών πράξεων θα αυξηθεί.

#### Παρατηρήσεις

1. Αν τα στοιχεία του  $A$  και  $b$  δεν είναι όλα αριθμοί της μηχανής αλλά πραγματικοί αριθμοί μέσα στο εύρος των αριθμών της μηχανής, η (1) δείχνει ότι τα δεδομένα  $A, b$  θα παρασταθούν στον υπολογιστή από τα  $\tilde{A}, \tilde{b}$  (με στοιχεία αριθμούς μηχανής), όπου  $\tilde{b}=b+\delta b$ ,  $\|\delta b\|_\infty \leq u\|b\|_\infty$ ,  $\tilde{A}=A+\delta A$ ,  $\|\delta A\|_\infty \leq u\|A\|_\infty$ . Τέτοια "αρχικά" εφάλματα μπορούν εύκολα να ενσωματωθούν στην ανάλυση που κάναμε ε' αυτήν την παράγραφο. Ας υποθέσουμε όμως, χάριν παιδιάς, ότι δεν γίνονται άλλα εφάλματα ετρογχύλευσης κατά τις πράξεις της απαλοιφής πέρα από τα αρχικά αυτά εφάλματα παράστασης. Τότε η υπολογιστική λύση  $\tilde{x}$  θα είναι ακριβής

λύση του ευστήματος  $\tilde{A}\tilde{x}=\tilde{b}$ . Χρησιμοποιώντας τώρα τις πιο πάνω εκτιμήσεις για τα  $\|b\|_\infty$ ,  $\|bA\|_\infty$  και την θεωρία της παρ. 1.2 μπορούμε να δούμε (βλ. 'Αεκ. 4) ότι αν  $\kappa_\infty(A) \leq r < 1$ , τότε π.χ.

$\|x-\tilde{x}\|_\infty/\|x\|_\infty \leq 2r/(1-r)$ , δηλ. ότι αν ο  $r$  είναι κοντά στην μονάδα, τότε θα πρέπει να περιμένουμε μεγάλα εφάλματα για την  $\tilde{x}$  έστω και μόνο λόγω της προεχχιστικής παράστασης των αριθμών στον υπολογιστή.

2. Μιά έστω και πρόχειρη ματιά στην βιβλιογραφία (π.χ. συζητάτε τα αποτελέσματα αυτής της παραγράφου - που στηρίζονται στην ανάλυση του Wilkinson όπως παρουσιάζεται στο βιβλίο [1.2] των Forsythe και Moler - με ανάλογα αποτελέσματα των βιβλίων [0.2], [0.4], [0.5], [1.4], [1.7], [1.9] κλπ.) μας πείθει ότι ακόμη και για τον ίδιο ακριβώς αλγόριθμο είναι δυνατόν με διαφορετικές αποδείξεις ή διαφορετικές υποθέσεις για την αριθμητική κλπ. να πάρουμε διαφορετικές σταθερές  $c_n$  σε εκτιμήσεις π.χ. της μορφής (37) - στην (37)  $c_n = 1.01(n^3 + 3n^2)$  - ή ακόμα και διαφορετικού τύπου φράγματα. (Βέβαια, για άλλο αλγόριθμο ή διαφορετική σειρά των πράξεων κάτι τέτοιο θα συμβεί εύκολα). Γιαυτό συνήθως δεν δίνουμε και μεγάλη σημασία π.χ. σε τέτοιες σταθερές - στην πράξη όπως είδαμε είναι συνήθως ασφαλές να υποθέσουμε ότι  $c_n = O(n)$ , βλ. (38') -. Πάνω ε' αυτό το θέμα ο Wilkinson, βλ. [1.4, σελ. 36], λέει:

"Εξακολουθεί να υπάρχει μία τάση να αποδίδεται υπερβολική σημασία στις λεπτομέρειες της μορφής των φραγμάτων που δίνει μια a priori ανάλυση του εφάλματος. Κατά την γνώμη μου το φράγμα καθ'αυτό είναι συνήθως το λιγώτερο σημαντικό κομμάτι μιάς τέτοιας ανάλυσης, ο κύριος στόχος της οποίας είναι να αποκαλύψει πιθανούς μηχανισμούς αστάθειας του αλγορίθμου - αν υπάρχουν - έτσι ώστε ευδεχομένως να μπορέσουμε να βελτιώσουμε τον αλγόριθμο. Συνήθως το φράγμα είναι πιο αδύνατο απ' ότι θα ήταν αν δεν είμαστε υποχρεωμένοι να περιορίσουμε ε' ένα λογικό επίπεδο το πλήθος των λεπτομερειών της απόδειξης και δεν είχαμε τον περιορισμό να εκφράσουμε τα εφάλματα χρησιμοποιώντας νόρμες πινάκων. Τα a priori φράγματα δεν είναι

εν χέσει μεγέθη που πρέπει να χρησιμοποιούμε στην πράξη. Φράγματα που έχουν πρακτική αξία συνήθως προκύπτουν από κάποιας μορφής *a posteriori*-ανάλυση του σφάλματος [Σημ. Μετ: Δηλ. χρησιμοποιώντας ποσότητες που υπολογίζονται από τον αλγόριθμο]. Μιά τέτοια ανάλυση έχει το πλεονέκτημα ότι παίρνει υπ' όψη της την στατιστική κατανομή των σφαλμάτων ετροχύλευσης και τυχόν ειδικά χαρακτηριστικά του πίνακα όπως π.χ. την δομή των μηδενικών του".

3. Σε ορισμένες περιπτώσεις είναι δυνατόν πολλαπλασιάζοντας, πριν λύσουμε το σύστημα, από αριστερά και δεξιά του πίνακα  $A$  επί διαγωνίους πίνακες (που αντιστοιχεί σε πολλαπλασιασμό των γραμμών και των στηλών του  $A$  με σταθερές, δηλ. σε μία εκ των προτέρων αλλαγή κλίμακας (scaling) των στοιχείων των  $A, b, x$ ) να πάρουμε ένα σύστημα με μικρότερο δείκτη κατάστασης, βλ. π.χ. [1.2], [1.4]. Τέτοιοι διαγωνίοι πίνακες μπορούν να βρεθούν για ειδικές περιπτώσεις πίνακων  $A$  αλλά το γενικό πρόβλημα της κατασκευής τους δεν έχει διελευκασθεί ακόμα αρκετά.

4. Αν ο δείκτης κατάστασης ενός πίνακα δεν είναι πολύ μεγάλος εν ευκρίβει προς την ακρίβεια των πράξεων, δηλ. αν το γινόμενο  $\kappa(A)u$  είναι αρκετά μικρό, τότε είναι δυνατόν, χρησιμοποιώντας τεχνικές επαναληπτικής βελτίωσης (iterative improvement) να βελτιώσουμε την ακρίβεια της υπολογιστικής λύσης  $\tilde{x}$ . Μιά τέτοια τεχνική υπολογίζει το υπόλοιπο  $r$  της  $\tilde{x}$ , λύνει ένα νέο σύστημα  $Ay=r$  και βρίσκει μία "διόρθωση"  $y$ , η οποία, προστιθέμενη στην  $\tilde{x}$  δίνει μία νέα προσέγγιση που συνήθως είναι ακριβέστερη της  $\tilde{x}$ , εφ' όσον το υπόλοιπο  $r$  έχει υπολογισθεί με μεγαλύτερη ακρίβεια (π.χ. διπλή) απ' ό,τι οι υπόλοιπες πράξεις. Βλ. π.χ. [0.2], [1.2], [1.4].

5. Στον υπολογισμό στην πράξη του δείκτη κατάστασης  $\kappa(A) = \|A\| \|A^{-1}\|$  θέλουμε να αποφύγουμε τον υπολογισμό του  $A^{-1}$  κατά τα γνωστά. Συνήθως λοιπόν δεν υπολογίζουμε ακριβώς τον  $\kappa(A)$  αλλά κάνουμε μία προσεγγιστική του εκτίμηση: αυτό που χρειαζόμαστε είναι ουσιαστικά η τάξη μεγέθους του  $\kappa(A)$  ώστε να μπορούμε να εκτιμήσουμε την αξιοπιστία των αποτελεσμάτων του υπολογισμού. Για το πώς εκτι-

μούμε του  $\kappa(A)$  - με υπολογισμούς κόστους  $O(n^2)$  - βλ. π.χ. [5.4, σελ. 67].

### Ασκήσεις 1.3

1. (α) Βυμθείτε την απόδειξη της (1) και αποδείξτε με κατάλληλη παραλλαγή της την σχέση

$$(1') f(x) = x/(1+\delta'), \text{ όπου } |\delta'| \leq \alpha,$$

από την οποία έπεται η (2') αν κάνουμε τις ίδιες παραδοχές που οδήγησαν από την (1) στην (2).

(β) Αποδείξτε τους ισχυρισμούς στην απόδειξη του Λήμματος 1.

(γ) Υπό τις προϋποθέσεις του Λήμματος 2 δείξτε ότι

$$|f(\sum_{i=1}^n x_i y_i) - \sum_{i=1}^n x_i y_i| \leq 1.01 \text{ αν } \sum_{i=1}^n |x_i| |y_i|.$$

Συμπεώς αν π.χ.  $|\sum_{i=1}^n x_i y_i| \ll \sum_{i=1}^n |x_i| |y_i|$ , δηλ. αν π.χ. πολλά ζευγάρια όρων  $x_i y_i$  έχουν αντίθετα πρόσημα και αλληλοακυρώνονται, τότε είναι δυνατόν το σχετικό σφάλμα του

$f(\sum_{i=1}^n x_i y_i)$  να μην είναι μικρό. (Να ένα χρήσιμο συμπέρασμα από μία a priori-ανάλυση του σφάλματος ετρογχύλευσης που αναδεικνύει ένα πιθανό μηχανισμό αποσταθεροποίησης του αλγορίθμου για τον υπολογισμό του  $\sum_{i=1}^n x_i y_i$  βλ. παρατήρηση 2).

2. (α). Έστω  $A, B$   $n \times n$  πίνακες με στοιχεία αριθμούς μηχανής και έστω  $f(AB)$  το υπολογιστικό τους γινόμενο που παίρνουμε με τον συνηθισμένο τρόπο υπολογίζοντας τα εσωτερικά γινόμενα (γραμμή επί στήλη) με τον αλγόριθμο (4). Έστω  $E \in \mathbb{R}^{n \times n}$  τέτοιος ώστε  $f(AB) = AB + E$ . Βρείτε φράγμα για τον  $|E|$  συναρτήσει των  $|A|$ ,  $|B|$ ,  $n, u$ . (Παράδειγμα απ' ευθείας ανάλυσης του σφάλματος του πολλαπλασιασμού πινάκων). Πότε

είναι δυνατόν το  $f_1(AB)$  να έχει μεγάλο σχετικό εφάλμα;

(β) (Παράδειγμα αντίστροφης ανάλυσης του εφάλματος): 'Εστω  $A, B$   $2 \times 2$  άνω τριγωνικοί πίνακες. Βρείτε άνω τριγωνικούς  $2 \times 2$  πίνακες  $\tilde{A}, \tilde{B}$  (με στοιχεία "κουτά" ετά αντίστοιχα στοιχεία των  $A, B$ ) τέτοιους ώστε  $f_1(AB) = \tilde{A}\tilde{B}$  (ακριβές γινόμενο!). Εκτιμήστε τους πίνακες  $|\delta A|$ , αυτετχ.  $|\delta B|$ , όπου  $\delta A = \tilde{A} - A$ ,  $\delta B = \tilde{B} - B$ , συνάρτησει των  $n, \mu$  και  $|A|$ , αυτετχ.  $|B|$ .

3. Αν αντί του αλγορίθμου (29) χρησιμοποιηθεί ο (30) για την επίλυση του κάτω τριγωνικού συστήματος  $Ly=b$ , βρείτε το αντίστοιχο  $\delta L$  (ώστε να ισχύει η (31')) και εκτιμήστε τον πίνακα  $|\delta L|$  και την νόρμα  $\|\delta L\|_\infty$ .

4. (α) Αποδείξτε τους ισχυρισμούς στην παρατήρηση 1.

(β) Στην περίπτωση που τα στοιχεία του  $A$  είναι αριθμοί μηχανής (δηλ. όταν  $\delta A=0$ ) και το μόνο εφάλμα ετρογχύλευσης προέρχεται από την παράσταση του  $b$ , δείξτε ότι η εκτίμηση του ερωτήματος (α) απλουετείται σε  $\|x-\tilde{x}\|_\infty/\|x\|_\infty \leq r$ . Εφαρμογή: 'Εστω ότι τα στοιχεία των  $A, b$  είναι ακέραιοι με εξαίρεση το  $b_1$  που είναι ίσο με  $10^{-1}$ . Παριστάνουμε τα στοιχεία των  $A, b$  στον VAX 11/780 και υποθέτουμε ότι οι υπόλοιπες πράξεις της απαλοιφής γίνονται ακριβώς. Ποιό είναι το αναμενόμενο  $\|x-\tilde{x}\|_\infty/\|x\|_\infty$  αν  $\kappa_\infty(A)=10^4$ ; Το ίδιο ερώτημα για απλή και διπλή ακρίβεια στον IBM 4361. (Η αριθμητική του 4361 είναι ίδια με την αριθμητική του συστήματος IBM 370. Για τις τιμές των  $\beta, t$  -βλ. [5.4, εελ. 18]).

### 1.4 Η ΑΝΑΛΥΣΗ CHOLESKY ΓΙΑ ΣΥΜΜΕΤΡΙΚΟΥΣ, ΘΕΤΙΚΑ ΟΡΙΣΜΕΝΟΥΣ ΠΙΝΑΚΕΣ

Στην παράγραφο αυτή θα ασχοληθούμε με την επίλυση γραμμικών συστημάτων  $Ax=b$ ,  $b \in \mathbb{R}^n$ , όπου ο  $A$  είναι ένας  $n \times n$  πραγματικός, συμμετρικός ( $A=A^T$ ) και θετικά ορισμένος (δηλ. με την ιδιότητα ότι  $x^T Ax > 0$  για κάθε  $0 \neq x \in \mathbb{R}^n$ ) πίνακας. Ένας θετικά ορισμένος πίνακας είναι προφανώς αντιστρέψιμος. Τέτοιοι πίνακες εμφανίζονται συχνά στις εφαρμογές. Θα εξετάσουμε μία άμεση μέθοδο για την λύση του γραμμικού συστήματος, που μπορεί να θεωρηθεί σαν ειδική μορφή της ανάλυσης LU στην περίπτωση μας, την λεγόμενη ανάλυση Cholesky.

Παίρνοντας  $x = e^i \in \mathbb{R}^n$  ( $e_i = 1$ ,  $e_j = 0$ ,  $i \neq j$ ) βλέπουμε ότι  $x^T Ax = a_{ii} > 0$ . γενικότερα κάθε κύριος υποπίνακας (βλ. θεκ. 1) ενός συμμετρικού θετικά ορισμένου πίνακα είναι επίσης συμμετρικός και θετικά ορισμένος. Από αυτήν την παρατήρηση (βλ. θεκ. 3) προκύπτει ότι συστήματα με συμμετρικούς και θετικά ορισμένους πίνακες μπορούν να επιλυθούν με

(i)

απλή απαλοιφή Gauss χωρίς εναλλαγές γραμμών· μάλιστα οι οδηγοί  $a_{ii}$  είναι όλοι θετικοί αριθμοί. Ο αλγόριθμος όμως που θα χρησιμοποιήσουμε εδώ θα είναι μία παραλλαγή της απαλοιφής και στηρίζεται στο εξής.

**ΘΕΩΡΗΜΑ 1** (Ανάλυση Cholesky). Έστω  $A \in \mathbb{R}^{n \times n}$  συμμετρικός και θετικά ορισμένος πίνακας. Τότε υπάρχει μοναδικός κάτω τριγωνικός πίνακας  $L$  με θετικά διαγώνια στοιχεία (όχι αναγκαστικά μονάδες), τέτοιος ώστε να ισχύει η ανάλυση Cholesky:

$$(1) \quad A = LL^T$$

Απόδειξη: Με επαγωγή. Το θεώρημα προφανώς ισχύει για  $1 \times 1$  θετικά ορισμένους πίνακες:  $a_{11} > 0$  και  $L_{11} = (a_{11})^{1/2}$ . Έστω ότι το θεώρημα ισχύει για  $(n-1) \times (n-1)$  συμμετρικούς, θετικά ορισμένους πίνακες και έστω  $A$  ένας  $n \times n$  τέτοιος πίνακας. Χωρίζουμε τον  $A$  σε υποπίνακες ως εξής:

$$A = \begin{pmatrix} d & u^T \\ u & H \end{pmatrix}$$

όπου  $d = a_{11} > 0$ ,  $u$  ένα (στηλο)διάνυσμα  $(n-1) \times 1$  και  $\tilde{H}$   $(n-1) \times (n-1)$  πίνακας. Μπορούμε να επαληθεύσουμε ότι ο  $A$  γράφεται τότε σαν γινόμενο ως εξής:

$$(2) \quad A = \begin{pmatrix} d^{1/2} & 0 \\ u/d^{1/2} & I_{n-1} \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & H \end{pmatrix} \begin{pmatrix} d^{1/2} & u^T/d^{1/2} \\ 0 & I_{n-1} \end{pmatrix},$$

όπου  $H = \tilde{H} - uu^T/d$  και  $I_{n-1}$  είναι ο  $(n-1) \times (n-1)$  μοναδιαίος πίνακας. Ο πίνακας  $H$  είναι συμμετρικός και θετικά ορισμένος, γιατί για κάθε  $x \in \mathbb{R}^{n-1}$ ,  $x \neq 0$ :

$$x^T H x = x^T (\tilde{H} - uu^T/d) x = y^T \begin{pmatrix} d & u^T \\ u & \tilde{H} \end{pmatrix} y = y^T A y > 0,$$

όπου το  $y \in \mathbb{R}^n$  δίνεται από

$$y = \begin{pmatrix} -x^T u/d \\ x \end{pmatrix}.$$

Από την υπόθεση της επαγωγής τώρα ο  $(n-1) \times (n-1)$  συμμετρικός και θετικά ορισμένος

Τ

θετικά ορισμένος πίνακας  $H$  γράφεται στη μορφή  $H = L_H L_H^T$ , όπου  $L_H$  κάτω τριγωνικός με θετικά διαγώνια στοιχεία. Συνεπώς, λόγω της (2) ο  $A$  ικανοποιεί

$$\begin{aligned}
 A &= \begin{pmatrix} d^{1/2} & 0 \\ u/d^{1/2} & I_{n-1} \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & L_H \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & L_H^T \end{pmatrix} \begin{pmatrix} d^{1/2} & u^T/d^{1/2} \\ 0 & I_{n-1} \end{pmatrix} = \\
 &= \begin{pmatrix} d^{1/2} & 0 \\ u/d^{1/2} & L_H \end{pmatrix} \begin{pmatrix} d^{1/2} & u^T/d^{1/2} \\ 0 & L_H^T \end{pmatrix} = LL^T
 \end{aligned}$$

Το επαγωγικό βήμα τελείωσε. Για να δείξουμε την μοναδικότητα του  $L$  υποθέτουμε ότι υπάρχει κάτω τριγωνικός πίνακας  $M$  με θετικά διαγώνια στοιχεία τέτοιος ώστε  $A=LL^T=MM^T$ . Άρα  $L^{-1}M=L^T(M^T)^{-1}$ . Επειδή  $L^{-1}, M$  κάτω τριγωνικοί,  $L^T, (M^T)^{-1}$  άνω τριγωνικοί, έπεται ότι  $L^{-1}M=L^T(M^T)^{-1}=D$  όπου  $D$  είναι διαγώνιος πίνακας. Τώρα  $L^{-1}M=0 \Rightarrow D_{ii}=M_{ii}/L_{ii}$  και  $L^T(M^T)^{-1}=0 \Rightarrow D_{ii}=L_{ii}/M_{ii}$ . Άρα  $(L_{ii})^2=(M_{ii})^2 \Rightarrow L_{ii}=M_{ii}$ , γιατί  $L_{ii}, M_{ii} > 0$ . Συνεπώς  $D=I_n \Leftrightarrow L=M$  ό.έ.δ. @

Για να βρούμε τώρα έναν αποτελεσματικό αλγόριθμο για την κατασκευή των στοιχείων  $L_{ij}$  του  $L$  συνάρτησε των  $a_{ij}$  εξισώνουμε ένα προς ένα τα στοιχεία του (κάτω τριγώνου του)  $A$  με τα αντίστοιχα στοιχεία του γινομένου  $LL^T$ . Το στοιχείο  $(1,1)$  δίνει  $a_{11}=(L_{11})^2 \Rightarrow L_{11}=(a_{11})^{1/2}$ . Τα στοιχεία της δεύτερης γραμμής δίνουν  $a_{21}=L_{21} \cdot L_{11} \Rightarrow L_{21}=a_{21}/L_{11}$  και  $a_{22}=(L_{21})^2+(L_{22})^2 \Rightarrow L_{22}=(a_{22}-(L_{21})^2)^{1/2}$ .

Γενικά, εξισώνοντας τα στοιχεία της  $i$ -ετής γραμμής του  $A$  με τα αντίστοιχα του  $LL^T$  παίρνουμε τις ισότητες

$$L_{ij} = (a_{ij} - \sum_{k=1}^{j-1} L_{ik} L_{jk}) / L_{jj} \quad \text{για } j=1, 2, \dots, i-1$$

και

$$L_{ii} = (a_{ii} - \sum_{k=1}^{i-1} (L_{ik})^2)^{1/2}$$





εξαρτάται από το  $n$  αλλά αυξάνει αργά με το  $n$ , και  $u$  το μοναδιαίο εφάλμα ετρογχύλευσης. Βλέπουμε δηλ. ότι η ανάλυση Cholesky είναι ευεταθής αλγόριθμος - παρατηρείστε ότι ο "ευσταθής μεγέθυνσης" των εφαλμάτων είναι  $p=1$ . (Τα αποτελέσματα αυτά έχουν αποδειχθεί και δεν αποτελούν εμπειρικές εκτιμήσεις). Επιπλέον, αν  $q_n$  ή  $\kappa_2(A) < 1$ , όπου  $q_n$  είναι άλλη μία μικρή σταθερά εξαρτώμενη από το  $n$ , ο Wilkinson δείχνει ότι ο αλγόριθμος τερματίζεται ομαλά, δηλ. ότι δεν εμφανίζεται ποτέ στους υπολογισμούς λόγω εφάλματος ετρογχύλευσης ένα  $L_{ii}$  μηδέν ή φανταστικό. Υπάρχουν δηλ. περιπτώσεις θεωρητικά θετικά ορισμένων πινάκων που είναι "εχεδόν" μη αντιστρέψιμοι, δηλ. που έχουν  $\kappa_2(A) \gg 1$ , για τους οποίους, αν η ακρίβεια στις πράξεις δεν είναι επαρκής (δηλ. αν το  $u$  δεν είναι αρκετά μικρό), είναι δυνατόν ο αλγόριθμος της ανάλυσης του Cholesky να μην ολοκληρωθεί (βλ. θεκ. 5)

#### Παρατηρήσεις

1. Ένας συμμετρικός και θετικά ορισμένος πίνακας  $A$  μπορεί να αναλυθεί και σε γινόμενο παραγόντων της μορφής:

$$A = MDM^T$$

όπου  $M$  κάτω τριγωνικός με μονάδες στην διαγώνιο και  $D$  διαγώνιος με θετικά στοιχεία. Η ανάλυση αυτή (ανάλυση Crout) προκύπτει π.χ. από την ανάλυση Cholesky αν  $M = LD^{-1}$  όπου  $\Delta$  διαγώνιος πίνακας με  $\Delta_{ii} = L_{ii}$  και  $D = \Delta \Delta^T$ . Προφανώς μπορούμε να βρούμε κατ' ευθείαν αλγόριθμο, ανάλογο του (3) για τον υπολογισμό των  $M_{ij}$ ,  $i > j$  και των  $D_{ii}$ .

2. Υπάρχουν και διάφοροι άλλοι τρόποι εκτός από του (3) για την κατασκευή των στοιχείων του πίνακα  $L$ . Π.χ. εξισώνοντας στην  $A = LL^T$  τα στοιχεία κατά ετήλες παίρνουμε τον αλγόριθμο της "ανάλυσης Cholesky κατά ετήλες" κατά του οποίου πρώτα υπολογίζουμε τα στοιχεία  $L_{11}$ ,  $i \geq 1$ , μετά τα  $L_{12}$ ,  $i \geq 2$  κ.ο.κ. Μία άλλη μορφή του αλγορίθμου, ή μορφή γινόμενου πινάκων, στηρίζεται σε μία κατασκευαστική μορφή της απόδειξης του θεωρήματος 1. Η προτίμηση για τον ένα ή τον άλλο αλγόριθμο υπαγορεύεται π.χ. από το πώς είναι αποθηκευμένος ο  $A$ , από

του  $A$  και του  $L$  (π.χ. αν είναι αραιοί) κ.τ.λ. Για περισσότερες λεπτομέρειες βλ. [5.4, παρ. 5,6].

3. Φυσικά λόγω του θεωρήματος 1 ύπαρξης-μοναδικότητας του  $L$ , είμαστε εκ των προτέρων βέβαιοι ότι στον αλγόριθμο (3) για κάθε  $i$ ,

$$a_{ii} > \sum_{k=1}^{i-1} (L_{ik})^2, \text{ δηλ. ότι το } L_{ii} \text{ που θα προκύψει θα είναι}$$

πραγματικό και θετικό. (Αυτό γιατί  $A = LL^T \Rightarrow a_{ii} = L_{ii}^2 + \sum_{k=1}^{i-1} (L_{ik})^2$  και  $L_{ii} > 0$ ). Συνεπώς, αν έχουμε ένα (πραγματικό) συμμετρικό πίνακα  $A$  και θέλουμε να ελέγξουμε υπολογιστικά αν είναι θετικά ορισμένος, εκτελούμε τον αλγόριθμο του Cholesky. Αν ο αλγόριθμος ολοκληρωθεί ομαλά, δηλ. αν  $L_{ii} > 0, 1 \leq i \leq n$ , τότε κατασκευάσαμε κάτω τριγωνικό πίνακα  $L$  με  $L_{ii} > 0$  τέτοιου ώστε  $A=LL^T$ . Συνεπώς (θεκ. 2γ) ο  $A$  είναι θετικά ορισμένος. (Αν όμως όπως είδαμε παραπάνω ο πίνακας  $A$  έχει μεγάλο δείκτη κατάστασης σχετικά με την ακρίβεια των πράξεων τότε είναι δυνατόν, λόγω εσφαλμάτων ετρογχύλευσης, ο αλγόριθμος στην πράξη να μην τερματισθεί, βλ. θεκ. 5).

4. Στα κεφάλαια 6-10 των [5.4] - που ακολουθούν πιστά το βιβλίο [1.3] των George και Liu - γίνεται μία λεπτομερής μελέτη της μεθόδου Cholesky για αραιούς, μεγάλους, συμμετρικούς και θετικά ορισμένους πίνακες  $A$ . Συστήματα με τέτοιους πίνακες εμφανίζονται επιχυστά στις εφαρμογές. Στην περίπτωση ενός μεγάλου αραιού πίνακα  $A$  (του οποίου δηλ. τα περισσότερα στοιχεία είναι μηδέν) μας ενδιαφέρει να αναδιατάξουμε τις γραμμές και τις ετήλες του  $A$ , να αποθηκεύσουμε με κατάλληλες δομές δεδομένων τα (μη μηδενικά) στοιχεία του και να χρησιμοποιήσουμε μία κατάλληλη μορφή του αλγορίθμου Cholesky έτσι ώστε να ελαχιστοποιήσουμε κατά το δυνατόν το "χέμπερα" του  $L$  (δηλ. το πλήθος των μη μηδενικών στοιχείων  $L_{ij}$  του  $L$  για τα οποία τα αντίστοιχα στοιχεία  $a_{ij}$  του  $A$  ήταν μηδέν), το κόστος αποθήκευσης του  $A$  και  $L$  καθώς και το πλήθος των πράξεων. Τα προβλήματα που εμφανίζονται είναι εξαιρετικά ενδιαφέροντα και από μαθηματική άποψη

(χρησιμοποιούνται οι μέθοδοι της θεωρίας γραφημάτων) και από άποψη προγραμματισμού, αριθμητικής ανάλυσης και θεωρίας δομών δεδομένων.

5. Τέλος, μία σημείωση ιστορικού ενδιαφέροντος: ο André-Louis Cholesky (1875-1918) ήταν Γάλλος αξιωματικός του μηχανικού που έκανε γεωδαιτικές και τοπογραφικές μετρήσεις στην ΚΡΗΤΗ και στην Βόρεια Αφρική πριν από τον Α' Παγκόσμιο Πόλεμο. Ανακάλυψε την ομώνυμη μέθοδο για να υπολογίζει λύσεις γραμμικών συστημάτων που προκύπτουν από τις λεγόμενες "κανονικές εξισώσεις" της μεθόδου ελαχίστου τετραγώνων που εφαρμόζεται για την προσέγγιση δεδομένων σε γεωδαιτικά προβλήματα. Η μέθοδος του δημοσιεύθηκε μετά θάνατον στο Bulletin Geodesique το 1924.

#### Ασκήσεις 1.4

1. Έστω  $A=(a_{ij})$  ένας  $n \times n$  συμμετρικός και θετικά ορισμένος πίνακας. Να δείξει ότι

(α) ο  $A^{-1}$  είναι επίσης θετικά ορισμένος.

(β) Κάθε τετραγωνικός υποπίνακας του  $A$  του οποίου η κύρια διαγώνιος βρίσκεται πάνω στην κύρια διαγώνιο του  $A$  (δηλ. κάθε πίνακας με στοιχεία  $k \leq i, j \leq m, 1 \leq k \leq m \leq n$ ) είναι θετικά ορισμένος.

(γ) Για κάθε  $k, 1 \leq k \leq n$ , ισχύει ότι

$$\max_{1 \leq i \leq k} a_{ii} = \max_{1 \leq i, j \leq k} |a_{ij}|.$$

2. (α) Δείξτε ότι ο τριδιαγώνιος συμμετρικός πίνακας  $A$  με  $a_{ii} = \lambda, a_{i, i \pm 1} = 1, \lambda \geq 2$  είναι θετικά ορισμένος.

(β) Ένας συμμετρικός πραγματικός πίνακας είναι θετικά ορισμένος αν και μόνο αν οι ιδιοτιμές του είναι θετικές.

(γ) Έστω  $L$  πραγματικός, κάτω τριγωνικός πίνακας με θετικά διαγώνια στοιχεία. Τότε ο πίνακας  $A=LL^T$  είναι συμμετρικός και θετικά ορισμένος.

(δ) Αν ο  $A$  είναι συμμετρικός και θετικά ορισμένος τότε ο πίνακας  $B^T A B$  έχει τις ίδιες ιδιότητες αν και μόνο αν ο  $B$  είναι αντιστρέψιμος.

3. Έστω  $A \in \mathbb{R}^{n \times n}$ , συμμετρικός και θετικά ορισμένος. Δείξτε ότι η τριγωνοποίησή του κατά την απαλοιφή Gauss μπορεί να γίνει χωρίς

(1)

εναλλαγές γραμμών και μάλιστα ότι όλοι οι οδηγοί  $a_{ii}$  είναι θετικοί

(1)

(2)

(Υπόδειξη: Προφανώς  $a_{ii} = a_{ii} > 0$ . Θεωρείστε τον υποπίνακα  $\tilde{A} = (a_{ij})_{2 \leq i, j \leq n}$ ,

που προκύπτει από τον  $A^{(2)}$  αν αφαιρέσουμε την πρώτη γραμμή και την πρώτη στήλη του. Δείξτε ότι ο  $\tilde{A}$  είναι συμμετρικός και θετικά

(2)

ορισμένος. Συνεπώς  $a_{22} > 0$ . Η επαγωγική συνέχεια της απόδειξης προφανής).

4. (α) Δείξτε ότι ο αλγόριθμος (3) απαιτεί  $(n^3 + 3n^2 - 4n)/6$  πράξεις (= πολ/εμούς και διαιρέσεις) και η τετραγωνικές ρίζες.

(β) Έστω ότι ο μιγαδικός πίνακας  $H = A + iB$  (όπου  $A, B \in \mathbb{R}^{n \times n}$ ) είναι αυτοσυζυγής ( $H = H^*$ ) και θετικά ορισμένος ως μιγαδικός πίνακας (δηλ. ισχύει  $(z, Hz)_2 > 0$  για κάθε  $0 \neq z \in \mathbb{C}^n$ ). Δείξτε ότι ο  $2n \times 2n$  πραγματικός πίνακας

$$C = \begin{pmatrix} A & -B \\ B & A \end{pmatrix}$$

είναι συμμετρικός και θετικά ορισμένος. Βρείτε έναν αλγόριθμο με  $4n^3/3 + O(n^2)$  (πραγματικές) πράξεις για την λύση του συστήματος  $H(x + iy) = (b + ic)$  όπου  $x, y, b, c \in \mathbb{R}^n$ . Πόση μνήμη απαιτείται;

5. Θεωρείστε τον πίνακα

$$A = \begin{pmatrix} 100 & .15 & .01 \\ .15 & 2.3 & .01 \\ .01 & .01 & 1 \end{pmatrix}$$

Δείξτε ότι με αριθμητική απεριόριστης ακρίβειας, η μέθοδος του Cholesky τερματίζεται κανονικά (= απόδειξη ότι ο  $A$  είναι θετικά ορισμένος). Δείξτε όμως ότι αν οι πράξεις γίνουν με  $\beta=10$ ,  $t=2$  και ετροχύλευση, τότε οι υπολογιστικές τιμές των στοιχείων του  $L$  είναι  $\tilde{L}_{11}=10$ ,  $\tilde{L}_{21}=1.5$ ,  $\tilde{L}_{31}=10^{-3}$ ,  $\tilde{L}_{22}=0.0$ , δηλ. ότι ο αλγόριθμος σταματά επιχειρώντας να υπολογίσει το  $\tilde{L}_{32}$ .

## 1.5. ΜΕΘΟΔΟΙ ΕΛΑΧΙΣΤΟΠΟΙΗΣΗΣ

Μία εμμαντική κατηγορία μεθόδων για την επίλυση γραμμικών ευστημάτων (κυρίως με συμμετρικούς, θετικά ορισμένους πίνακες) είναι μέθοδοι ελαχιστοποίησης καταλλήλων συναρτησιακών  $\varphi(x)$ ,  $x \in \mathbb{R}^n$ . Αν το  $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$  είναι ομαλό, τότε από τον Απειροστικό Λογισμό, ξέρουμε ότι τα τοπικά του ελάχιστα πρέπει να αναζητηθούν μέσα στο σύνολο των κριτίμων σημείων του  $\varphi$ , δηλ. μεταξύ των σημείων  $x$  όπου  $\nabla \varphi(x) = 0$ , όπου ως γνωστόν η κλίση  $\nabla \varphi$  του  $\varphi$  είναι το διάνυσμα

$$\nabla \varphi = \left( \frac{\partial \varphi}{\partial x_1}, \frac{\partial \varphi}{\partial x_2}, \dots, \frac{\partial \varphi}{\partial x_n} \right)^T.$$

Επίσης ξέρουμε ότι αν  $\nabla \varphi(x) \neq 0$ , τότε η συνάρτηση  $g(\alpha) = \varphi(x + \alpha u)$ ,  $x, u \in \mathbb{R}^n$ ,  $\alpha \in \mathbb{R}$ , που περιγράφει την συμπεριφορά του  $\varphi$  κοντά στο σημείο  $x$  κατά την κατεύθυνση  $u$ , ελαττώνεται με μέγιστο ρυθμό ελάττωσης όταν  $u = -\nabla \varphi(x)$  λόγω του ότι  $g'(0) = (\nabla \varphi(x), u)_2$ .

Η τελευταία παρατήρηση μας οδηγεί σε μία προεχχιστική επαναληπτική μέθοδο για τον υπολογισμό των ελαχίστων  $x^*$  της  $\varphi(x)$ , την λεγόμενη μέθοδο της "καθόδου μεγίστης κλίσεως" (steepest descent) του Cauchy, η οποία ανάγει το πρόβλημα σε επανειλημμένους υπολογισμούς ελαχίστων συναρτήσεων μίας μεταβλητής. Η μέθοδος παράγει μία ακολουθία  $\{x^j\}$ ,  $j \geq 0$  προεχχίσεων ενός (τοπικού) ελαχίστου  $x^*$  της  $\varphi(x)$  ως εξής: Έστω  $x^k$  μία προεχχιση. Η  $x^{k+1}$  ορίζεται ως το πλησιέστερο στο  $x^k$  ελάχιστο της  $\varphi(x)$  όταν το  $x$  περιορίζεται πάνω στην ακτίνα που διέρχεται από το  $x^k$  και έχει την διεύθυνση του  $-\nabla \varphi(x^k)$ . Βρίσκουμε δηλ. (γενικά προεχχιστικά) τον μικρότερο θετικό αριθμό  $\alpha_k$  όπου η συνάρτηση  $g(\alpha) = \varphi(x^k - \alpha \nabla \varphi(x^k))$  έχει ελάχιστο και ορίζουμε  $x^{k+1} = x^k - \alpha_k \nabla \varphi(x^k)$ . Συνεπώς η τακτική μας σε κάθε βήμα  $k$  του αλγορίθμου είναι να ελαχιστοποιήσουμε την  $\varphi(x)$  κοντά στο  $x^k$  πάνω στη διεύθυνση  $-\nabla \varphi(x^k)$  κατά την οποία τοπικά η  $\varphi(x)$  ελαττώνεται γρηγορότερα.

Θεωρούμε τώρα το πχη πραγματικό ευστημα με το ελαχιστοποιοεισ προβλημα

$$(1) \quad Ax = b$$

όπου ο  $A$  είναι συμμετρικός και θετικά ορισμένος. Συμβολίζοντας με  $(\cdot, \cdot) = (\cdot, \cdot)_2$  το ευκλείδειο εσωτερικό γινόμενο στον  $\mathbb{R}^n$  και με  $\|\cdot\| = \|\cdot\|_2$  την αντίστοιχη νόρμα, θεωρούμε το πρόβλημα ελαχιστοποίησης για  $x \in \mathbb{R}^n$  του ευαρτησιακού

$$(2) \quad \varphi(x) = (Ax, x)/2 - (b, x).$$

Έστω  $z = A^{-1}b$  η λύση του (1). Τότε επειδή για κάθε  $y \in \mathbb{R}^n$  ισχύει  $\varphi(z+y) = \varphi(z) + (Ay, y)/2$ , συμπεραίνουμε ότι  $\varphi(z+y) > \varphi(z) \quad \forall 0 \neq y \in \mathbb{R}^n$ , δηλ. ότι το  $\varphi(x)$  παίρνει την ελάχιστη τιμή του στον  $\mathbb{R}^n$  (ίση με  $-(b, z)/2$ ) στο μοναδικό σημείο  $x = z = A^{-1}b$ . Συνεπώς το πρόβλημα της επίλυσης του συστήματος (1) είναι ισοδύναμο με το πρόβλημα ελαχιστοποίησης χωρίς περιορισμούς

$$(3) \quad \min_{x \in \mathbb{R}^n} \varphi(x),$$

που έχει μοναδική λύση. Εξ άλλου επειδή η κλίση του  $\varphi$  είναι ίση με

$$(4) \quad \nabla \varphi(x) = Ax - b,$$

το μοναδικό κρίσιμο σημείο του  $\varphi$  είναι το σημείο για το οποίο  $\nabla \varphi(x) = 0$ , δηλ. το  $A^{-1}b$  στο οποίο το  $\varphi$  έχει ελάχιστο επειδή είναι "αυστηρά κυρτό".

Στο σημείο  $x^k$  η συνάρτηση  $\varphi(x)$  ελαττώνεται ταχύτερα στην κατεύθυνση  $-\nabla \varphi(x^k) = b - Ax^k$ , την οποία ταυτίζουμε με το υπόλοιπο  $r^k$  του  $x^k$ , δηλ. θέτουμε

$$(5) \quad r^k = -\nabla \varphi(x^k) = b - Ax^k.$$

Αν  $r^k \neq 0$  (αλλιώς  $x^k = A^{-1}b$ ), η μέθοδος της καθόδου μεγίστης κλίσεως υπολογίζει συνεπώς το  $x^{k+1}$  ελαχιστοποιώντας το διάνυσμα ως προς  $a$

$$\varphi(x^k + ar^k) = (Ar^k, r^k)a^2/2 - (r^k, r^k)a + \varphi(x^k)$$



του οποίου το ελάχιστο λαμβάνεται για  $a = a_* > 0$  όπου

$$a_* = (r^k, r^k) / (Ar^k, r^k).$$

Συμπεώς αν  $r^k \neq 0$  θα έχουμε  $\varphi(x^k + a_* r^k) < \varphi(x^k)$ . Πάλι είναι δυνατό να υπολογιστεί το ελάχιστο του  $\varphi$  στην ευθεία που περνάει από το  $x^k$  και παράλληλη προς  $r^k$ .

$$(6) \quad \varphi(x^k + a_* r^k) = \varphi(x^k) - (r^k, r^k)^2 / 2 (Ar^k, r^k)$$

Οδηγούμαστε λοιπόν στον εξής αλγόριθμο της μεθόδου της καθόδου μεγίστης κλίσεως για την ελαχιστοποίηση του (2):

$$(7) \quad \left[ \begin{array}{l} x^0 = 0 \\ \text{Για } k=1, 2, \dots \\ \quad r^{k-1} = b - Ax^{k-1} \\ \quad \text{Αν } r^{k-1} = 0 \\ \quad \quad \text{τότε τέλος, } x = x^{k-1}. \\ \quad \text{αλλιώς} \\ \quad \quad a_k = (r^{k-1}, r^{k-1}) / (Ar^{k-1}, r^{k-1}) \\ \quad \quad x^k = x^{k-1} + a_k r^{k-1} \end{array} \right.$$

Ο αλγόριθμος τερματίζεται όταν ικανοποιηθούν ένα ή περισσότερα από τα ευρήθη κριτήρια τερματισμού επαναληπτικών μεθόδων (βλ. [5.4, βελ. 159], [1.8], [1.11]). Προφανώς, η αρχική τιμή  $x^0 = 0$  δεν είναι δεσμευτική.

Θα μελετήσουμε την σύγκλιση της μεθόδου στην "φωτεινή" νόρμα του προβλήματος δηλ. στην νόρμα  $x \mapsto (Ax, x)^{1/2}$  που παράγεται από το εσωτερικό γινόμενο  $(Ax, y)$  στον  $\mathbb{R}^n$  (βλ. θεκ. 2(a)). Αποδεικνύουμε το εξής θεώρημα σύγκλισης που δίνει επίσης ένα μέτρο του πόσο γρήγορα συγκλίνει η ακολουθία  $\{x^j\}$  στην νόρμα  $(A \cdot, \cdot)^{1/2}$ .

**ΘΕΩΡΗΜΑ 1.** Έστω  $\{x^j\}$ ,  $j \geq 0$  η ακολουθία που παράγει ο αλγόριθμος (7) της μεθόδου της καθόδου μεγίστης κλίσεως για οποιοδήποτε  $x^0 \in \mathbb{R}^n$  και έστω  $x$  η λύση του συστήματος (1). Έστω  $\kappa = \lambda_{\max} / \lambda_{\min}$ , όπου  $\lambda_{\max}$ , αντιστ.  $\lambda_{\min}$ , η μέγιστη, αντιστ. ελάχιστη, ιδιοτιμή του  $A$ . Τότε  $x^j \rightarrow x$ ,  $j \rightarrow \infty$ . Πέραν, αν  $e^j = x - x^j$  είναι το εφάλμα της προσέγγισης  $x^j$ , ισχύει ότι

$$(8) \quad (Ae^j, e^j)^{1/2} \leq [(\kappa-1)/(\kappa+1)]^j (Ae^0, e^0)^{1/2}, \quad j=0, 1, 2, \dots$$

Απόδειξη. Έχουμε  $Ae^j = Ax - Ax^j = b - Ax^j = r^j$ . Επίσης χρησιμοποιώντας τον ορισμό του  $x^{j+1}$  έχουμε  $e^{j+1} = x - x^{j+1} = x - x^j + x^j - x^{j+1} = e^j - a_{j+1} r^j$ . Συνεπώς  $Ae^{j+1} = Ae^j - a_{j+1} Ar^j$  από την οποία, χρησιμοποιώντας και τον ορισμό του  $a_{j+1}$  από τον (7)

$$(Ae^{j+1}, r^j) = (Ae^j, r^j) - a_{j+1} (Ar^j, r^j) = (r^j, r^j) - a_{j+1} (Ar^j, r^j) = 0.$$

Οι πιο πάνω σχέσεις δίνουν τώρα για οποιοδήποτε  $a \in \mathbb{R}$

$$\begin{aligned} (Ae^{j+1}, e^{j+1}) &= (Ae^{j+1}, e^j - a_{j+1} r^j) = (Ae^{j+1}, e^j) - a_{j+1} (Ae^{j+1}, r^j) = \\ &= (Ae^{j+1}, e^j - ar^j). \end{aligned}$$

Χρησιμοποιώντας τώρα την ανισότητα των Cauchy-Schwarz για το εσωτερικό γινόμενο  $(A, \cdot)$  έχουμε

$$(Ae^{j+1}, e^{j+1}) \leq (Ae^{j+1}, e^{j+1})^{1/2} (A(e^j - ar^j), e^j - ar^j)^{1/2}, \quad \forall a \in \mathbb{R}.$$

Τέλος, από την σχέση  $Ae^j = r^j$ , έχουμε  $e^j - ar^j = (1-aA)e^j$ , και  $A(e^j - ar^j) = A(1-aA)e^j = (1-aA)Ae^j$ . Συμπεραίνουμε ότι

$$(9) \quad (Ae^{j+1}, e^{j+1}) \leq \inf_{a \in \mathbb{R}} ((1-aA)Ae^j, (1-aA)e^j), \quad j \geq 0.$$

Χρησιμοποιούμε τώρα την λεγόμενη "φασματική παράσταση" του  $A$ . Ο  $A$ , αν συμμετρικός πίνακας, έχει η πραγματικές ιδιοτιμές (θετικές

γιατί ο  $A$  είναι θετικά ορισμένος)  $0 < \lambda_1 \equiv \lambda_{\min} \leq \lambda_2 \leq \dots \leq \lambda_n \equiv \lambda_{\max}$  και  $n$  αντίστοιχα, ορθοκανονικά ως προς το  $(\cdot, \cdot)$ , ιδιοδιανύσματα  $u^j$ ,  $1 \leq j \leq n$ . Έχουμε δηλ.  $Au^j = \lambda_j u^j$ ,  $1 \leq j \leq n$ ,  $(u^i, u^j) = \delta_{ij}$ . Άρα, για κάθε  $u \in \mathbb{R}^n$  ισχύει

$u = \sum_{j=1}^n (u, u^j) u^j$ , και συνεπώς  $Au = \sum_{j=1}^n \lambda_j (u, u^j) u^j$ . Γενικότερα, για κάθε πραγματικό πολυώνυμο  $p$ , η τελευταία σχέση δίνει

$p(A)u = \sum_{j=1}^n p(\lambda_j) (u, u^j) u^j$ ,  $\forall u \in \mathbb{R}^n$ . Συνεπώς για κάθε  $\alpha \in \mathbb{R}$  έχουμε

$$\begin{aligned} ((1-\alpha A)Ae^j, (1-\alpha A)e^j) &= \sum_{i=1}^n (1-\alpha \lambda_i)^2 \lambda_i (e^j, u^i)^2 \leq \\ &\leq \max_{1 \leq i \leq n} (1-\alpha \lambda_i)^2 \sum_{i=1}^n \lambda_i (e^j, u^i)^2 = \max_{1 \leq i \leq n} (1-\alpha \lambda_i)^2 (Ae^j, e^j). \end{aligned}$$

Συμπεραίνουμε, επειδή  $\lambda_j \in [\lambda_{\min}, \lambda_{\max}]$ , ότι λόγω της (9)

$$(10) \quad (Ae^{j+1}, e^{j+1})^{1/2} \leq \inf_{\alpha \in \mathbb{R}} (\max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |1-\alpha \lambda|) (Ae^j, e^j)^{1/2}, \quad j \geq 0.$$

Το "min-max" πρόβλημα του υπολογισμού του

$$\epsilon = \inf_{\alpha \in \mathbb{R}} (\max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |1-\alpha \lambda|)$$

μπορεί βέβαια να λυθεί ως ειδική περίπτωση ενός γενικότερου προβλήματος που θα αναλύσουμε στην παρ. 3.6 χρησιμοποιώντας πολυώνυμο Chebyshev με την παρατήρηση ότι

$$\epsilon = \inf_{p \in \Pi_1} \max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |p(\lambda)| \quad \text{όπου } \Pi_1 \text{ το σύνολο των πραγματικών γραμμικών πολυνομίων } p(\lambda) \text{ τέτοιων ώστε } p(0)=1. \text{ Λύεται όμως και στοιχειωδώς ως εξής: Προφανώς } \epsilon = \inf_{\alpha \in \mathbb{R}} [\max (|1-\alpha \lambda_{\min}|, |1-\alpha \lambda_{\max}|)].$$

Στρέφοντας στο επίπεδο  $(\lambda, \mu)$  την ευθεία  $\mu=1-\alpha\lambda$  περί το σημείο  $(0,1)$  καθώς το  $\alpha$  μεταβάλλεται στο  $\mathbb{R}$ , παρατηρούμε ότι η συνάρτηση

$\alpha \mapsto \max(|1-\alpha \lambda_{\min}|, |1-\alpha \lambda_{\max}|)$  παίρνει την ελάχιστη τιμή της όταν

$1 - a\lambda_{\min} > 0$ ,  $1 - a\lambda_{\max} < 0$  και  $|1 - a\lambda_{\min}| = |1 - a\lambda_{\max}|$  δηλ. όταν

$a = 2/(\lambda_{\min} + \lambda_{\max})$  έχουμε λοιπόν ότι  $\epsilon = |1 - 2\lambda_{\min}/(\lambda_{\min} + \lambda_{\max})| = (\kappa - 1)/(\kappa + 1)$ .

Η (10) δίνει τότε

$$(\theta^j) (Re^{j+1}, e^{j+1})^{1/2} \leq [(\kappa - 1)/(\kappa + 1)] (Re^j, e^j)^{1/2}, \quad j \geq 0$$

από την οποία προκύπτει αμέσως η  $(\theta)$  και η εύκλιση  $e^j \rightarrow 0, j \rightarrow \infty$ . @

Η άσκηση 4 δείχνει ότι η ταχύτητα εύκλισης, δηλ. ο "λόγος"  $\rho = (\kappa - 1)/(\kappa + 1)$  της "γεωμετρικής" εύκλισης της ακολουθίας  $e^j$  στην νόρμα  $(A^j, \cdot)^{1/2}$ , είναι η καλύτερη δυνατή. Αλλά μπορεί να γίνει πολύ μικρή αν ο λόγος  $\kappa = \lambda_{\max}/\lambda_{\min}$  είναι μεγάλος (δηλ. αν ο δείκτης κατάστασης  $\kappa = \kappa_2(A)$  του  $A$  είναι μεγάλος - φυσικά!). Ενδιαφέρον έχει εδώ η γεωμετρική ερμηνεία της μεθόδου. (Βλκ. 5). Οι "ισοψείς" επιφάνειες  $\varphi(x) = \text{σταθ.}$  του συναρτησιακού  $\varphi$  είναι ελλειψοειδή στον  $R^n$  με κέντρο το σημείο  $z = A^{-1}b$  και άξονες παράλληλους προς τα ορθογώνια ιδιοδιανύσματα  $u^i$  του  $A$  με μήκη αξόνων ανάλογα προς τους αριθμούς  $\lambda_i^{-1/2}$ , όπου  $\lambda_i$  οι ιδιοτιμές του  $A$ . (Αλλάξτε ευτεταχμένες με νέο εύστημα αξόνων κέντρου  $z$  και άξονες παράλληλους πρό τα  $u^i$ ). Από τον ορισμό του  $\nabla \varphi(x^k)$ , έπεται ότι η προσέγγιση  $x^{k+1}$  βρίσκεται πάνω στην κάθετο στο  $x^k$  της επιφάνειας του ελλειψοειδούς που διέρχεται από το σημείο  $x^k$ . Τα διανύσματα  $x^{k+1} - x^k$  και  $x^k - x^{k-1}$  είναι εξ' άλλου ορθογώνια (Βλκ. 1(γ)). Αν  $\lambda_{\max} \gg \lambda_{\min}$  τα ελλειψοειδή είναι υπερβολικά "στενόμακρα" και η κάθοδος προς το κέντρο μιάς βάρδιας χαράδρας με απότομες πλαγιές γίνεται με ζίγκ-ζάγκ στις διαδοχικές ορθογώνιες κατευθύνσεις  $r^k = -\nabla \varphi(x^k)$  μεταξύ των πλαγιών πράγμα που καθυστερεί πολύ την κάθοδο (βλ. και Βλκ. 4). Έτσι η τακτική της ελαχιστοποίησης τοπικά κατά μήκος των διευθύνσεων των υπολοίπων  $r^k$  (δηλ. των διευθύνσεων καθόδου μεγίστης κλίσεως) αποδεικνύεται κακή στρατηγική.

Επιχειρούμε λοιπόν τώρα να ελαχιστοποιήσουμε το  $\varphi$  διαδοχικά σε κατευθύνσεις  $p^1, p^2, \dots$  (που δεν ευρίσκονται αναγκαστικά με τα υπόλοιπα  $r^0, r^1, r^2, \dots$ ) ελπίζοντας ότι η επιλογή καταλλήλων  $p^i$  θα βελτιώσει την ταχύτητα σύγκλισης. Προκαταρκτικά, γενικεύοντας λίγο υπολογισμούς που κάναμε προηγουμένως, βλέπουμε ότι, για δεδομένα  $x^{k-1}$  και  $p^k \neq 0$ , η συνάρτηση  $\varphi(x^{k-1} + \alpha p^k)$  ελαχιστοποιείται όταν

$$(11) \quad \alpha = \alpha_k = (p^k, r^{k-1}) / (A p^k, p^k),$$

οπότε, αν  $x^k = x^{k-1} + \alpha_k p^k$ , θα ισχύει  $\varphi(x^k) < \varphi(x^{k-1})$  εφ' όσον  $(p^k, r^{k-1}) \neq 0$ .

Πάλι στα,

$$(12) \quad \varphi(x^k) = \varphi(x^{k-1}) - (p^k, r^{k-1})^2 / 2 (A p^k, p^k).$$

(Η μέθοδος της καθόδου μεγίστης κλίσεως επιλέγει  $p^k = r^{k-1}$ ). Ο αλγόριθμος (7) γενικεύεται λοιπόν ως εξής:

$$(13) \quad \left[ \begin{array}{l} x^0 = 0 \\ \text{Για } k=1, 2, \dots \\ r^{k-1} = b - A x^{k-1} \\ \text{Αν } r^{k-1} = 0 \\ \text{τότε τέλος, } x = x^{k-1} \\ \text{αλλιώς} \\ \text{διαλέξε } p^k \text{ τέτοιο ώστε } (p^k, r^{k-1}) \neq 0 \\ \alpha_k = (p^k, r^{k-1}) / (A p^k, p^k) \\ x^k = x^{k-1} + \alpha_k p^k. \end{array} \right.$$

Το ερώτημα βέβαια είναι πως θα διαλέξουμε τα διανύσματα  $p^i$  έτσι ώστε ο αλγόριθμος (13) να συγκλίνει και να μην έχει τα μειονεκτήματα της μεθόδου της καθόδου μεγίστης κλίσεως. Παρατηρούμε, λόγω της αναδρομικής σχέσης  $x^k = x^{k-1} + \alpha_k p^k$ , ότι για κάθε  $k$  το διάνυσμα  $x^k$  που

δίνει ο αλγόριθμος (13) δίδεται από ένα γραμμικό συνδυασμό  $x^k = a_1 p^1 + \dots + a_k p^k$  των  $p^i$ ,  $1 \leq i \leq k$ . Θα ήταν λοιπόν ιδανικό π.χ. να διαλέγουμε τα διανύσματα  $p^i$  έτσι ώστε να είναι γραμμικά ανεξάρτητα και έτσι ώστε το  $x^k$  να λύνει το πρόβλημα ελαχιστοποίησης:

$$(14) \quad \min_{x \in \langle p^1, \dots, p^k \rangle} \varphi(x)$$

όπου με  $\langle a^1, \dots, a^N \rangle$ ,  $a^i \in \mathbb{R}^n$ , συμβολίζουμε τον υπόχωρο του  $\mathbb{R}^n$  που παράγεται από τα  $a^i$ ,  $1 \leq i \leq N$ , δηλ. το εύλογο των διανυσμάτων που είναι γραμμικοί συνδυασμοί των  $a^i$ ,  $1 \leq i \leq N$ . (Το πρόβλημα ελαχιστοποίησης του  $\varphi(x)$  για  $x \in M$ , όπου  $M$  υπόχωρος του  $\mathbb{R}^n$ , έχει μοναδική λύση: βλ. θεκ. 6).

Με μιά τέτοια κατασκευή των  $x^k$ , ο αλγόριθμος (13) θα ενέκλινε σε  $n$  βήματα στη λύση του συστήματος  $Ax=b$ . Πράγματι το  $x^n$  θα ελαχιστοποιούσε το  $\varphi(x)$  για  $x \in \langle p^1, \dots, p^n \rangle = \mathbb{R}^n \Rightarrow x^n = A^{-1}b$ . Βέβαια, από την κατασκευή των  $a_k, x^k$  ο αλγόριθμος (13) παράγει για  $k=1, 2, \dots$  ένα  $x^k$  που λύνει κάθε φορά το μονοδιάστατο πρόβλημα ελαχιστοποίησης.

$$(15) \quad \min_{a \in \mathbb{R}} \varphi(x^{k-1} + ap^k).$$

Είναι δυνατόν να διαλέξουμε τα  $p^k$  έτσι ώστε το  $x^k$ , η λύση του (15), να είναι ευχρόνως και λύση του (14);

Η απάντηση είναι θετική. Ας κάνουμε μιά ανάλυση του προβλήματος επαγωγικά. θεωρούμε τον  $n \times j$  πίνακα  $P_j = [p^1, \dots, p^j]$  με στήλες  $p^1, \dots, p^j$ . Ένα διάνυσμα  $z \in \mathbb{R}^n$  ανήκει στο πεδίο τιμών  $\mathcal{R}(P_j)$  του πίνακα  $P_j$  αν και μόνο αν  $z \in \langle p^1, \dots, p^j \rangle$ . Αν  $x \in \mathcal{R}(P_k)$ , τότε το  $x$  μπορεί να γραφτεί στην μορφή  $x = P_{k-1}y + ap^k$  για κάποιο  $y \in \mathbb{R}^{k-1}$  και  $a \in \mathbb{R}$ . Έχουμε επίσης τότε

$$(16) \quad \varphi(x) = \varphi(P_{k-1}y) + [a^2(Ap^k, p^k)/2 - a(p^k, b)] + a(P_{k-1}y, Ap^k).$$

Αν ο τελευταίος όρος του δευτέρου μέλους της (16) ήταν μηδέν, τότε το πρόβλημα της ελαχιστοποίησης του  $\varphi(x)$  για  $x \in \mathcal{R}(P_k)$  θα αναχόταν: (α) στην ελαχιστοποίηση του  $\varphi(P_{k-1}y)$  για  $y \in \mathbb{R}^{k-1}$ , δηλ. στην ελαχιστοποίηση του  $\varphi$  πάνω στον υπόχωρο  $\langle p^1, \dots, p^{k-1} \rangle$  - πρόβλημα που η λύση του  $x^{k-1}$  υποτίθεται ότι έχει βρεθεί και (β) στην ελαχιστοποίηση του δευτέρου όρου του δευτέρου μέλους της (16), δηλ. ε' ένα πρόβλημα που λύνεται φυσικά όταν

$$(17) \quad a = a_k \equiv (p^k, b) / (Ap^k, p^k).$$

Ένας προφανής τρόπος για να μηδενίσουμε τον τελευταίο όρο της (16) είναι να υπολογίσουμε το  $p^k$ , με δεδομένα τα  $p^i$ ,  $1 \leq i \leq k-1$ , έτσι ώστε

$$(18) \quad P_{k-1}^T Ap^k = 0.$$

Αν ισχύει η (18), το  $x^{k-1}$  ελαχιστοποιεί το  $\varphi$  πάνω στον υπόχωρο  $\langle p^1, \dots, p^{k-1} \rangle$  και το  $a$  δίνεται από την (17), συμπεραίνουμε ότι το πρόβλημα (14) λύνεται, για  $x = x^k \equiv x^{k-1} + a_k p^k$ . Επειδή  $x^{k-1} \in \langle p^1, \dots, p^{k-1} \rangle$  έχουμε εξ' άλλου ότι  $x^{k-1} = P_{k-1} z$  για κάποιο  $z \in \mathbb{R}^{k-1}$ . Άρα από την (18)

προκύπτει ότι  $(p^k, Ax^{k-1}) = (Ap^k, P_{k-1} z) = z^T P_{k-1}^T Ap^k = 0$ , δηλ. ότι το  $a_k$  που δίδεται από την (17) ικανοποιεί και την σχέση (19) για τον αριθμό

$$(19) \quad a_k = (p^k, b - Ax^{k-1}) / (Ap^k, p^k) = (p^k, r^{k-1}) / (Ap^k, p^k),$$

δηλ. λύνει και το πρόβλημα της μονοδιάστατης ελαχιστοποίησης (15).

Συνοψίζουμε: αν το  $p^k$  επιλεγεί έτσι ώστε να ισχύει η (18), δηλ. έτσι ώστε να ισχύουν οι σχέσεις

$$(20) \quad (Ap^j, p^k) = 0, \quad j=1, 2, \dots, k-1,$$

τότε ο αριθμός (13) κατασκευάζει  $a_k$  και  $x^k$  που λύνουν τα προβλή-

ματα (14) και (15). Αν η (20) ισχύει, τότε λέμε ότι το  $p^k$  είναι A-ευζυγές (ή A-ορθογώνιο) προς τα  $p^1, \dots, p^{k-1}$ .

Απομένει να εξετάσουμε τρία ζητήματα:

- (α) Είναι τα  $p^i$  γραμμικά ανεξάρτητα;
- (β) Ισχύει ότι  $(p^k, r^{k-1}) \neq 0$ , έτσι ώστε, βλ. (12), να έχουμε ελάττωση του  $\varphi$  στο βήμα  $k$ ;
- (γ) Πώς υπολογίζουμε στην πράξη τα  $p^i$ ;

Η απάντηση στο ερώτημα (α) είναι καταφατική: Αν τα διανύσματα  $p^i \neq 0$ ,  $1 \leq i \leq m$  του  $\mathbb{R}^n$  είναι A-ευζυγή, δηλ. αν

$$(21) \quad (Ap^i, p^j) = 0, \text{ αν } i \neq j,$$

τότε τα  $p^i$ ,  $1 \leq i \leq m$  είναι γραμμικά ανεξάρτητα. (θεμ. 7(β)). Συνεπώς, αν κατασκευάσουμε τα  $p^k \neq 0$  για  $k=1, 2, \dots, n$  έτσι ώστε να ισχύει η (20), τότε θα είναι γραμμικά ανεξάρτητα, πράγμα που εγγυάται ότι ο αλγόριθμος δίνει την ακριβή λύση του προβλήματος σε  $n$  βήματα, δηλ. ότι  $x^n = A^{-1}b$ .

Στο ερώτημα (β), αν για κάποιο  $k$ ,  $r^{k-1} = 0$ , τότε ο αλγόριθμος τερματίζει και έχουμε ότι  $x^{k-1} = x = A^{-1}b$ . Αν  $r^{k-1} \neq 0$ , τότε υπάρχει  $p^k$  που ικανοποιεί την (20) τέτοιο ώστε  $(p^k, r^{k-1}) \neq 0$ . Πράγματι, αν για κάθε  $p \in \mathbb{R}^n$  που είναι A-ευζυγές προς τα  $p^i$ ,  $1 \leq i \leq k-1$ , ισχύει  $(p, r^{k-1}) = 0$ , θα έχουμε, επειδή  $r^{k-1} = b - Ar^{k-1}$  και  $x^{k-1} \in \langle p^1, \dots, p^{k-1} \rangle$ , ότι  $(p, b) - (p, Ar^{k-1}) = 0$  για κάποιο  $z \in \mathbb{R}^{k-1}$ , δηλ. ότι  $(p, b) = 0$  για κάθε  $p$  A-ευζυγές προς τα  $p^i$ ,  $1 \leq i \leq k-1$ . (Χρησιμοποιήσαμε την (18)). Συνεπώς  $(p, b) = 0$  για κάθε  $p \in \langle Ap^1, Ap^2, \dots, Ap^{k-1} \rangle^{\perp} \Leftrightarrow b \in \langle Ap^1, \dots, Ap^{k-1} \rangle \Leftrightarrow A^{-1}b \in \langle p^1, \dots, p^{k-1} \rangle \Leftrightarrow x^{k-1} = x = A^{-1}b \Leftrightarrow r^{k-1} = 0$ , άτοπο. Συνεπώς μπορούμε να βρούμε  $p^k$ , A-ευζυγές προς τα  $p^i$ ,  $1 \leq i \leq k-1$ , τέτοιο ώστε  $(p^k, r^{k-1}) \neq 0$ , αν βέβαια  $r^{k-1} \neq 0$ .

Όπως παρατηρήσαμε και παραπάνω "p A-ευζυγές προς τα  $p^i$ ,  $1 \leq i \leq k-1$ "  $\Leftrightarrow p \in \langle Ap^1, \dots, Ap^{k-1} \rangle^{\perp}$ . Η λεγόμενη μέθοδος των ευζυγών κλίσεων (conjugate gradients) των Hestenes και Stiefel (1952) επιλέγει ως  $p^k$  το πλησιέστερο προς το  $r^{k-1}$  διάνυσμα του υπόχωρου



$\langle Ar^1, \dots, Ar^{k-1} \rangle^\perp$ , δηλ. κατασκευάζει τα  $x^k$  με βάση την εξής εξειδίκευση του γενικού προγράμματος (13) για τον υπολογισμό της λύσης  $x$  του  $Ax=b$ :

$$(22) \quad \left. \begin{array}{l} x^0=0 \\ \text{Γιά } k=1, 2, \dots, n \\ r^{k-1}=b-Ax^{k-1} \\ \text{Αν } r^{k-1}=0 \\ \text{τότε, εκχώρησε } x=x^{k-1} \text{ και τερμάτισε.} \\ \text{αλλιώς, όρισε} \\ r^k = \begin{cases} r^0, & \text{αν } k=1 \\ \text{ορθή προβολή του } r^{k-1} \text{ στον υπόχωρο} \\ \langle Ar^1, \dots, Ar^{k-1} \rangle^\perp, & \text{αν } k>1. \end{cases} \\ a_k = (r^k, r^{k-1}) / (Ar^k, r^k) \\ x^k = x^{k-1} + a_k r^k \end{array} \right\} x=x^n$$

Ουσιαστικά όμως δεν απαντήσαμε στο ερώτημα (γ), δηλ. στο πώς υπολογίζουμε τα  $r^k$ ,  $k>1$  με αποτελεσματικό τρόπο στην πράξη. Με το πρόβλημα αυτό καθώς και με το πρόβλημα της εκτίμησης του εφάλματος  $e^j = x - x^j$  για κάθε  $j$ , θ' ασχοληθούμε στην επόμενη παράγραφο.

### Ασκήσεις 1.5

1. (α) Επαληθεύστε την (δ).

(β) Βρείτε αναδρομική σχέση μεταξύ των  $r^k, r^{k-1}$  που μας επιτρέπει να χράφουμε τον αλγόριθμο (7) έτσι ώστε να απαιτεί ένα μόνο πολλαπλασιασμό πίνακα επί διάνυσμα για κάθε  $k$ . Πάρες πράξεις απαιτεί ο νέος αλγόριθμος ανά βήμα  $k$ ;

(γ) Δείξτε ότι δύο διαδοχικά υπόλοιπα της μεθόδου της καθόδου μέγιστης κλίσεως για την λύση του  $Ax=b$  είναι ορθογώνια, δηλ. ότι  $(r^j, r^{j-1})=0$ ,  $j \geq 1$ . Επίσης ότι τα διανύσματα  $x^{j+1}-x^j$ ,  $x^j-x^{j-1}$ ,  $j \geq 1$  είναι ορθογώνια.

2. (α) Δείξτε ότι η παράσταση  $(Ax, y)$ ,  $x, y \in \mathbb{R}^n$  ορίζει ένα εσωτερικό γινόμενο στο  $\mathbb{R}^n$  και ότι συνεπώς η συνάρτηση  $x \mapsto (Ax, x)^{1/2}$  είναι νόρμα. Βρείτε τις (καλύτερες) σταθερές σύγκρισης μεταξύ αυτής της νόρμας και της  $\|x\|_2$ .

(β) Δείξτε ότι η  $(\delta')$  είναι ισοδύναμη με την  $\varphi(x^{j+1}) + (b, x)/2 \leq ((\kappa-1)/(\kappa+1))^2 [\varphi(x^j) + (b, x)/2]$ .

3. (α) Χρησιμοποιώντας τον συμβολισμό του θεωρήματος 1, δείξτε ότι για  $j \geq 0$

$$(Ae^{j+1}, e^{j+1}) = \{1 - [(r^j, r^j)^2 / (Ar^j, r^j)(A^{-1}r^j, r^j)]\} (Ae^j, e^j)$$

(β) Δείξτε την ανισότητα του Καντοβιτς:

$$(Ax, x)(A^{-1}x, x) \leq [(\lambda_{\max} + \lambda_{\min})^2 / 4\lambda_{\max}\lambda_{\min}] \|x\|^4$$

(γ) Χρησιμοποιώντας την ανισότητα του Καντοβιτς και την ταυτότητα του ερωτήματος (α) δείξτε την  $(\delta')$  δίνοντας έτσι μία άλλη απόδειξη του θεωρήματος 1.

4. Θεωρούμε το  $2 \times 2$  σύστημα  $Ax=b$  όπου  $A = \begin{pmatrix} 1 & 0 \\ 0 & \lambda \end{pmatrix}$  με  $\lambda > 0$  και  $x=b=0$ .

(α) Αν  $x^k = (x_1, x_2)^T$ , δείξτε ότι η μέθοδος της καθόδου μέγιστης κλίσεως δίνει

$$x^{k+1} = [x_1 x_2 (\lambda-1) / (x_1^2 + \lambda^3 x_2^2)] \begin{pmatrix} \lambda^2 x_2 \\ -x_1 \end{pmatrix}$$

Σχεδιάστε για  $\lambda=5$  τις ισοψείς καμπύλες του συναρτησιακού  $\varphi(x)$  (δηλ. τις καμπύλες  $\varphi(x_1, x_2)=\text{σταθ}$ ) και μερικές διαδοχικές προεγγύσεις  $x^k$ ,  $k=1, 2, 3, \dots$  με  $x^0=(5, 1)^T$ .

(β) Δείξτε ότι αν  $x^j = c(\lambda, \pm 1)^T$ , τότε η ανισότητα (8') ισχύει ως ισότητα.

(γ) Αν  $\lambda = 100$ , πόσα βήματα της μεθόδου θα χρειαστούν έτσι ώστε το εφάλμα  $(Re^j, e^j)^{1/2}$  να γίνει μικρότερο από  $\epsilon (Re^0, e^0)^{1/2}$  όπου  $\epsilon > 0$  δεδομένο;

5. Αποδείξτε τους ισχυρισμούς περί γεωμετρικής ερμηνείας της μεθόδου της καθόδου μεγίστης κλίσης που περιλαμβάνονται στο κείμενο (Σελ 1.5.6) από "Οι "ισοψείς" επιφάνειες ..." μέχρι "... του ελλειψοειδούς που διέρχεται από το σημείο  $x^k$ ".

6. Έστω  $\Pi$  ένας υπόχωρος του  $\mathbb{R}^n$  και έστω  $y \in \mathbb{R}^n$  ένα δεδομένο διάνυσμα. Δείξτε ότι το πρόβλημα της ελαχιστοποίησης του  $\varphi(x)$  για  $x \in \Pi$  έχει μοναδική λύση.

7. (α) Επαληθεύστε τις (11), (12) και δείξτε ότι  $(r^j, p^j) = 0$ ,  $j \geq 1$ .

(β) Έστω ότι  $0 \neq p^i \in \mathbb{R}^n$ ,  $1 \leq i \leq m$  και ότι ισχύει η (21), δηλ. ότι τα  $p^i$  είναι  $\mathbb{R}$ -επιζυγή. Δείξτε ότι είναι γραμμικά ανεξάρτητα.

## 1.6 Η ΜΕΘΟΔΟΣ ΤΩΝ ΣΥΖΥΓΩΝ ΚΛΙΣΕΩΝ

Ο σκοπός μας ε' αυτήν την παράγραφο θα είναι να μελετήσουμε λεπτομερώς τις ιδιότητες των  $x^k, r^k, p^k$  που κατασκευάζει ο αλγόριθμος (1.5.22), να εκτιμήσουμε τα σφάλματα  $e^k = x - x^k$  και να βρούμε πρακτικούς και αποτελεσματικούς τρόπους για να υπολογίζουμε σε κάθε βήμα  $k$  την νέα κατεύθυνση ελαχιστοποίησης  $p^k$  που ορίζεται στον (1.5.22) ως η ορθή προβολή της κλίσεως (υπολοίπου)  $r^{k-1}$  πάνω στον υπόχωρο  $\langle Rr^1, \dots, Rr^{k-1} \rangle$ .

Θα αρχίσουμε αποδεικνύοντας μία σειρά χρήσιμων ιδιοτήτων των υπολοίπων  $r^i$  και των συζυγών κατευθύνσεων  $p^i$ . Λέμε ότι "ο αλγόριθμος (1.5.22) ολοκληρώνει  $k$  βήματα" αν  $r^i \neq 0$ ,  $0 \leq i \leq k-1$ . Τότε τα  $x^i$ ,  $1 \leq i \leq k$  θα έχουν υπολογισθεί επαγωγικά από τις εκθέσεις:  $x^0 = 0$  και

$$(1) \quad x^i = x^{i-1} + a_i p^i, \quad i \geq 1.$$

**Λήμμα 1.** Αν ο αλγόριθμος (1.5.22) ολοκληρώνει  $k$  βήματα, τότε για  $i=1, 2, \dots, k$  έχουμε

$$(2) \quad r^i = r^{i-1} - a_i p^i,$$

$$(3) \quad p^T r^i = 0,$$

όπου  $P_i$  είναι ο  $n \times i$  πίνακας  $[p^1, \dots, p^i]$  με στήλες  $p^1, \dots, p^i$ .

Απόδειξη: Η (2) είναι προφανής απόρροια της (1). Η (3) μπορεί να προκύψει "αλγεβρικά" από την (1) και το γεγονός ότι τα  $p^i$  είναι  $R$ -συζυγή (βλ. Άσκηση 1.(α)) αλλά και ως εξής: Γνωρίζουμε ότι το πρόβλημα της ελαχιστοποίησης

$$(4) \quad \min_{x \in \langle p^1, \dots, p^i \rangle} \varphi(x).$$

λύεται μοναδικά για  $x = x^i$ . Εξ άλλου, επειδή  $x \in \langle p^1, \dots, p^i \rangle \Leftrightarrow x \in \mathcal{R}(P_i)$ ,

το (4) είναι ισοδύναμο με το πρόβλημα

$$(5) \quad \min \varphi(P, y), \\ y \in \mathbb{R}^i$$

το οποίο, επειδή  $\varphi(P, y) = (AP, y, P, y) / 2 - (b, P, y)$ , λύνεται (γιατί); για

$y = y^i$ , όπου  $y^i$  είναι η λύση του συστήματος  $(P_i^T AP_i^T) y^i = P_i^T b$ . Άρα από την σχέση

$$x^i = P_i y^i \text{ έπεται η } P_i^T A x^i = P_i^T b \Leftrightarrow P_i^T r^i = 0. @$$

**Λήμμα 2.** Για  $k \geq 2$ , τα διανύσματα  $p^k$  που παράγονται από τον αλγόριθμο (1.5.22) (εφ' όσον ολοκληρώνει  $k$  βήματα) προεδιορίζονται από τις σχέσεις

$$(6) \quad p^k = r^{k-1} - AP_{k-1} z^{k-1},$$

όπου το  $z^{k-1} \in \mathbb{R}^{k-1}$  είναι η (μοναδική) λύση του προβλήματος ελαχιστοποίησης

$$(7) \quad \min_{z \in \mathbb{R}^{k-1}} \|r^{k-1} - AP_{k-1} z\|.$$

Απόδειξη: Έξ' ορισμού το  $p^k$  είναι η (ορθή) προβολή του  $r^{k-1}$  στον υπόχωρο  $\mathcal{R}(AP_{k-1})^\perp$ . Συνεπώς το  $r^{k-1} - p^k$  είναι η ορθή προβολή του  $r^{k-1}$  στον υπόχωρο  $\mathcal{R}(AP_{k-1})$ . (Για επανάληψη περί προβολών βλ. Άσκηση 2). Συνεπώς  $r^{k-1} - p^k = AP_{k-1} z^{k-1}$  για  $z^{k-1} \in \mathbb{R}^{k-1}$ . Η μοναδικότητα του  $z^{k-1}$  είναι τώρα απόρροια της γραμμικής ανεξαρτησίας των  $p^i$ . Επιπλέον το  $z^{k-1}$  λύνει το πρόβλημα της ελαχιστοποίησης (ελαχίστων τετραγώνων)

$$\begin{aligned}
 (\|r^{k-1} - (r^{k-1} - p^k)\| &= \|p^k\| =) \|r^{k-1} - AP_{k-1}z^{k-1}\| = \\
 &= \min_{z \in \mathbb{R}^{k-1}} \|r^{k-1} - AP_{k-1}z\|. \textcircled{*} \\
 u &\in \mathcal{R}(AP_{k-1})
 \end{aligned}$$

Θα χρησιμοποιήσουμε τα παραπάνω δύο λήμματα για να αποδείξουμε τα εξής εμφαντικό θεωρητικό αποτέλεσμα:

**ΠΡΟΤΑΣΗ 1** Αν ο αλγόριθμος (1.5.22) ολοκληρώνει  $k$  βήματα, τότε:

(α) Τα υπόλοιπα  $r^0, r^1, \dots, r^{k-1}$  είναι ορθογώνια μεταξύ τους.

(β) Για  $j=1, 2, \dots, k$  ισχύει

$$(8) \langle r^1, \dots, r^j \rangle = \langle r^0, \dots, r^{j-1} \rangle = \langle b, Ab, \dots, A^{j-1}b \rangle.$$

Απόδειξη:

(α) Από την (2) για  $1 \leq i \leq k-1$  έχουμε ότι  $Ap^i = (r_i - r_{i-1})/a_i$  (γιατί  $a_i \neq 0$  ;). Άρα  $Ap^i \in \langle r^0, \dots, r^i \rangle$ . Συνεπώς η (6) δίνει ότι  $p^j = r^{j-1} - [Ap^1, \dots, Ap^{j-1}]z^{j-1} \in \langle r^0, \dots, r^{j-1} \rangle$  για  $1 \leq j \leq k$ . Έπεται λοιπόν ότι για  $1 \leq j \leq k$   $p^j = c_{j0}r^0 + \dots + c_{j,j-1}r^{j-1}$ , όπου  $c_{j,j-1} \neq 0$  γιατί τα  $p^j$  είναι γραμμικά ανεξάρτητα (ως  $A$ -ευσυχή). Άρα, εισάγοντας του  $n \times k$  πίνακα  $R_k = [r^0, \dots, r^{k-1}]$  έχουμε  $P_k = R_k C$ , όπου  $C$  ένας  $k \times k$  αντιστρέψιμος, άνω τριγωνικός πίνακας. Συνεπώς  $R_k = P_k C^{-1}$  και επειδή ο  $C^{-1}$  είναι άνω τριγωνικός συμπεραίνουμε ότι για  $1 \leq j \leq k$ ,  $r^{j-1} \in \langle p^1, \dots, p^j \rangle$ . Από την (3) έχουμε όμως ότι  $(p^j, r^i) = 0$ ,  $1 \leq j \leq i \leq k-1$ . Από τις δύο αυτές εκθέσεις έχουμε ότι  $(r^j, r^i) = 0$  για  $0 \leq j \leq i \leq k-1$ , ά.έ.δ.

(β) Η (8) αποδεικνύεται με επαγωγή ως προς  $j$  ως εξής: Επειδή  $p^1 = r^0 = b$ , ισχύει για  $j=1$ . Έστω τώρα ότι ισχύει για κάποιο  $j$ ,  $1 \leq j \leq k-1$ . Τότε η (2) ευενάχεται ότι  $r^j = r^{j-1} - a_j Ap^j \in \langle b, Ab, \dots, A^j b \rangle$ .

Εξ άλλου από την (6) έχουμε τότε ότι  $p^{j+1} = r^j - R(z_j p^1 + \dots + z_j p^j)$

$\in \langle b, Ab, \dots, A^j b \rangle$ . Συνεπώς και οι δύο υπόχωροι  $\langle r^0, \dots, r^j \rangle$  και  $\langle p^1, \dots, p^{j+1} \rangle$  (και οι δύο διαστάσεως  $j+1$  - γιατί  $r^j \neq 0$ ) περιέχονται στον υπόχωρο  $\langle b, Ab, \dots, A^j b \rangle$ , ο οποίος έχει συνεπώς διάσταση  $j+1$ . Άρα η (8) ισχύει για  $j+1$ . (Χρησιμοποιήσαμε βεβαίως το γεγονός ότι  $r^j \neq 0$ ,  $0 \leq j \leq k-1$ ). @

Τα μέχρι τώρα θεωρητικά αποτελέσματά μας για την μέθοδο των ευζυγών κλίσεων αρκούν για να εκτιμήσουμε το εφάλμα  $x-x^k$  σε κάθε βήμα. Κατ' αρχήν μία παρατήρηση: μπορούμε εύκολα να γενικεύσουμε την μέθοδο για οποιαδήποτε αρχική τιμή  $x^0 \neq 0$ . Για οποιαδήποτε  $x^0 \in \mathbb{R}^n$  ορίζουμε στον αλγόριθμο (1.5.22)  $r^0 = b - Ax^0$  και τα υπόλοιπα  $p^k, a^k, x^k, r^k$  όπως πριν. Στην γενική περίπτωση  $x^0 \neq 0$  όμως, ορισμένα προφανή πράγματα αλλάζουν στις αποδείξεις. Π.χ. το  $x^j$  λύνει τώρα το πρόβλημα ελαχιστοποίησης

$$(4') \quad \min_{y \in x^0 + \langle p^1, \dots, p^j \rangle} \varphi(x).$$

Οι εκθέσεις (1), (2), (3) εξακολουθούν να ισχύουν όπως επίσης και το συμπέρασμα (α) της Πρότασης 1. Η (8) πρέπει να αντικατασταθεί από την

$$(8') \quad \langle p^1, \dots, p^j \rangle = \langle r^0, \dots, r^{j-1} \rangle = \langle r^0, Ar^0, \dots, A^{j-1}r^0 \rangle, \quad 1 \leq j \leq k.$$

**ΘΕΩΡΗΜΑ 1.** Έστω  $\{x^j\}$ ,  $j \geq 0$  η ακολουθία που παράγει η μέθοδος των ευζυγών κλίσεων με οποιαδήποτε αρχική τιμή  $x^0 \in \mathbb{R}^n$ , έστω  $x$  η λύση του ευστήρατος (1.5.1) και  $e^j = x - x^j$ . Έστω  $\kappa = \lambda_{\max} / \lambda_{\min}$  όπου  $\lambda_{\max}$ , αντιστ.  $\lambda_{\min}$ , είναι η μέγιστη, αντιστ. ελάχιστη, ιδιοτιμή του  $A$ . Για  $j \geq 1$  έχουμε τότε ότι

$$(9) \quad (Ae^j, e^j)^{1/2} \leq 2[(\kappa^{1/2}-1)/(\kappa^{1/2}+1)]^j (Ae^0, e^0)^{1/2}.$$

**Απόδειξη:** Κατ' αρχήν μία προκαταρκτική παρατήρηση: για οποιαδήποτε  $y \in \mathbb{R}^n$  έχουμε  $(A(x-y), x-y) = 2\varphi(y) + (Ax, x)$ . Συνεπώς το πρόβλημα ελαχιστοποίησης  $\min_{y \in S} \varphi(y)$  για  $S \subset \mathbb{R}^n$  είναι ισοδύναμο με το πρόβλημα ελαχιστοποίησης  $\min_{y \in S} (A(x-y), x-y)$ . Συμπεραίνουμε, επειδή το  $x^j$  είναι η (μοναδική) λύση του (4'), ότι

$$\begin{aligned}
 (10) \quad (Re^j, e^j) &= (A(x-x^j), x-x^j) = \min_{y \in x^0 + \langle p^1, \dots, p^j \rangle} (A(x-y), x-y) \\
 &= \min_{z \in \langle p^1, \dots, p^j \rangle} (A(e^0+z), e^0+z).
 \end{aligned}$$

Λόγω της (8') έχουμε ότι  $z \in \langle p^1, \dots, p^j \rangle \Leftrightarrow z \in \langle \pi_0, \pi_1, \dots, \pi^{j-1} \rangle \Leftrightarrow z = \pi_{j-1}(A)r^0$  για κάποιο πολυώνυμο  $\pi_{j-1} \in \mathcal{P}_{j-1}$ , όπου με  $\mathcal{P}_k$  θα συμβολίζουμε τον χώρο των πραγματικών πολυωνύμων βαθμού το πολύ  $k$ . Επειδή  $r^0 = Ae^0$ , συμπεραίνουμε από τον (10) ότι

$$(11) \quad (Re^j, e^j) = \min_{\pi \in \mathcal{P}_{j-1}} ((1+\pi(A))^2 Ae^0, e^0).$$

(Υποθέσαμε έμμεσα ότι  $j \geq 1$  και ότι ο αλγόριθμος της μεθόδου των ευζυγών κλίσεων ολοκληρώνει  $j$  βήματα. Η σχέση (11) ισχύει και για  $j=0$ , αν ορίσουμε  $\mathcal{P}_{-1} = \{0\}$ .) Χρησιμοποιώντας τώρα την φασματική παράσταση του  $A$  παίρνουμε, με ανάλογες πράξεις με εκείνες της απόδειξης του θεωρήματος 1.5.1. ότι

$$(Re^j, e^j) \leq \min_{\pi \in \mathcal{P}_{j-1}} \left( \max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} (1+\pi(\lambda))^2 \right) (Ae^0, e^0).$$

Συνοπώς έχουμε τελικά ότι

$$(12) \quad (Re^j, e^j)^{1/2} \leq \varepsilon_j (Ae^0, e^0)^{1/2}, \quad j \geq 0$$

όπου για  $j \geq 0$  ( $\varepsilon_0 = 1$ )

$$(13) \quad \varepsilon_j = \min_{\pi \in \mathcal{P}_j, \pi(0)=1} \left( \max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |\pi(\lambda)| \right)$$



θα αποδείξουμε στο Λήμμα 4 πιο κάτω ότι

$$(14) \quad \epsilon_j \leq 2[(\kappa^{1/2}-1)/(\kappa^{1/2}+1)]^j, j \geq 1.$$

Από τις (12) και (14) έπεται αμέσως η (9). @

Η λύση του προβλήματος "min-max" (12) είναι κλασική και γίνεται με χρήση πολυωνύμων Chebyshev: Τα πολυώνυμα Chebyshev  $T_j(z)$ , βαθμού  $j \geq 0$ , μίας πραγματικής μεταβλητής  $z$ , ορίζονται ως γνωστόν από τις αναδρομικές σχέσεις

$$(14) \quad \begin{aligned} T_0(z) &= 1 \\ T_1(z) &= z \\ T_j(z) &= 2zT_{j-1}(z) - T_{j-2}(z), \quad j \geq 2. \end{aligned}$$

Από τον ορισμό αυτό μπορούμε εύκολα να αποδείξουμε επαγωγικά τις εξής ιδιότητες των πολυωνύμων Chebyshev  $T_m(z)$ ,  $m \geq 0$ :

(α) βαθμ  $(T_m(z)) = m$ ,  $T_m(1) = 1$ ,  $T_m(z)$  περιττή, αντιστ. άρτια, συνάρτηση του  $z$  αν  $m$  περιττός, αντιστ. άρτιος.

$$(β) \quad T_m(z) = \begin{cases} \cos(m \cos^{-1} z) & \text{αν } -1 \leq z \leq 1 \\ \cosh(m \cosh^{-1} z) & \text{αν } z \geq 1, \end{cases}$$

δηλ.

$$(\beta_1) \quad T_m(z) = \cos(m\theta), \quad \text{όπου } z = \cos\theta, \quad \theta \in [0, \pi], \quad \text{αν } -1 \leq z \leq 1.$$

$$(\beta_2) \quad T_m(z) = \cosh(mu), \quad \text{όπου } z = \cosh u, \quad u \geq 0, \quad \text{αν } z \geq 1.$$

$$(\gamma) \quad T_m(z) = [(z + (z^2 - 1)^{1/2})^m + (z - (z^2 - 1)^{1/2})^m] / 2.$$

Τα πολυώνυμα Chebyshev είναι χρήσιμα στην περίπτωση μας λόγω της εξής σημαντικής ιδιότητάς τους:

Λήμμα 4. Έστω  $0 < \alpha < \beta$ . Τότε το πρόβλημα min-max

$$(16) \min_{p \in \mathcal{P}_m, p(0)=1} \left( \max_{\alpha \leq z \leq \beta} |p(z)| \right)$$

λύεται μοναδικά από το πολυώνυμο

$$(17) \tilde{p}_m(z) = T_m\left(\frac{\beta+\alpha-2z}{\beta-\alpha}\right) / T_m\left(\frac{\beta+\alpha}{\beta-\alpha}\right),$$

για το οποίο

$$(18) \max_{\alpha \leq z \leq \beta} |\tilde{p}_m(z)| = 1 / T_m\left(\frac{\beta+\alpha}{\beta-\alpha}\right),$$

Απόδειξη: Από την παράσταση  $(\beta_1)$  των πολυωνύμων Chebyshev για  $-1 \leq y \leq 1$  έχουμε για  $m \geq 0$  ότι  $\max_{-1 \leq y \leq 1} |T_m(y)| = 1$  και ότι η μέγιστη αυτή απόλυτη τιμή λαμβάνεται στα σημεία  $y_j = \cos(j\pi/m)$ ,  $j=0, 1, \dots, m$ , όπου  $T_m(y_j) = (-1)^j$ , δηλ. όπου το  $T_m(y)$  έχει εναλλασσόμενο πρόσημο. Θεωρούμε τώρα τον γραμμικό μετασχηματισμό  $z \mapsto y$ ,  $y = (\beta + \alpha - 2z) / (\beta - \alpha)$  που απεικονίζει αμφιμονοσήμαντα το διάστημα  $\alpha \leq z \leq \beta$  πάνω στο  $-1 \leq y \leq 1$ . Για  $\alpha \leq z \leq \beta$  θεωρούμε το πολυώνυμο  $\tilde{p}_m(z)$  που δίνεται από την (17). Έχουμε  $\tilde{p}_m(0) = 1$ . Επίσης από τα παραπάνω,  $\max_{\alpha \leq z \leq \beta} |\tilde{p}_m(z)| = \max_{-1 \leq y \leq 1} |T_m(y)| / T_m\left(\frac{\beta+\alpha}{\beta-\alpha}\right) = 1 / T_m\left(\frac{\beta+\alpha}{\beta-\alpha}\right)$ , δηλ. ισχύει η (18). (Επειδή  $(\beta+\alpha)/(\beta-\alpha) > 1$  χρησιμοποιήθηκε η  $(\beta_2)$ ). Επίσης, η μέγιστη αυτή τιμή του  $|\tilde{p}_m(z)|$  για  $\alpha \leq z \leq \beta$  λαμβάνεται σε  $m+1$  διακριτά σημεία  $z_j$ ,  $0 \leq j \leq m$ , του  $[a, b]$ , (στις προεικόνες των  $y_j$ ), όπου το  $\tilde{p}_m$  έχει εναλλασσόμενο πρόσημο. Έστω τώρα ότι υπάρχει άλλο πραγματικό πολυώνυμο  $q_m(z)$ , βαθμού  $\leq m$ , τέτοιο ώστε  $q_m(0) = 1$  και  $\max_{\alpha \leq z \leq \beta} |q_m(z)| < \max_{\alpha \leq z \leq \beta} |\tilde{p}_m(z)|$ . Θεωρούμε το πολυώνυμο  $p_m(z) = \tilde{p}_m(z) - q_m(z)$ , βαθμού  $\leq m$ . Το  $p_m(z)$  έχει εναλλασσόμενο πρόσημο στα σημεία  $z_j$ ,  $0 \leq j \leq m$  και ευνεπώς υπάρχουν  $m$  διακριτά σημεία  $s_j$ ,  $1 \leq j \leq m$ ,

τέτοια ώστε  $0 < \alpha \leq z_{j-1} < s_j < z_j \leq \beta$ , όπου  $p_m(s_j) = 0$ . Επίσης  $p_m(0) = \tilde{p}_m(0) - q_m(0) = 0$ , δηλ. το πολυώνυμο  $p_m(z)$  έχει  $m+1$  διακριτές ρίζες  $\Rightarrow p_m(z) = 0$ , άτοπο. Συνεπώς το πολυώνυμο  $\tilde{p}_m(z)$  λύνει το πρόβλημα min-max (16). Η μοναδικότητα της λύσης του (16) αφήνεται ως άσκηση). @

Με την βοήθεια του λήμματος αυτού μπορούμε να αποδείξουμε τώρα την ισχύ των (14).

**Λήμμα 4.** Αν οι αριθμοί  $\epsilon_j$  ορίζονται από την (13), τότε ισχύει η (14).

Απόδειξη: Από την (13) και το Λήμμα 3 συμπεραίνουμε ότι  $\epsilon_j = 1/T_j((\lambda_{\max} + \lambda_{\min})/(\lambda_{\max} - \lambda_{\min}))$ , αν  $\lambda_{\max} \neq \lambda_{\min}$ . (Η (13) δίνει ότι  $\epsilon_j = 0$  αν  $\lambda_{\min} = \lambda_{\max}$ ,  $j \geq 1$ ). Συνεπώς για  $k \neq 1$ ,  $j \geq 1$  έχουμε

αυτά οπότε από θεωρήματα στο κεφάλαιο 4 προκύπτει ότι  $\epsilon_j = 1/T_j((k+1)/(k-1))$ .  $T_j(x)$  είναι ο πολυώνυμος Chebyshev βαθμού  $j$  με κεντρικό μετασχηματισμό  $z \rightarrow y$ ,  $y = (\alpha + \beta - 2z)/(\beta - \alpha)$ . Χρησιμοποιώντας την παράσταση  $T_j(x)$  των πολυωνύμων Chebyshev για  $z = (k+1)/(k-1)$ , παίρνουμε μετά από λίγες πράξεις ότι

Έχουμε  $\tilde{p}_m(0) = 1$ . Επίσης από τα παραπάνω προκύπτει  $\tilde{p}_m(0) =$

(20)  $\epsilon_j = f(k^{-1/2}) [(k^{1/2}-1)/(k^{1/2}+1)]^j$ ,  $j \geq 1$

όπου για  $x \in [0, 1]$  η συνάρτηση  $f(x)$  ορίζεται ως

για την οποία ισχύει  $1 = f(0) \leq f(x) \leq f(1) = 2$ ,  $x \in [0, 1]$ . Συνεπώς, λόγω της (20) και του ότι  $k \geq 1$ , μία εκτίμηση του  $\epsilon_j$  δίνεται από την (14). @

Ας έρθουμε τώρα στο πρόβλημα της αποτελεσματικής υλοποίησης του αλγορίθμου (1.5.22) στην πράξη. Θα στηριχθούμε στα αποτελέσματα των Λημμάτων 1 και 2 και της Πρότασης 1 και στο εξής εμφαντικό πόρισμα

τους: Για  $k > 2$  το διάνυσμα  $p^k$  είναι γραμμικός συνδυασμός των  $p^{k-1}$  και  $r^{k-1}$ . Για  $k=2$  αυτό είναι προφανές λόγω της (8). Αν  $k > 2$ , υποθέτουμε πάντα ότι  $r^{k-1} \neq 0$ , θεωρούμε το διάνυσμά  $z^{k-1} \in \mathbb{R}^{k-1}$  της (6) το οποίο γράφουμε ως  $z^{k-1} = (w, \mu)^T$ , όπου  $w \in \mathbb{R}^{k-2}$  και  $\mu \in \mathbb{R}^1$  ( $\mu \neq 0$  γιατί;) Χρησιμοποιώντας τώρα την (6) και την (2) για  $i=k-1$  έχουμε

$$(21) \quad p^k = (1 + (\mu/a_{k-1}))r^{k-1} + s^{k-1},$$

όπου

$$(22) \quad s^{k-1} = -\mu r^{k-2}/a_{k-1} - AP_{k-2}w.$$

Η (22) δίνει τώρα ότι  $(s^{k-1}, r^{k-1}) = 0$ , επειδή τα  $r^i$  είναι ορθογώνια και επειδή, λόγω της (6),  $AP_{k-2}w = A(w_1 p^1 + \dots + w_{k-2} p^{k-2}) \in \langle Ar^0, Ar^2, \dots, Ar^{k-2} \rangle \subset \langle r^0, \dots, r^{k-2} \rangle$ . Συνεπώς, το Πυθαγόρειο θεώρημα και η (21) δίνουν

$$(23) \quad \|p^k\|^2 = (1 + (\mu/a_{k-1}))^2 \|r^{k-1}\|^2 + \|s^{k-1}\|^2.$$

Θυμόμαστε τώρα τον χαρακτηρισμό του  $z^{k-1} = (w, \mu)^T$  ως εκείνου του στοιχείου του  $\mathbb{R}^{k-1}$  που λύνει το πρόβλημα ελαχίστων τετραγώνων (?): ευμπεραίνουμε λοιπόν ότι η κατασκευή του  $z^{k-1}$  είναι ισοδύναμη με τον καθορισμό εκείνων των  $w, \mu$  που ελαχιστοποιούν το  $\|p^k\| (= \|r^{k-1} - AP_{k-1}z^{k-1}\|, \text{βλ. (6)})$ , δηλ. που ελαχιστοποιούν το δεύτερο μέλος της (23). Επειδή, από την (22) η ποσότης

$$\|s^{k-1}\|^2 = (\mu/a_{k-1})^2 \|r^{k-2} - AP_{k-2}w'\|^2, \quad w' = a_{k-1}w/\mu$$

ελαχιστοποιείται ως προς  $w'$  για  $r^{k-2} - AP_{k-2}w' = p^{k-1}$ . (βλ. Λήμμα 2), θα πρέπει το  $s^{k-1}$  να είναι πολλαπλάσιο του  $p^{k-1}$ . Η (21) τώρα αποδεικνύει τον ισχυρισμό.

Συνεπώς  $p^k \in \langle p^{k-1}, r^{k-1} \rangle$  για  $k \geq 2$ , δηλ. η ευθεία  $x^{k-1} + ap^k$ ,  $a \in \mathbb{R}$ , που διέρχεται από το  $x^{k-1}$  και πάνω στην οποία ελαχιστοποιείται το  $\varphi(x)$  για  $x = x^k$ , βρίσκεται στο επίπεδο  $x^{k-1} + \langle p^{k-1}, r^{k-1} \rangle$  η τομή του οποίου με το ελλειψοειδές  $\varphi(x) = \varphi(x^{k-1})$  (δηλ. με την "ισοψή" επιφάνεια του  $\varphi$  που διέρχεται από το σημείο  $x^{k-1}$ ) είναι μία έλλειψη  $C_k$  που διέρχεται φυσικά από το  $x^{k-1}$ , όπου εφάπτεται στην ευθεία  $x^{k-1} + ap^{k-1}$  (αφού το  $x^{k-1}$  ήταν ακρότατο του  $\varphi(x)$  πάνω στην ευθεία αυτή) και όπου έχει κάθετο παράλληλη προς το  $r^{k-1}$ . Οι ελλείψεις - τομές των ισοψών  $\varphi(x) = \text{σταθ.}$  του  $\varphi$  με το επίπεδο  $x^{k-1} + \langle p^{k-1}, r^{k-1} \rangle$  είναι ομόκεντρες και ομοιόθετες της  $C_k$ . Άρα το ελάχιστο του  $\varphi(x)$  πάνω ε' αυτό το επίπεδο λαμβάνεται στο σημείο  $x^k$  που είναι το κοινό κέντρο των ελλείψεων αυτών. Συνεπώς οι κατευθύνσεις  $p^k$  και  $p^{k-1}$  είναι "ευζυχείς" ως προς την έλλειψη  $C_k$ , εξ ου και η ονομασία τους "H-ευζυχείς" (επειδή ικανοποιούν την (1.5.21)) είναι συμβιβαστή με την γνωστή από την επίπεδη γεωμετρία έννοια. Για περαιτέρω πάνω στην γεωμετρική ερμηνεία της μεθόδου των ευζυγών κλίσεων (καθώς και για την πληρέστερη ανάλυση των κλασικών ιδιοτήτων των μεθόδων ελαχιστοποίησης), βλ. το βιβλίο του M.R. Hestenes, "Conjugate direction methods in optimization", Springer-Verlag, Berlin 1980.

Επαναρχόμαστε στο θέμα της υλοποίησης του αλγορίθμου (1.5.22). Χωρίς περιορισμό της γενικότητας (βλ. θεκ. 5(γ)) υποθέτουμε ότι

$$(24) \quad p^k = r^{k-1} + \beta_k p^{k-1}, \quad k \geq 2.$$

Επειδή τα  $p^k, p^{k-1}$  είναι H-ευζυγή η (24) δίνει

$$(25) \quad \beta_k = -(p^{k-1}, Ar^{k-1}) / (Ar^{k-1}, p^{k-1}).$$

Οι τύποι (24) και (25) ορίζουν πλήρως το  $p^k$ . Εξ άλλου επειδή  $(p^{k-1}, r^{k-1}) = 0$ , η (24) και η (1.5.19) δίνουν ότι

$$(26) \quad \alpha_k = \|r^{k-1}\|^2 / (Ar^k, p^k).$$

Οδηγούμαστε λοιπόν στην εξής υλοποίηση του αλγορίθμου (1.5.22):

$$\begin{array}{l}
 x^0 = 0 \\
 \left. \begin{array}{l}
 \text{Για } k=1, 2, \dots, n \\
 r^{k-1} = b - Ar^{k-1} \\
 \text{Αν } r^{k-1} = 0 \\
 \text{τότε, εκχώρησε } x = x^{k-1} \text{ και τερμάτισε.} \\
 \text{αλλιώς,} \\
 \text{αν } k=1 \\
 \text{τότε, } p^k = r^0 \\
 \text{αλλιώς,} \\
 \beta_k = -(r^{k-1}, Ar^{k-1}) / (Ar^{k-1}, r^{k-1}) \\
 p^k = r^{k-1} + \beta_k r^{k-1} \\
 a_k = \|r^{k-1}\|^2 / (Ar^k, p^k) \\
 x^k = x^{k-1} + a_k p^k
 \end{array} \right\} x = x^n
 \end{array}
 \tag{27}$$

Ο αλγόριθμος απαιτεί δύο πολλαπλασιασμούς του πίνακα  $A$  επί διάνυσμα σε κάθε βήμα  $k$ . Παρατηρώντας όμως ότι τα υπόλοιπα μπορούν να υπολογισθούν αναδρομικά μέσω της  $r^j = r^{j-1} - a_j Ar^j$ , έχουμε χρησιμοποιώντας την σχέση αυτή για  $j=k-1$  και την ορθογωνιότητα των  $r^{k-1}, r^{k-2}$  - ότι  $\|r^{k-1}\|^2 = -a_{k-1}(r^{k-1}, Ar^{k-1})$  και  $\|r^{k-2}\|^2 = a_{k-1}(r^{k-2}, Ar^{k-1}) = a_{k-1}(r^{k-1}, Ar^{k-1})$  - λόγω της  $p^{k-1} = r^{k-2} + \beta_{k-1} r^{k-2}$ . Συνεπώς  $\beta_k = -(r^{k-1}, Ar^{k-1}) / (Ar^{k-1}, p^{k-1}) = \|r^{k-1}\|^2 / \|r^{k-2}\|^2$ , αποφεύγοντας έτσι τον πολλαπλασιασμό  $Ar^{k-1}$ . Καταλήγουμε λοιπόν στην εξής τελική μορφή του αλγορίθμου της μεθόδου των ευζυγών κλίσεων όπως ουσιαστικά τον παρουσίασαν οι Hestenes και Stiefel (1952):

$$\begin{aligned}
 & x^0 = 0 \\
 & r^0 = b \\
 & \text{Γι} \acute{\alpha} \ k=1, 2, \dots, n \\
 & \text{Αν } r^{k-1} = 0 \\
 & \text{τότε, εκχώρησε } x = x^{k-1} \text{ τερμάτισε.} \\
 & \text{αλλιώς} \\
 & \text{αν } k=1 \\
 & \text{τότε, } p^k = r^0 \\
 & \text{αλλιώς,} \\
 & \beta_k = \|r^{k-1}\|^2 / \|r^{k-2}\|^2 \\
 & p_k = r^{k-1} + \beta_k p^{k-1} \\
 & a_k = \|r^{k-1}\|^2 / (A p^k, p^k) \\
 & x^k = x^{k-1} + a_k p^k \\
 & r^k = r^{k-1} - a_k A p^k \\
 & x = x^n
 \end{aligned}
 \tag{28}$$

### Παρατηρήσεις

1. Η απόδειξη της βασικής ιδιότητας της μεθόδου των ευζυγών κλίσεων, ότι δηλ. βρίσκει την λύση  $x^n = x$  σε  $n$  βήματα (ισοδύναμη με την ορθογωνιότητα των υπολοίπων  $r^0, r^1, \dots, r^{n-1}$ , οπότε  $r^n = 0$ ) έγινε με την προϋπόθεση βέβαια ότι όλοι οι υπολογισμοί, π.χ. στον αλγόριθμο (27), γίνονται ακριβώς. Στην πράξη όμως (όπως έγινε χρήχαρα εαφές μετά την ανακάλυψη της μεθόδου) τα εσφάλματα εστροχχύλεισης που οφείλονται στην αριθμητική πεπερασμένης ακρίβειας καταστρέφουν την ορθογωνιότητα των  $r^j$  που υπολογίζουμε, έτσι ώστε τελικά το  $r^n$  να μην είναι μηδέν. Πάλιετα, όπως περιμένουμε, όσο αυξάνει ο δείκτης του  $A$ , τόσο το φαινόμενο αυτό, δηλ. ότι  $(r^j, r^j) \neq 0$ , γίνεται πιά έντονο. Εν τούτοις, η μέθοδος των ευζυγών κλίσεων, σαν μέθοδος ελαχιστοποίησης που είναι, όπως ελαττώνει την τιμή του συναρτησιακού  $\varphi(x)$  από βήμα σε βήμα. Συνεπώς μπορεί να χρησιμοποιηθεί (και έτσι χρησιμοποιείται σήμερα στην πράξη) σαν επαναληπτική μέθοδος που παράγει προσεγγίσεις

$x^0, x^1, \dots, x^k, \dots$  της λύσης  $x$  του συστήματος  $Ax=b$  η επανάληψη δεν σταματά στο βήμα  $n$  αλλά τερματίζεται όταν ικανοποιηθούν κάποια από τα γνωστά μας κριτήρια τερματισμού. Για περισσότερες λεπτομέρειες βλ. [5.4, Παρ.16]. Θα θέλαμε όμως απλώς να παρατηρήσουμε εδώ ότι επειδή η κύρια και επαναλαμβανόμενη σε κάθε βήμα πράξη του αλγορίθμου είναι ο πολλαπλασιασμός του πίνακα  $A$  επί διάνυσμα, η μέθοδος των ευζυγών κλίσεων είναι ιδιαίτερα αποτελεσματική για αραιούς πίνακες  $A$ . Για τέτοιους πίνακες υπάρχουν δομές δεδομένων για την αποθήκευση των μη μηδενικών στοιχείων τους ιδιαίτερα κατάλληλες για την πράξη του πολλαπλασιασμού πίνακα επί διάνυσμα, βλ. [5.4, παρ. 16]

2. Στην πράξη λοιπόν χρησιμοποιούμε την μέθοδο ευζυγών κλίσεων σαν επαναληπτική. Μας ενδιαφέρει ευνεπώς η ταχύτητα εύγκλισης (δηλ. ελάττωσης του σφάλματος), που μας δίνει π.χ. η (9) ιδίως για μεγάλο  $k$ . Αν και η ταχύτητα εύγκλισης της μεθόδου των ευζυγών κλίσεων είναι μεγαλύτερη της ταχύτητας της μεθόδου της καθόδου μεγίστης κλίσεως (γιατί  $(k-1)/(k+1) = 1-2k^{-1}+O(k^{-2})$  ενώ  $(k^{1/2}-1)/(k^{1/2}+1) = 1-2k^{-1/2}+O(k^{-1})$  όταν  $k \rightarrow \infty$ , βλ. (1.5.8), (9), εξακολουθεί να είναι πολύ μικρή για μεγάλο  $k$ . Μια σημαντική τεχνική που χρησιμοποιείται για την επιτάχυνση της εύγκλισης είναι η λεγόμενη προϋψμιση (preconditioning), για την υλοποίηση της οποίας ο αλγόριθμος απαιτεί σε κάθε βήμα του εκτός ενός πολλαπλασιασμού του πίνακα  $A$  επί διάνυσμα και την λύση (με την μέθοδο Cholesky π.χ.) ενός  $n \times n$  γραμμικού συστήματος με ένα συμμετρικό θετικά ορισμένο πίνακα, του λεγόμενου προϋψμιστή (preconditioner)  $M$ . Η επιλογή κατάλληλου προϋψμιστή είναι σημαντικό πρόβλημα: πρέπει να έχει απλή δομή έτσι ώστε η λύση συστημάτων με πίνακα  $M$  να μην απαιτεί πολλές πράξεις· επιπλέον πρέπει ο πίνακας  $\tilde{A}=M^{-1}A$  να έχει λόγο ιδιοτιμών  $\tilde{\lambda}_{\max}/\tilde{\lambda}_{\min}$  μικρότερο του  $\kappa=\lambda_{\max}/\lambda_{\min}$  ώστε να επιταχύνεται πράγματι η εύγκλιση. Βλέπε [5.4, Παρ. 16] για περισσότερα επ' αυτού.



### Ασκήσεις 1.6

1. (α) Χρησιμοποιώντας την (1) και τις (1.5.11), (1.5.21), δείξτε ότι  $(r^j, r^j) = 0$ ,  $j=1, 2, \dots, i$ , δηλ. ότι ισχύει η (3), υπό τις προϋποθέσεις του Λήμματος 1.

(β) Αναπτύξτε ετά "γιατί;" των αποδείξεων του Λήμματος 1 και της Πρότασης 1.

2. Έστω  $r \in \mathbb{R}^n$  και  $\Pi$  υπόχωρος του  $\mathbb{R}^n$ . Λέμε ότι το  $r \in \Pi$  είναι η (ορθή) προβολή του  $r$  στον  $\Pi$  (ή προβολή ως προς το Ευκλείδιο εσωτερικό γινόμενο  $(\cdot, \cdot)$ ) αν

$$(r, x) = (p, x), \quad x \in \Pi.$$

(α) Δείξτε ότι η προβολή  $p$  του  $r$  στον  $\Pi$  υπάρχει, είναι μοναδική και ικανοποιεί  $\|p\|^2 = \|r\|^2 - \|r-p\|^2$ .

(β) Δείξτε ότι το  $r-p$  είναι η προβολή του  $r$  στον υπόχωρο  $\Pi^\perp$ .

(γ) Δείξτε ότι  $\|r-p\| = \min_{x \in \Pi} \|r-x\|$  και ότι η ιδιότητα αυτή χαρακτηρίζει μονασήμαντα την προβολή  $p$  του  $r$ .

(δ) Δείξτε ότι για κάθε  $B \in \mathbb{R}^{m \times n}$ ,  $\mathcal{R}(B)^\perp = \mathcal{N}(B^T)$ , όπου  $\mathcal{N}(B^T) = \{x \in \mathbb{R}^n : B^T x = 0\}$ .

(ε) Τι σημαίνει το αποτέλεσμα (δ) στην περίπτωση του Λήμματος 2, όπου  $p^k \in \mathcal{R}(A P_{k-1})^\perp$ ;

3. Αναπτύξτε την μέθοδο του ευζυγώου κλίσεων για  $x^0 \neq 0$  (τα  $x^j$  λύνει τώρα το πρόβλημα (4')). Βρείτε τα ανάλογα των Λημμάτων 1 και 2 και της Πρότασης 1.

4. (α). Αποδείξτε τις ιδιότητες (α), (β), (γ) των πολυωνύμων Chebyshev.

(β). Δείξτε ότι το πολυώνυμο  $\tilde{p}_m(z)$  που ορίζεται από την (17) είναι η μοναδική λύση του προβλήματος min-max (16).

5. (α). Γιατί στην αρχή της απόδειξης του ισχυρισμού ότι το  $p^k \in \langle p^{k-1}, r^{k-1} \rangle$  (σελ. 1.6.9) ισχύει ότι  $\mu \neq 0$ ;

(β) Στην ίδια απόδειξη - στο τελευταίο βήμα της, σελ. 1.6.9 - γιατί αγνοήσαμε το  $\mu$  στο πρόβλημα της ελαχιστοποίησης του δευτέρου μέλους της (23) και μελετήσαμε μόνο την ελαχιστοποίηση ως προς  $w$ ;

(γ) Δείξτε ότι η υπόθεση (24) δεν αποτελεί περιορισμό της γενικότητας, δηλ. δείξτε ότι ένας γενικός γραμμικός συνδυασμός της μορφής  $p^k = \gamma_k r^{k-1} + \beta_k p^{k-1}$  οδηγεί, μέσω π.χ. του αλγορίθμου (27) στον υπολογισμό των ίδιων  $x^k$ , όπως και προηγουμένως.

## 2. ΑΡΙΘΜΗΤΙΚΗ ΛΥΣΗ ΜΗ ΓΡΑΜΜΙΚΩΝ ΣΥΣΤΗΜΑΤΩΝ

### 2.1 ΠΑΡΑΓΩΓΙΣΙΜΕΣ ΣΥΝΑΡΤΗΣΕΙΣ ΣΤΟΝ $\mathbb{R}^n$

Στο κεφάλαιο αυτό θα ασχληθούμε με μεθόδους για την αριθμητική επίλυση μη γραμμικών ευστημάτων η εξίσωση με η αγνώστους, δηλ. ευστημάτων της μορφής

$$(1) f_i(x_1, \dots, x_n) = 0, \quad i=1, 2, \dots, n,$$

τα οποία γράφουμε ευθέως στην διανυσματική μορφή

$$(1') F(x) = 0,$$

όπου  $F$  είναι μία (μη γραμμική εν γένει) απεικόνιση ενός υποσυνόλου  $D$  του  $\mathbb{R}^n$  στον  $\mathbb{R}^n$  (θα γράφουμε  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ ) με συνιστώσες  $F(x) = (f_1(x), \dots, f_n(x))^T$  και όπου  $x = (x_1, \dots, x_n)^T$  κατά τα χωστά.

Προβλήματα που οδηγούν στην επίλυση μη γραμμικών ευστημάτων της μορφής (1) εμφανίζονται πολύ συχνά στις εφαρμογές. Μία σημαντική πηγή προβλημάτων είναι π.χ. ο υπολογισμός τοπικών ακροτάτων ενός συναρτησιακού  $g: \mathbb{R}^n \rightarrow \mathbb{R}^1$ , οπότε  $F = \nabla g$ , υπό την προϋπόθεση βέβαια ότι η κλίση  $\nabla g$  υπάρχει και μπορεί να υπολογισθεί για κάθε  $x$ . Θα εξετάσουμε όμως το πρόβλημα στην γενική μορφή (1) χωρίς να υπεισέλθουμε εδώ σε ειδικές μεθόδους για προβλήματα βελτιστοποίησης.

Σημαντικό ρόλο τόσο στην θεωρία όσο και στις αριθμητικές μεθόδους για την λύση του (1) παίζουν οι ιδιότητες παραγωγισιμότητας της  $F$  καθώς και διαφόρου τύπου θεωρήματα "μέσης τιμής". Τα θεωρήματα αυτά θα εξετάσουμε ε' αυτήν την εισαγωγική παράγραφο.

Λέμε ότι η απεικόνιση  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη ε' ένα σημείο  $x \in \text{Int} D$  (ακριβέστερα παραγωγίσιμη με την έννοια του Frechet ή  $F$  - παραγωγίσιμη) αν υπάρχει γραμμικός τελεστής  $A_x: \mathbb{R}^n \rightarrow \mathbb{R}^n$  τέτοιος ώστε για κάποια νόρμα  $\|\cdot\|$  του  $\mathbb{R}^n$  να ισχύει

$$(2) \lim_{h \rightarrow 0} \frac{\|F(x+h) - F(x) - A_x h\|}{\|h\|} = 0$$

λόγω της ισοδυναμίας των νορμών στον  $\mathbb{R}^n$  η ύπαρξη του  $A_x$  είναι ανεξάρτητη της νόρμας  $\|\cdot\|$ . Είναι επίσης προφανές ότι ο ορισμός (2) γενικεύει την έννοια της παραγωγισιμότητας παραχματικών συναρτήσεων μίας μεταβλητής.

Αν η  $F$  είναι παραγωγίσιμη στο  $x$  τότε ο γραμμικός τελεστής  $A_x$  είναι μοναδικός. Πράγματι αν υπήρχαν δύο γραμμικοί τελεστές  $A_1, A_2$  τέτοιοι ώστε για τον κάθε ένα να ισχύει η (2) θα είχαμε από την τριγωνική ανισότητα ότι για κάθε  $0 \neq y \in \mathbb{R}^n$ ,  $t \neq 0$

$$\|(A_1 - A_2)y\| / \|y\| \leq (\|F(x+ty) - F(x) - A_1(ty)\| / \|ty\|) + (\|F(x+ty) - F(x) - A_2(ty)\| / \|ty\|),$$

από την οποία, θέτοντας  $h=ty$  και παίρνοντας  $t \rightarrow 0$  έχουμε λόγω της (2) ότι  $A_1 y = A_2 y \quad \forall y \in \mathbb{R}^n$ , δηλ. ότι  $A_1 = A_2$ . Αν λοιπόν η  $F$  είναι παραγωγίσιμη στο  $x$  λέμε ότι ο τελεστής  $A_x$  είναι η παράγωγος της  $F$  (ακριβέστερα ή "παράγωγος της  $F$  με την έννοια του Fréchet" ή η " $F$ -παράγωγος της  $F$ ") στο σημείο  $x$  και συμβολίζουμε  $A_x = F'(x)$ . Γενικά λέμε ότι η

$F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη ε' ένα εύνολο  $D_0 \subset D$  αν  $D_0 \subset \text{Int}(D)$  και αν η  $F$  είναι παραγωγίσιμη  $\forall x \in D_0$ . Τότε η  $F'$  μπορεί να θεωρηθεί ως απεικόνιση του  $D_0$  στο εύνολο  $L(\mathbb{R}^n)$  των γραμμικών τελεστών από του  $\mathbb{R}^n$  στον εαυτό του. Εύκολα επίσης μπορούμε να αποδείξουμε την γραμμικότητα της πράξης της παραγωγίσιμης: Έστω ότι οι απεικονίσεις  $F_1, F_2: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμες στο  $x \in \text{Int}(D)$ . Τότε για  $\alpha, \beta \in \mathbb{R}$ , η  $\alpha F_1 + \beta F_2$  είναι παραγωγίσιμη στο  $x$  και  $(\alpha F_1 + \beta F_2)'(x) = \alpha F_1'(x) + \beta F_2'(x)$ .

Αν η παράγωγος  $F'(x)$  υπάρχει στο σημείο  $x \in \text{Int}(D)$  τότε υπάρχουν όλες οι μερικές παράγωγοι  $\partial_j f_i (= \partial f_i / \partial x_j)$ ,  $1 \leq i, j \leq n$  των συνιστωσών  $f_i$  της  $F$  στο σημείο  $x$ , ο δέ  $n \times n$  πίνακας που παριετάνει τον γραμμικό τελεστή  $F'(x)$  ως προς την κανονική βάση  $\{e^j\}$ ,  $1 \leq j \leq n$ , του  $\mathbb{R}^n$

$(e_i = \delta_{ij})$  είναι ο ιακωβιανός πίνακας  $J(x)$ :  $J_{ij}(x) = \partial_j f_i(x)$ ,  $1 \leq i, j \leq n$ .

Πράγματι, θέτοντας, για  $j=1, \dots, n$ ,  $h=te^j$  στην (2) (με  $\|\cdot\|=\|\cdot\|_2$  π.χ.) και υποθέτοντας ότι στο σημείο  $x$  η  $F'(x)=A_x$  παριετάνεται ως προς την βάση  $\{e^j\}$  από τον πίνακα  $(a_{ij})$  έχουμε για  $1 \leq i, j \leq n$

$$\lim_{t \rightarrow 0} |(f_i(x+te^j) - f_i(x))/t - a_{ij}| = 0,$$

δηλ. ότι ούτως  $a_{ij} = \partial_j f_i(x) = J_{ij}(x)$ .

Η ύπαρξη μόνο των μερικών παραγώγων  $\partial_j f_i(x)$  δεν εγγυάται όμως ότι η  $F$  είναι παραγωγίσιμη στο  $x$ . Αυτό φαίνεται αμέσως από την εξής σημαντική συνέπεια της παραγωγισιμότητας:

**ΠΡΟΤΑΣΗ 1.** Έστω ότι η  $F : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη στο σημείο  $x \in \text{Int}(D)$ . Τότε η  $F$  είναι συνεχής στο  $x$ .

Απόδειξη: Επειδή  $x \in \text{Int}(D)$ ,  $\exists \delta_1 > 0$  τέτοιο ώστε  $x+h \in D$  αν  $\|h\| < \delta_1$ . Η (2) ευενάχεται τώρα ότι για δεδομένο  $\epsilon > 0$  υπάρχει  $\delta > 0$ , το οποίο μπορούμε να πάρουμε  $\leq \delta_1$ , τέτοιο ώστε  $\|F(x+h) - F(x) - F'(x)h\| \leq \epsilon \|h\|$  αν  $\|h\| \leq \delta$ , από την οποία έπεται ότι  $\|F(x+h) - F(x)\| \leq (\|F'(x)\| + \epsilon) \|h\|$ . Σταθεροποιώντας το  $\epsilon$  λοιπόν συμπεραίνουμε ότι για δεδομένο  $x \in \text{Int}(D)$   $\exists \delta > 0$  και  $c \geq 0$  τέτοια ώστε  $x+h \in D$  και  $\|F(x+h) - F(x)\| \leq c \|h\|$  αν  $\|h\| \leq \delta$ , που είναι μάλιστα ένα συμπέρασμα ισχυρότερο από την συνέχεια της  $F$  στο  $x$ . ©

Θα διερευνήσουμε σχέσεις μεταξύ της ύπαρξης και της συνέχειας της  $F'(x)$  και των αναλόγων ιδιοτήτων του Ιακωβιανού πίνακα  $J(x)$  καθώς και μιάς "αδευέστερης" παραγώγου της  $F$ , της λεγόμενης παραγώγου Gateaux, σε μιά σειρά παρατηρήσεων και ασκήσεων στο τέλος της παραγράφου. Προς το παρόν θα συνεχίσουμε με την μελέτη ορισμένων θεωρημάτων "μέσης τιμής".

Το γνωστό μας θεώρημα μέσης τιμής για παραγωγίσιμες συναρτήσεις  $f: \mathbb{R}^1 \rightarrow \mathbb{R}^1$ , ότι δηλ.  $\forall x, y \in \mathbb{R}^1$ ,  $f(x) - f(y) = f'(z)(x-y)$  για κάποιο  $z$  μεταξύ των  $x$  και  $y$ , δεν ισχύει για απεικονίσεις  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  αν  $n \geq 2$  (βλ. θεορ. 7(β)). Υπάρχουν όμως εναλλακτικά αποτελέσματα του τύπου "μέσης

τιμής" πολύ χρήσιμα στην μη γραμμική ανάλυση. Παραδείχματος χάριν πολλές φορές ενδιαφερόμαστε απλώς να φράξουμε την ποσότητα  $\|F(x)-F(y)\|$  συναρτήσει της  $F'$ :

**ΠΡΟΤΑΣΗ 2.** Υποθέτουμε ότι η  $F:D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη σε ένα κυρτό εύνολο  $D_0 \subset D$ . Τότε αν  $x, y \in D_0$

$$(3) \|F(x)-F(y)\| \leq \sup_{0 \leq t \leq 1} \|F'(x+t(y-x))\| \|x-y\|.$$

Απόδειξη: Εξ υποθέσεως  $x+t(y-x) \in D_0$  για  $t \in [0,1]$ . Έστω ότι  $M = \sup_{0 \leq t \leq 1} \|F'(x+t(y-x))\| < \infty$ . Για δεδομένο  $\varepsilon > 0$ , έστω  $\Gamma_\varepsilon$  το εύνολο των  $t \in [0,1]$  τέτοιων ώστε

$$(4) \|F(x+t(y-x)) - F(x)\| \leq Mt\|y-x\| + \varepsilon t\|x-y\|.$$

Προφανώς το  $\Gamma_\varepsilon$  δεν είναι κενό γιατί  $0 \in \Gamma_\varepsilon$ . Έστω  $\delta_\varepsilon = \sup_{t \in \Gamma_\varepsilon} t$ . Τότε

$0 \leq \delta_\varepsilon \leq 1$  και επειδή λόγω της Πρότασης 1 η συνάρτηση  $t \mapsto F(x+t(y-x))$  είναι συνεχής στο  $[0,1]$ , παίρνοντας το όριο στην (4) μιάς ακολουθίας  $t_i \rightarrow \delta_\varepsilon$  σημείων του  $\Gamma_\varepsilon$ , έχουμε

$$(5) \|F(x+\delta_\varepsilon(y-x))-F(x)\| \leq M\delta_\varepsilon \|y-x\| + \varepsilon\delta_\varepsilon \|x-y\|.$$

Αν για κάθε  $\varepsilon > 0$   $\delta_\varepsilon = 1$ , τότε η (5) δίνει την ζητούμενη ανισότητα (4).

Αν για κάποιο  $\varepsilon > 0$ ,  $0 \leq \delta_\varepsilon < 1$  επειδή η  $F'$  υπάρχει στο σημείο  $x+\delta_\varepsilon(y-x)$  έχουμε, από τον ορισμό της παραγώγου (2), παίρνοντας ως  $h$  ένα κατάλληλο (μικρό) πολλαπλάσιο του  $y-x$ , ότι υπάρχει  $\beta_\varepsilon \in (\delta_\varepsilon, 1)$  τέτοιο ώστε

$$\|F(x+\beta_\varepsilon(y-x))-F(x+\delta_\varepsilon(y-x))-F'(x+\delta_\varepsilon(y-x))(\beta_\varepsilon-\delta_\varepsilon)(y-x)\| \leq \varepsilon(\beta_\varepsilon-\delta_\varepsilon)\|y-x\|,$$

από την οποία έπεται ότι

$$(6) \|F(x+\beta_\epsilon(y-x))-F(x+\delta_\epsilon(y-x))\| \leq \|(\beta_\epsilon-\delta_\epsilon)\| \|y-x\| + \epsilon(\beta_\epsilon-\delta_\epsilon)\|y-x\|.$$

Οι (5) και (6) δίνουν τότε μέσω της τριγωνικής ανισότητας ότι

$$\|F(x+\beta_\epsilon(y-x))-F(x)\| \leq \| \beta_\epsilon \| \|y-x\| + \epsilon \beta_\epsilon \|y-x\|,$$

δηλ. ότι για αυτό το  $\epsilon$  η (4) ισχύει για κάποιο  $\beta_\epsilon: \delta_\epsilon < \beta_\epsilon < 1$ , πράγμα που αντιφάσκει στον ορισμό του  $\delta_\epsilon$ . Συνεπώς  $\delta_\epsilon = 1 \quad \forall \epsilon > 0$  και το αποτέλεσμα προκύπτει όπως παραπάνω. @

Ενός άλλου τύπου αποτελέσματα είναι "ολοκληρωτικές" μορφές του θεωρήματος μέσης τιμής. Κατά τα γνωστά, αν η διανυσματική συνάρτηση μιάς μεταβλητής  $G: [a, b] \rightarrow \mathbb{R}^n$  έχει συνιστώσες  $G = (g_1, \dots, g_n)^T$ , λέμε ότι η  $G$  είναι ολοκληρώσιμη (με την έννοια του Riemann) αν και μόνο αν για κάθε  $i$  οι συναρτήσεις  $g_i: [a, b] \rightarrow \mathbb{R}^1$  είναι ολοκληρώσιμες κατά Riemann στο  $[a, b]$  οπότε και ορίζουμε

$$\int_a^b G(t) dt = \left( \int_a^b g_1(t) dt, \dots, \int_a^b g_n(t) dt \right)^T.$$

**ΠΡΟΤΑΣΗ 3.** Υποθέτουμε ότι η  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συνεχώς παραγωγίσιμη ε' ένα κυρτό εύνολο  $D_0 \subset D$ . Τότε, αν  $x, y \in D_0$

$$(7) F(y) - F(x) = \int_0^1 F'(x+t(y-x))(y-x) dt.$$

Απόδειξη Επειδή η  $F'$  είναι συνεχής στο  $D_0$ , τότε, για  $x, y \in D_0$  η συνάρτηση  $t \mapsto F'(x+t(y-x))$  είναι συνεχής στο  $[0, 1]$ . Συνεπώς οι διανυσματικές συναρτήσεις  $\nabla f_i(x+t(y-x))$ ,  $1 \leq i \leq n$  (γραμμές του Ιακωβιανού πίνακα που περιτάνει την  $F'(x+t(y-x))$ ) είναι συνεχείς συναρτήσεις του  $t$  για  $t \in [0, 1]$  και συνεπώς ολοκληρώσιμες κατά Riemann στο  $[0, 1]$ . Επειδή τώρα για  $x, y \in D_0$

$df_i(x+t(y-x))/dt = (\nabla f_i(x+t(y-x)))^T (y-x)$ ,  
 έχουμε, ολοκληρώνοντας ως προς  $t$  από 0 έως 1, ότι για  $1 \leq i \leq n$ :

$$f_i(y) - f_i(x) = \int_0^1 \nabla f_i(x+t(y-x))^T (y-x) dt,$$

που είναι ακριβώς η (7) γραμμένη κατά συνιστώσες. @

Θα χρησιμοποιήσουμε στην συνέχεια μία ενδιαφέρουσα εφαρμογή της Πρότασης 3. Πρώτα ένα εύκολο λήμμα:

**ΛΗΜΜΑ 1.** Έστω ότι η  $G: [a, b] \rightarrow \mathbb{R}^n$  είναι συνεχής στο  $[a, b]$ .  
 Τότε

$$(8) \quad \left\| \int_a^b G(t) dt \right\| \leq \int_a^b \|G(t)\| dt.$$

Απόδειξη: Η συνάρτηση  $t \mapsto \|G(t)\|$  είναι συνεχής στο  $[a, b]$  λόγω της υπόθεσής μας και της συνέχειας της  $x \mapsto \|x\|$ . Συνεπώς η  $t \mapsto \|G(t)\|$  είναι ολοκληρώσιμη κατά Riemann στο  $[a, b]$ . Λόγω της ολοκληρωσιμότητας και της  $G(t)$  έχουμε ότι  $\forall \epsilon > 0$  υπάρχει διαμερισμός  $a \leq t_0 < t_1 < \dots < t_s \leq b$  τέτοιος ώστε

$$\left\| \int_a^b G(t) dt - \sum_{j=1}^s G(t_j) (t_j - t_{j-1}) \right\| \leq \epsilon$$

και

$$\left| \int_a^b \|G(t)\| dt - \sum_{j=1}^s \|G(t_j)\| (t_j - t_{j-1}) \right| \leq \epsilon.$$

Συνεπώς από την τριγωνική ανισότητα και τις δύο αυτές εκθέσεις έπεται ότι



$$\begin{aligned} \left\| \int_a^b G(t) dt \right\| &\leq \left\| \sum_{j=1}^s G(t_j)(t_j - t_{j-1}) \right\| + \epsilon \leq \sum_{j=1}^s \|G(t_j)\| (t_j - t_{j-1}) \\ &+ \epsilon \leq \int_a^b \|G(t)\| dt + 2\epsilon. \end{aligned}$$

Επειδή το  $\epsilon > 0$  ήταν αυθαίρετο, έπεται η (8). @

**ΠΡΟΤΑΣΗ 4.** Υποθέτουμε ότι η  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συνεχώς παραγωγίσιμη σ' ένα κυρτό σύνολο  $D_0 \subset D$  και ότι επιπλέον υπάρχουν σταθερές  $a, p \geq 0$  τέτοιες ώστε

$$(9) \quad \|F'(u) - F'(v)\| \leq a \|u - v\|^p, \quad u, v \in D_0.$$

Τότε για κάθε  $x, y \in D_0$

$$(10) \quad \|F(y) - F(x) - F'(x)(y-x)\| \leq a \|x-y\|^{p+1}/(p+1).$$

Απόδειξη: Η  $F$  ικανοποιεί τις υποθέσεις της Πρότασης 3. Συνεπώς για  $x, y \in D_0$

$$F(y) - F(x) = \int_0^1 F'(x+t(y-x))(y-x) dt.$$

Άρα λόγω της (8) και της (9)

$$\begin{aligned} \|F(y) - F(x) - F'(x)(y-x)\| &= \left\| \int_0^1 [F'(x+t(y-x)) - F'(x)](y-x) dt \right\| \\ &\leq \|x-y\| \int_0^1 \|F'(x+t(y-x)) - F'(x)\| dt \leq a \|x-y\|^{p+1} \int_0^1 t^p dt \\ &= a \|x-y\|^{p+1}/(p+1). \quad @ \end{aligned}$$

### Παρατηρήσεις

1. Τα αποτελέσματα αυτής της παραγράφου εύκολα γενικεύονται σε ευστήματα  $m$  εξισώσεων με  $n$  αγνώστους, σε απεικονίσεις δηλ.  $F$  της μορφής  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $F = (f_1, f_2, \dots, f_m)^T$ ,  $x = (x_1, \dots, x_n)^T$ . Ο ορισμός της παραγώγου (2) είναι ο ίδιος (αν με το ίδιο σύμβολο  $\|\cdot\|$  παραστήσουμε δύο οποιεσδήποτε νόρμες στον  $\mathbb{R}^n$  και στον  $\mathbb{R}^m$ ). Η παράγωγος  $F'(x)$  είναι για κάθε  $x \in D$ , γραμμικός τελεστής από τον  $\mathbb{R}^n$  στον  $\mathbb{R}^m$  (γράφουμε  $F'(x) \in L(\mathbb{R}^n, \mathbb{R}^m)$ ), και παριστάνεται πάλι - ως προς τις κανονικές βάσεις του  $\mathbb{R}^n$  και του  $\mathbb{R}^m$  - από τον  $m \times n$  Ιακωβιανό πίνακα  $J_{ij}(x) = \partial_j f_i(x)$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n$ . Στην ειδική περίπτωση ενός ευναρτησιακού, δηλ. μιάς απεικόνισης  $g: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^1$ , η παράγωγος  $g'(x)$  είναι γραμμικός τελεστής από τον  $\mathbb{R}^n$  στον  $\mathbb{R}^1$  και παριστάνεται από το γραμμοδιάνυσμα  $(\partial_1 g(x), \dots, \partial_n g(x))$ , δηλ. από το γραμμοδιάνυσμα  $\nabla g(x)$ . Είναι φανερό ότι οι Προτάσεις 1-4 ισχύουν για απεικονίσεις  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ , *mutatis mutandis*.

2. Για πολλές εφαρμογές η έννοια της παραγώγου Frechet που ορίσαμε ε' αυτήν την παράγραφο είναι πιά ισχυρή απ' ό,τι χρειάζεται. Μία ασθενέστερη έννοια είναι η λεγόμενη παράγωγος με την έννοια του Gateaux που ορίζεται ως εξής: Λέμε ότι μία απεικόνιση  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  είναι παραγωγίσιμη ε' ένα σημείο  $x \in \text{Int}(D)$  με την έννοια του Gateaux (ή  $\delta$ -παραγωγίσιμη) αν υπάρχει γραμμικός τελεστής  $A_x \in L(\mathbb{R}^n, \mathbb{R}^m)$  τέτοιος ώστε για κάθε  $h \in \mathbb{R}^n$  να ισχύει

$$(11) \quad \lim_{t \rightarrow 0} \|F(x+th) - F(x) - tA_x h\| / t = 0,$$

δηλ. αν είναι κατά κάποιο τρόπο "παραγωγίσιμη" σε κάθε κατεύθυνση  $h$ . Όπως και προηγουμένως, η ύπαρξη του τελεστή  $A_x$  είναι ανεξάρτητη της νόρμας  $\|\cdot\|$  και ο  $A_x$  είναι μοναδικός. Ορίζουμε λοιπόν για μία τέτοια ευναρτησιότητα την παραγωγή της  $F'(x)$  με την έννοια του Gateaux (ή  $\delta$ -παράγωγο) ως  $F'(x) = A_x$ . Είναι προφανές ότι η παραγωγή Gateaux

είναι γραμμική πράξη και ότι αν υπάρχει η  $G$ -παράγωγος  $F'(x)$ , τότε υπάρχουν οι μερικές παράγωγοι  $\partial_j f_i(x)$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n$ , η δε  $F'(x)$  παριετώνεται από τον Ιακωβιανό πίνακα  $J_{ij} = \partial_j f_i(x)$ . Είναι επίσης προφανές ότι αν η  $F$  είναι παραγωγίσιμη (κατά Frechet) και έχει  $(F-)$  παράγωγο  $F'(x)$ , τότε είναι παραγωγίσιμη κατά Gateaux και η  $G$ -παράγωγός της στο  $x$  συμπίπτει με την  $F'(x)$ . Το αντίστροφο όμως δεν είναι αληθινό: Παραγωγισιμότητα κατά Gateaux δεν συνεπάγεται παραγωγισιμότητα κατά Frechet (βλ. Ασκ. 3). Επίσης παραγωγισιμότητα κατά Gateaux ε' ένα σημείο  $x$  δεν συνεπάγεται συνέχεια της  $F$  στο  $x$  σε αντίθεση με ό,τι αποδείξαμε στην Πρόταση 1 για την παραγωγισιμότητα Frechet, (βλ. Ασκ. 3). Πάντως εύκολα βλέπουμε ότι οι Προτάσεις 2, 3 και 4 ισχύουν, αν στις υποθέσεις τους η  $F'(x)$  είναι η παράγωγος κατά Gateaux.

Για σειρά ασκήσεων (Ασκ. 2-6) στο τέλος της παραγράφου διερευνά έχουμε μεταξύ των εννοιών της παραγωγίσιμης Frechet, Gateaux και της συνεχιζόμενης (μερικής) παραγωγίσιμης ως προς  $x_1$ .

3. Οι ορισμοί και τα περιεωότερα αποτελέσματα αυτής της παραγράφου ισχύουν γενικά και για απεικονίσεις μεταξύ χώρων Banach. Έστω  $F: D \subset X \rightarrow Y$  όπου  $X, Y$  χώροι Banach. Τότε οι παράγωγοι  $F'(x)$  κατά Frechet ή κατά Gateaux ορίζονται ως (φραχμένοι) γραμμικοί τελεστές  $A_x: X \rightarrow Y$  από τις (2) και (11), αντίστοιχα. Οι Προτάσεις 1-4 εξακολουθούν να ισχύουν βέβαια το ολοκλήρωμα μιάς συνάρτησης  $G: [a, b] \subset X$  δεν ορίζεται πιά μέσω των συνιστωσών της  $G$ .

### Ασκήσεις 2.1

1. Με συντομία επιβεβαιώστε τους ισχυρισμούς εκείνους της παρατήρησης 2 σχετικά με την παράγωγο Gateaux για τους οποίους δεν υπάρχει ειδική παραπομπή σε άλλη άσκηση.

2. (α) θεωρείστε την συνάρτηση  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^1$  δίνεται από

$$f(x_1, x_2) = \begin{cases} x_1 & \text{αν } x_2 = 0 \\ x_2 & \text{αν } x_1 = 0 \\ 1 & \text{αλλιώς} \end{cases}$$

Δείξτε ότι οι μερικές παράγωγοι  $\partial_1 f(0)$  και  $\partial_2 f(0)$  υπάρχουν αλλά ότι η  $f$  δεν είναι συνεχής και δεν είναι παραγωγίσιμη κατά Gateaux (και ευγενώς ούτε και κατά Frechet) στο 0.

(β) θεωρείστε την συνάρτηση  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^1$  που δίνεται από

$$f(x_1, x_2) = \begin{cases} 0 & \text{αν } x=0 \\ x_1 x_2^2 / (x_1^2 + x_2^4) & \text{αν } x \neq 0. \end{cases}$$

Δείξτε ότι το όριο  $\lim_{t \rightarrow 0} [f(th) - f(0)]/t$  υπάρχει για κάθε  $h \in \mathbb{R}^2$  αλλά ότι η  $f$  δεν είναι παραγωγίσιμη κατά Gateaux (ούτε συνεχής) στο 0.

3. θεωρείστε την συνάρτηση  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^1$  που δίνεται από

$$f(x_1, x_2) = \begin{cases} 0 & \text{αν } x_1 = 0 \\ 2x_2 \exp(-x_1^{-2}) / (x_2^2 + \exp(-2x_1^{-2})) & \text{αν } x_1 \neq 0. \end{cases}$$

Δείξτε ότι έχει παράγωγο κατά Gateaux στο 0 αλλά ότι δεν είναι συνεχής (και ευγενώς ούτε παραγωγίσιμη κατά Frechet) στο 0.

4. Αν η  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  είναι παραγωγίσιμη κατά Gateaux στο κυρτό σύνολο  $D_0 \subset D$  δείξτε ότι για  $x, y, z \in D_0$

$$\|F(y) - F(z) - F'(x)(y-z)\| \leq \sup_{0 \leq t \leq 1} \|F'(z+t(y-z)) - F'(x)\| \|y-z\|$$

(Υπόδειξη: θεωρείτε την απεικόνιση  $G(u) = F(u) - F'(x)u$ ,  $u \in D$  και εφαρμόζοντας την Πρόταση 2 - για παραγωγούς Gateaux και για απεικονίσεις από τον  $\mathbb{R}^n$  στον  $\mathbb{R}^m$  - στην  $G$ , φράξτε το  $\|G(y) - G(z)\|$ ).

5. Υποθέστε ότι η  $F: D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$  έχει παράγωγο κατά Gateaux  $F'$  σε κάθε σημείο μιάς ανοικτής μπάλας με κέντρο  $x$  και ότι η  $F'$  είναι συνεχής στο  $x$ . Τότε η  $F$  είναι παραγωγίσιμη (κατά Fréchet) στο  $x$  και οι δύο παράγωγοι ευρύνονται. (Υπόδειξη: χρησιμοποιείστε την Άσκηση 4).

6. Υποθέστε ότι η  $F: D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$  έχει παράγωγο κατά Gateaux  $F'$  σε κάθε σημείο μιάς ανοικτής μπάλας με κέντρο  $x$ . Τότε η  $F'$  είναι συνεχής στο  $x$  αν και μόνο αν όλες οι μερικές παράγωγοι  $\partial_j f_i$  είναι συνεχείς στο  $x$ .

7. (α) Το θεώρημα της μέσης τιμής στη ευνηθισμένη του μορφή ισχύει για ευαρτησιακά: Υποθέστε ότι η  $g: D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^1$  έχει παράγωγο  $g'$  κατά Gateaux σε κάθε σημείο ενός κυρτού συνόλου  $D_0 \subset D$ . Τότε αν  $x, y \in D_0$ , υπάρχει  $t \in (0, 1)$  τέτοιο ώστε  $g(y) - g(x) = g'(x + t(y-x))(y-x)$ .

(β) θεωρείτε την απεικόνιση  $F: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  με συνιστώσες  $f_1(x) = x_1^3$ ,  $f_2(x) = x_2^2$ . Αν  $x=0$  και  $y=(1,1)^T$  δείξτε ότι δεν υπάρχει  $z$  της μορφής  $x+t(y-x)$ ,  $t \in [0,1]$  τέτοιο ώστε

$$(12) \quad F(y) - F(x) = F'(z)(y-x)$$

(γ) Έστω  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  μία διαγώνια απεικόνιση (δηλ. μία απεικόνιση  $F = (f_1, \dots, f_n)^T$  τέτοια ώστε η  $f_i$  να είναι συνάρτηση μόνο της μεταβλητής  $x_i$ ). Υποθέστε ότι η  $F$  είναι παραγωγίσιμη κατά Gateaux στον  $\mathbb{R}^n$ . Δείξτε ότι ισχύει η (12) αλλά με  $z$  όχι αναγκαστικά στο ευθύγραμμο τμήμα  $x+t(y-x)$ ,  $0 \leq t \leq 1$ .

8. (Σε ασκήσεις, επομένων παραγράφων θα αναφερθούμε ευχυνά στις υποθέσεις αυτής της άσκησης).

Στην θεωρία των μη γραμμικών ταλαντώσεων ευναντάμε ευχυνά το εξής πρόβλημα: Ζητάμε μία πραγματική απεικόνιση  $u(x)$ ,  $x \in [0,1]$ , δύο φορές ευνεχώς παραγωγίσιμη στο  $[0,1]$  που να ικανοποιεί το ευνοριακό πρόβλημα "δύο σημείων"

$$(*) \begin{cases} u''(x) = g(u(x)), & 0 \leq x \leq 1 \\ u(0) = \alpha, & u(1) = \beta, \end{cases}$$

όπου  $g: \mathbb{R}^1 \rightarrow \mathbb{R}^1$  μία δεδομένη δύο φορές ευνεχώς παραγωγίσιμη ευνάρτηση και  $\alpha, \beta$  δεδομένες πραγματικές σταθερές. Είναι γνωστό ότι αν η  $g$  ικανοποιεί την

$$(13) \quad g'(s) \geq \underline{g} > -n^2, \quad \forall s \in \mathbb{R}^1,$$

τότε υπάρχει μοναδική λύση  $u(x)$  του (\*).

Γιά να επιλύσουμε προερχηστικά το πρόβλημα (\*) με μία μέθοδο πεπερασμένων διαφορών, εισάχουμε τον ομοιόμορφο διαμερισμό  $x_i = ih$ ,  $0 \leq i \leq n+1$  όπου  $(n+1)h=1$ ,  $n \in \mathbb{N}$ . Ζητάμε να κατασκευάσουμε προερχήσεις  $U_i$ ,  $0 \leq i \leq n+1$ , των τιμών της λύσης  $u(x_i)$ ,  $0 \leq i \leq n+1$ , που ορίζονται ως εξής:

$$(14) \quad \begin{cases} (U_{i-1} - 2U_i + U_{i+1})/h^2 = g(U_i), & 1 \leq i \leq n, \\ U_0 = \alpha, & U_{n+1} = \beta. \end{cases}$$

Δηλ. το διάνυσμα  $U = (U_1, \dots, U_n)^T$  των αγνώστων ικανοποιεί το  $n \times n$  μη γραμμικό σύστημα

$$(15) \quad F(U) = 0,$$

όπου  $F(U) = [f_1(U), \dots, f_n(U)]^T$  και όπου οι ευναρτήσεις  $f_i: \mathbb{R}^n \rightarrow \mathbb{R}^1$  ορίζονται ως

$$(16) \quad \begin{cases} f_1(U) = -\alpha + 2U_1 - U_2 + h^2 g(U_1) \\ f_i(U) = -U_{i-1} + 2U_i - U_{i+1} + h^2 g(U_i), \quad 2 \leq i \leq n-1 \\ f_n(U) = -U_{n-1} + 2U_n - \beta + h^2 g(U_n). \end{cases}$$

Εισάγοντας τον τριδιαγώνιο πίνακα

$$(17) \quad A = [-1, 2, -1]$$

και την διαγώνια απεικόνιση  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^n$  που ορίζεται για  $x = (x_1, \dots, x_n)^T \in \mathbb{R}^n$  ως

$$(18) \quad \phi(x) = h^2 (g(x_1) - \alpha h^{-2}, g(x_2), \dots, g(x_{n-1}), g(x_n) - \beta h^{-2})^T,$$

βλέπουμε ότι το σύστημα (15) γράφεται και στην μορφή

$$(19) \quad F(U) \equiv AU + \phi(U) = 0.$$

Ερώτημα: Υποθέστε ότι η  $g: \mathbb{R}^1 \rightarrow \mathbb{R}^1$  είναι συνεχής. Δείξτε ότι η απεικόνιση  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συνεχώς παραγωγίσιμη (κατά Frechet) στον  $\mathbb{R}^n$  και ότι

$$(20) \quad F'(x) = A + \Phi'(x), \quad x \in \mathbb{R}^n$$

όπου η  $\Phi'$  είναι η παράγωγος (Frechet) της  $\phi$  που παριστάνεται από τον διαγώνιο Ιακωβιανό πίνακα

$$(21) \quad J_\phi(x) = h^2 \text{diag}(g'(x_1), \dots, g'(x_n)), \quad x \in \mathbb{R}^n.$$

## 2.2 ΤΟΠΙΚΑ ΘΕΩΡΗΜΑΤΑ ΣΥΓΚΛΙΣΗΣ. ΤΟ ΘΕΩΡΗΜΑ ΤΗΣ ΣΥΣΤΟΛΗΣ

Θα ασχοληθούμε στη συνέχεια με την κατασκευή προσεγγίσεων των λύσεων του μη γραμμικού συστήματος

$$(1) \quad F(x) = 0,$$

όπου  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι μία δεδομένη απεικόνιση. Οι αριθμητικές μας μέθοδοι θα είναι επαναληπτικές, θα παράχουν δηλ. μία ακολουθία  $\{x^k\}$ ,  $k \geq 0$ , διανυσμάτων του  $\mathbb{R}^n$ , τέτοια ώστε  $x^k \rightarrow x^*$ ,  $k \rightarrow \infty$ , όπου  $x^*$  είναι κάποια λύση του (1). Υπάρχουν βασικά τρία σημαντικά ερωτήματα για κάθε τέτοια μέθοδο:

Πρώτα - πρώτα, οι προσεγγίσεις  $x^k$ ,  $k \geq 0$ , πρέπει να είναι καλά ορισμένες. Υπάρχει δηλ. το ερώτημα της κατάλληλης επιλογής της αρχικής τιμής  $x^0$  (γενικότερα πιθανώς των αρχικών τιμών  $x^0, x^1, \dots, x^p$ ) και της κατασκευής των  $x^k$ ,  $k \geq 0$ , έτσι ώστε να βρίσκονται μέσα στα πεδία ορισμού των απεικονίσεων που πρόκειται να υπολογισθούν επί  $x^k$ .

Το δεύτερο ερώτημα αφορά την εύγκλιση της ακολουθίας  $\{x^k\}$  και το αν το όριό της είναι πράγματι λύση του (1). Εδώ διακρίνουμε διαφόρων τύπων θεώρημα εύγκλισης. Π.χ. ένα θεώρημα τοπικής εύγκλισης ("τοπικό θεώρημα") τυπικά υποθέτει ότι υπάρχει μία λύση  $x^*$  του (1) και διαβεβαιώνει ότι υπάρχει μία περιοχή  $S$  του  $x^*$  τέτοια ώστε αν η αρχική τιμή επιλεγεί μέσα στην  $S$ , τότε η ακολουθία  $\{x^k\}$  είναι καλά ορισμένη και ευγκλίνει στο  $x^*$ . Ένα θεώρημα περιορισμένης εύγκλισης δεν υποθέτει την ύπαρξη λύσης  $x^*$  του (1), αλλά διαβεβαιώνει ότι για μία ειδική (περιορισμένη γενικά) επιλογή αρχικών τιμών, η  $\{x^k\}$  είναι καλά ορισμένη και η εύγκλιση της εξασφαλισμένη σε κάποιο όριο  $x^*$  που είναι λύση της (1). Βέβαια τα πιο επιθυμητά αποτελέσματα είναι εκείνα που για ορισμένου τύπου απεικονίσεις  $F$  εξασφαλίζουν ότι οποιαδήποτε επιλογή τιμών στον  $\mathbb{R}^n$  (ή τουλάχιστον ε' ένα μεγάλο υποσύνολό του) δίνει μία καλά ορισμένη ακολουθία, της οποίας το όριο είναι λύση του (1).

Ένα τρίτο πρόβλημα είναι η ταχύτητα εύγκλισης της ακολουθίας  $\{x^k\}$  στο  $x^*$ . Σχετικές πληροφορίες μπορεί να μας δώσει μία εκτίμηση του εφάλματος  $\|x^k - x^*\|$  αν και συνήθως τέτοιες εκτιμήσεις είναι μάλλον



περιμετρικές. Γιαυτό ενδιαφερόμαστε πολλές φορές για την "ασυμπτωτική ταχύτητα σύγκλισης" δηλ: την συμπεριφορά του σφάλματος για μεγάλο  $k$ .

Εκτός από τα παραπάνω θεωρητικά ερωτήματα σημασία στην πράξη έχουν αποτελεσματικοί τρόποι αναζήτησης αρχικών τιμών, εκτέλεσης των βημάτων του αλγορίθμου και τερματισμού του με κατάλληλα κριτήρια.

θα αναλύσουμε πρώτα μία γενική επαναληπτική μέθοδο για την προεχχιστική επίλυση του (1). Υποθέτουμε ότι ένα  $x \in \mathbb{R}^n$  είναι λύση του (1) αν και μόνο αν ικανοποιεί την εξίσωση  $x = G(x)$  όπου  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι μία άλλη απεικόνιση, δηλ. ότι κάθε λύση του (1) είναι ένα εταθερό σημείο της απεικόνισης  $G$ . Τότε η γενική επαναληπτική μέθοδος

$$(2) \quad x^{k+1} = G(x^k), \quad k=0,1,2,\dots, \quad x^0 \in \mathbb{R}^n \text{ δεδομένο,}$$

είναι μία προφανής μέθοδος για την προέχχιση των εταθερών σημείων της  $G$ . Πράγματι αν  $x^k \in D$ ,  $k \geq 0$ , αν η  $G$  είναι συνεχής και αν  $x^k \rightarrow x^*$ ,  $k \rightarrow \infty$  όπου  $x^* \in D$  τότε το  $x^*$  είναι εταθερό σημείο της  $G$ .

Λέμε ότι ένα σημείο  $x^* \in D$  είναι σημείο έλξεως για την ακολουθία που παράγεται από την αναδρομική σχέση (επαναληπτική μέθοδο) (2) αν υπάρχει (ανοιχτή) περιοχή  $S \subset D$  του  $x^*$  τέτοια ώστε για κάθε  $x^0 \in S$  τα  $x^k$  που ορίζονται από την (2) να βρίσκονται στην  $S$  και να συγκλίνουν στο  $x^*$  όταν  $k \rightarrow \infty$ . Συμβολίζουμε με  $S(x, \delta)$  την ανοιχτή μπάλα (για κάποια νόρμα) κέντρου  $x$  και ακτίνας  $\delta$  και με  $\bar{S}(x, \delta)$  την αντίστοιχη κλειστή μπάλα. Έχουμε το εξής γενικό αποτέλεσμα τοπικής σύγκλισης:

**Λήμμα 1.** Έστω  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  και  $x^* \in D$ . Υποθέτουμε ότι υπάρχει ανοιχτή μπάλα  $S = S(x^*, \delta) \subset D$  και εταθερά  $\alpha < 1$  τέτοια ώστε

$$(3) \quad \|G(x) - x^*\| \leq \alpha \|x - x^*\|, \quad \forall x \in S.$$

Τότε, για κάθε  $x^0 \in S$ , οι όροι  $x^k$  της ακολουθίας που ορίζονται από την (2) βρίσκονται στην  $S$  και συγκλίνουν στο  $x^*$ , δηλ. το  $x^*$  είναι σημείο έλξεως της (2).

Απόδειξη: Έστω  $x^0 \in S$ . Η απόδειξη είναι επαγωγική. Για  $k=1$   $\|x^1 - x^*\| = \|G(x^0) - x^*\| \leq$  (από την (3))  $a\|x^0 - x^*\|$ . Επειδή  $a < 1$  και  $\|x^0 - x^*\| < \delta$  έχουμε ότι  $\|x^1 - x^*\| < \delta$  δηλ. ότι  $x^1 \in S$ . Υποθέτοντας τώρα ότι για κάποιο  $k$ ,  $x^k \in S$  και  $\|x^k - x^*\| \leq a^k \|x^0 - x^*\|$ , έχουμε ότι  $\|x^{k+1} - x^*\| = \|G(x^k) - x^*\| \leq a\|x^k - x^*\| \leq a^{k+1} \|x^0 - x^*\|$  επίσης  $x^{k+1} \in S$  ό.έ.δ. Άρα,  $x^0 \in S \Rightarrow x^k \in S$  και  $\|x^k - x^*\| \leq a^k \|x^0 - x^*\|$  για  $k \geq 0 \Rightarrow x^k \rightarrow x^*$ ,  $k \rightarrow \infty$ , δηλ. το  $x^*$  είναι σημείο έλξης για την (2). @

Πρίν χρησιμοποιήσουμε το Λήμμα 1 για να διατυπώσουμε μία γενική ικανή συνθήκη για τοπική εύγκλιση της (2) θα αποδείξουμε ένα βοηθητικό αποτέλεσμα από την γραμμική άλγεβρα. Έστω  $A \in L(\mathbb{C}^n)$  και  $\lambda_i(A)$ ,  $1 \leq i \leq n$  οι ιδιοτιμές του. Η φασματική ακτίνα  $\rho(A)$  του  $A$  ορίζεται ως  $\rho(A) = \max_i |\lambda_i(A)|$ . Είναι προφανές ότι για κάθε νόρμα  $\|\cdot\|$  ισχύει  $\rho(A) \leq \|A\|$ . Μας ενδιαφέρει το εξής, εφεδόν αντίστροφο, αποτέλεσμα:

**Λήμμα 2.** Έστω  $A \in L(\mathbb{C}^n)$ . Τότε για κάθε  $\epsilon > 0$  υπάρχει νόρμα  $\|\cdot\|$  του  $\mathbb{C}^n$  τέτοια ώστε

$$(4) \quad \|A\| \leq \rho(A) + \epsilon.$$

Απόδειξη: Ταυτίζοντας γραμμικούς τελεστές με τους πίνακες που τους παριστάνουν ως προς την κανονική βάση  $\{e^j\}$  του  $\mathbb{C}^n$ , έχουμε  $P^{-1}AP = J$  όπου  $J \in L(\mathbb{C}^n)$  είναι η μορφή Jordan του  $A$ , η οποία, ως γνωστόν, παριστάνεται από ένα διαγώνιο πίνακα τετραγωνικών υποπίνακων  $J = \text{diag}\{J_1, \dots, J_m\}$  όπου κάθε υποπίνακας  $J_i$  είναι είτε ο  $1 \times 1$  πίνακας  $(\lambda_i)$  ή είναι διδιαγώνιος πίνακας της μορφής

$$J_i = \begin{bmatrix} \lambda_i & 1 & & 0 \\ & & \ddots & \\ & & & 1 \\ & 0 & & & \lambda_i \end{bmatrix}$$

όπου  $\lambda_i$  οι ιδιοτιμές του  $A$ . θεωρούμε τον διαγώνιο  $n \times n$  πίνακα

$D = \text{diag}(1, \varepsilon, \varepsilon^2, \dots, \varepsilon^{n-1})$ . Εύκολα βλέπουμε ότι ο πίνακας  $\tilde{J} = D^{-1}JD$  είναι ο ίδιος με τον  $J$  εκτός από το ότι οι μονάδες της πρώτης υπερδιαγωνίου του  $J$  αντικαθίστανται από  $\varepsilon$ . Συνεπώς  $\|\tilde{J}\|_{\infty} \leq \rho(A) + \varepsilon$ .

Θέτουμε  $Q = PD$  και ορίζουμε την νόρμα  $\|\cdot\|$  στον  $\mathbb{C}^n$  από την εξέση  $\|x\| = \|Q^{-1}x\|_{\infty}$ . Τότε  $\|A\| = \max_{\|x\|=1} \|Ax\| = \max_{\|Q^{-1}x\|_{\infty}=1} \|Q^{-1}Ax\|_{\infty} = \max_{\|y\|_{\infty}=1} \|Q^{-1}AQy\|_{\infty} =$

$$= \max_{\|y\|_{\infty}=1} \|\tilde{J}y\|_{\infty} = \|\tilde{J}\|_{\infty} \leq \rho(A) + \varepsilon. \quad @$$

$$\|y\|_{\infty}=1$$

Το παρακάτω θεώρημα διατυπώνει μιά χρήσιμη ικανή συνθήκη για την ισχύ της (3).

**ΘΕΩΡΗΜΑ 1.** (Ostrowski). Υποθέτουμε ότι η  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  έχει ένα σταθερό σημείο  $x^* \in \text{Int}(D)$  και η  $G$  έχει παράγωγο (Frechet)  $G'(x^*)$  στο  $x^*$ . Τότε, αν  $\rho(G'(x^*)) \leq \varepsilon < 1$ , το  $x^*$  είναι σημείο έλξης της (2).

Απόδειξη: Από το λήμμα 2 και την υπόθεσή μας έπεται ότι για  $\varepsilon > 0$  υπάρχει νόρμα  $\|\cdot\|$  τέτοια ώστε  $\|G'(x^*)\| \leq \varepsilon + \varepsilon$ . Επειδή η  $G$  είναι παραγωγίσιμη στο  $x^*$  έχουμε ότι  $\exists \delta = \delta(x^*, \varepsilon) > 0$  τέτοιο ώστε  $S = S(x^*, \delta) \subset D$  και  $\|G(x) - G(x^*) - G'(x^*)(x - x^*)\| \leq \varepsilon \|x - x^*\|$ ,  $\forall x \in S$ . Από τις δύο αυτές εξέσεις και το γεγονός ότι  $G(x^*) = x^*$  παίρνουμε

$$\begin{aligned} \|G(x) - x^*\| &\leq \|G(x) - G(x^*) - G'(x^*)(x - x^*)\| + \|G'(x^*)(x - x^*)\| \\ &\leq (\varepsilon + 2\varepsilon) \|x - x^*\|. \end{aligned}$$

Διαλέγοντας λοιπόν  $\varepsilon > 0$  τέτοιο ώστε  $\alpha\varepsilon + 2\varepsilon < 1$  βλέπουμε ότι υπάρχει νόρμα  $\|\cdot\|$  τέτοια ώστε  $\alpha < 1$  και

$$\|G(x) - x^*\| \leq \alpha \|x - x^*\| \quad \forall x \in S(x^*, \delta),$$

δηλ. ότι ικανοποιούνται οι υποθέσεις του λήμματος 1. Συνεπώς το  $x^*$  είναι σημείο έλξης της (2). @

Ύστερα από αυτό το τυπικό δείγμα θεωρήματος τοπικής σύγκλισης προχωρούμε στην διατύπωση ενός αποτελέσματος του τύπου "περιορισμένης σύγκλισης", του γνωστού μας "θεωρήματος της συστολής". Λέμε ότι η απεικόνιση  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συστολή ε' ένα εύνολο  $D_0 \subset D$  αν υπάρχει νόρμα  $\|\cdot\|$  του  $\mathbb{R}^n$  και  $\alpha < 1$  τέτοια ώστε

$$(5) \quad \|G(x) - G(y)\| \leq \alpha \|x - y\|, \quad \forall x, y \in D_0.$$

**ΘΕΩΡΗΜΑ 2 (Συστολής).** Υποθέτουμε ότι η απεικόνιση  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συστολή ε' ένα κλειστό εύνολο  $D_0 \subset D$  και ότι  $G(D_0) \subset D_0$ . Τότε η  $G$  έχει ένα μοναδικό σταθερό σημείο  $x^*$  στο  $D_0$ . Επιπλέον, για οποιοδήποτε  $x^0 \in D_0$ , η ακολουθία  $\{x^k\}$ ,  $k \geq 0$ , που παράγεται από την (2), συγκλίνει στο  $x^*$ . Πάλιετα, για  $\|\cdot\|$ ,  $\alpha$  όπως στην (5), ισχύουν οι εκτιμήσεις

$$(6) \quad \|x^k - x^*\| \leq (\alpha / (1 - \alpha)) \|x^k - x^{k-1}\|, \quad k = 1, 2, 3, \dots,$$

$$(7) \quad \|x^k - x^*\| \leq (\alpha^k / (1 - \alpha)) \|G(x^0) - x^0\|, \quad k = 1, 2, 3, \dots$$

Απόδειξη: Έστω  $x^0 \in D_0$ . Επειδή  $G(D_0) \subset D_0$ , η ακολουθία  $x^{k+1} = G(x^k)$ ,  $k = 0, 1, 2, \dots$ , είναι καλά ορισμένη και ανήκει στο  $D_0$ .

Επιπλέον  $\|x^{k+1} - x^k\| = \|G(x^k) - G(x^{k-1})\| \leq \alpha \|x^k - x^{k-1}\|$ ,  $k \geq 1$  από την οποία για οποιοδήποτε  $m \geq 1$  ακέραιο παίρνουμε για  $k = 1, 2, 3, \dots$

$$(8) \quad \|x^{k+m} - x^k\| \leq \sum_{i=1}^m \|x^{k+i} - x^{k+i-1}\| \leq (\alpha^{m-1} + \dots + \alpha) \|x^{k+1} - x^k\|$$

$$\leq (1 - \alpha)^{-1} \|x^{k+1} - x^k\| \leq \alpha (1 - \alpha)^{-1} \|x^k - x^{k-1}\| \leq \dots$$

$$\leq (\alpha^k / (1 - \alpha)) \|x^1 - x^0\|.$$

Συμπεραίνουμε ότι η ακολουθία  $\{x^k\}$  είναι Cauchy ως προς την νόρμα  $\|\cdot\|$  στο κλειστό εύνολο  $D_0$ . Συνεπώς  $\exists x^* \in D_0$  τέτοιο ώστε  $x^k \rightarrow x^*$ ,  $k \rightarrow \infty$ . Το όριο αυτό είναι σταθερό σημείο της  $G$  γιατί

$$\begin{aligned} \|x^* - G(x^*)\| &\leq \|x^* - x^{k+1}\| + \|G(x^k) - G(x^*)\| \leq \|x^* - x^{k+1}\| \\ &+ a \|x^k - x^*\| \rightarrow 0, \quad k \rightarrow \infty, \end{aligned}$$

είναι δε το μοναδικό σταθερό σημείο της  $G$  στο  $D_0$  διότι αν υπήρχε και άλλο  $x^{**}$  θα είχαμε

$$\|x^* - x^{**}\| = \|G(x^*) - G(x^{**})\| \leq a \|x^* - x^{**}\| < \|x^* - x^{**}\|,$$

άτοπο. Οι εκτιμήσεις (6), (7) προκύπτουν τώρα από την (8) παίρνοντας τό όριο  $n \rightarrow \infty$ . @

### Παρατηρήσεις

1. Το θεώρημα της ευστολής (και η απόδειξή του ε' αυτήν την παράγραφο) ισχύει και για απεικονίσεις  $G: D \subset X \rightarrow Y$ , όπου  $X, Y$  χώροι Banach. Γενικότερα, αν ο  $X$  πλήρης μετρικός χώρος με μετρική  $d(\cdot, \cdot)$  και  $G: X \rightarrow X$  είναι μία απεικόνιση για την οποία υπάρχει  $a < 1$  τέτοιο ώστε  $d(G(x), G(y)) \leq ad(x, y)$ ,  $\forall x, y \in X$ , τότε η  $G$  έχει ένα μοναδικό σταθερό σημείο στον  $X$ . Υπάρχει και η εξής ενδιαφέρουσα επέκταση: λέμε ότι μία απεικόνιση  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  δεν είναι διαστολή στο εύνολο  $D_0 \subset D$  αν ικανοποιεί μία συνθήκη Lipschitz με σταθερά ίση με την μονάδα στο  $D_0$ , δηλ. αν

$$(9) \quad \|G(x) - G(y)\| \leq \|x - y\|, \quad \forall x, y \in D_0.$$

Αν ισχύει η (9), αν το  $D_0$  είναι κλειστό και κυρτό και αν  $G(D_0) \subset D_0$ , τότε η  $G$  έχει σταθερό σημείο  $x^*$  στο  $D_0$  αν και μόνο αν η ακολουθία

(2) είναι φραγμένη για τουλάχιστον ένα  $x^0 \in D_0$  (Στο  $x^*$  συγκλίνει τότε

μία υποκολουθία της  $x^k$ ). Για την απόδειξη του θεωρήματος αυτού (που αποτελεί ειδική περίπτωση γενικότερου θεωρήματος του Brouwer) βλ. [2.3, εελ. 121 και Παρ. 6.3].

2. Οι εκτιμήσεις του εφάλματος (6) και (7) είναι προφανώς πολύ χρήσιμες αν είναι γνωστή η τιμή του  $a$ . Τότε π.χ. μπορούμε να υπολογίσουμε ακριβώς εκ των προτέρων τον αριθμό των βημάτων που απαιτούνται έτσι ώστε το εφάλμα να γίνει μικρότερο δεδομένου  $\epsilon > 0$ . Το πρόβλημα είναι βέβαια στην πράξη ο ακριβής κατά το δυνατόν προσδιορισμός της σταθεράς συστολής  $a$ .

3. Ειδική περίπτωση απεικόνισης  $G$  είναι και η "αφινική" απεικόνιση (ε' όλου του  $\mathbb{R}^n$ )  $x \mapsto Hx+d$ , όπου  $d \in \mathbb{R}^n$  σταθερό. Σ' αυτήν την περίπτωση ικανή και αναγκαία συνθήκη για την σύγκλιση της ακολουθίας  $x^{k+1} = Hx^k + d$  στην (υποτιθέμενη μοναδική) λύση  $x$  του συστήματος  $x = Hx + d$  για οποιοδήποτε  $x^0 \in \mathbb{R}^n$  είναι να υπάρχει νόρμα  $\|\cdot\|$  του  $\mathbb{R}^n$  τέτοια ώστε  $\|H\| < 1$ , δηλ. να είναι η  $G$  συστολή (εδώ  $a = \|H\|$ ). Ισοδύναμη ικανή και αναγκαία συνθήκη είναι (βλ. Λήμμα 2) η  $\rho(H) < 1$ . (Αυτό είναι το βασικό αποτέλεσμα σύγκλισης των κλασικών επαναληπτικών μεθόδων για την λύση του γραμμικού συστήματος  $Ax=b$  το οποίο γιαυτόν του σκοπό γράφουμε στην μορφή  $Mx = Nx + b$ , όπου  $A = M - N$ ,  $M$  αντιστρέψιμος, δηλ. στην μορφή  $x = Hx + d$  με  $H = M^{-1}N$ ,  $d = M^{-1}b$ . Βλέπε [5.4, Παρ. 11-13]).

## Άσκησης 2.2

1. (α) Να αποδειχθούν οι ισχυρισμοί για την σύγκλιση της  $x^{k+1} = Hx^k + d$  της Παρατήρησης 3.

(β) Δώστε παράδειγμα ενός γραμμικού τελεστή  $H \in L(\mathbb{R}^2)$  που είναι συστολή ως προς μία νόρμα του  $\mathbb{R}^2$  αλλά δεν είναι ως προς μία άλλη.

2. (α) Ορίζουμε  $f: [0,1] \rightarrow \mathbb{R}^1$  ως  $f(x) = (x/2) + 2$ . Δείξτε ότι η  $f$  είναι συστολή στο  $[0,1]$  αλλά δεν έχει σταθερό σημείο στο  $[0,1]$ . Τι συμβαίνει;

(β) Υποθέστε ότι η  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  έχει παράγωγο με την έννοια του Gateaux που ικανοποιεί  $\|G'(x)\| \leq \alpha < 1$  για κάθε  $x \in D$  (συνήθως έπει εκτιμούμε στην πράξη την τιμή της σταθεράς ευστολής).

(γ) Υποθέστε ότι η  $G: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συνεχώς παραγωγίσιμη (κατά Frechet) σε μία (ανοιχτή) περιοχή  $S_1$  ενός σημείου  $x \in \text{Int}(D)$  και ότι  $\rho(F'(x)) < 1$ . Δείξτε ότι υπάρχει περιοχή  $S_2$  του  $x$  και νόρμα  $\|\cdot\|$  του  $\mathbb{R}^n$  τέτοιες ώστε η  $G$  να είναι ευστολή στην  $S_2$  ως προς  $\|\cdot\|$ .

3. Στο θεώρημα 1 του Ostrowski δείξτε ότι η συνθήκη  $\rho(G'(x^*)) < 1$  δεν είναι αναγκαία για σύγκλιση. (Υπόδειξη: ερευνήστε τις ιδιότητες του σημείου  $x^* = 0$  για τις απεικονίσεις  $x \mapsto x \pm x^3$  του  $\mathbb{R}^1$  στον  $\mathbb{R}^1$ ).

4. Υποθέστε ότι η  $G: \mathbb{R}^n \rightarrow \mathbb{R}^n$  έχει την ιδιότητα ότι για κάθε συμπαγές (κλειστό και φραγμένο) σύνολο  $C$  του  $\mathbb{R}^n$  υπάρχει σταθερά  $\alpha_C < 1$  τέτοια ώστε

$$(10) \quad \|G(x) - G(y)\| \leq \alpha_C \|x - y\|, \quad \forall x, y \in C$$

και υποθέστε ότι η  $G$  έχει σταθερό σημείο  $x^* \in \mathbb{R}^n$ . Δείξτε ότι το  $x^*$  είναι μοναδικό στον  $\mathbb{R}^n$  και ότι η ακολουθία (2) συγκλίνει στο  $x^*$  για κάθε  $x^0 \in \mathbb{R}^n$ . Μπορείτε να βρείτε ένα παράδειγμα (π.χ. στον  $\mathbb{R}^1$ ) που να δείχνει ότι μόνο η συνθήκη (10) δεν εξασφαλίζει την ύπαρξη σταθερού σημείου για την  $G$ ;

5. (α) Έστω ότι η (9) ισχύει σε ένα σύνολο  $D_0 \subset D$  ως αυστηρή ανισότητα. Δείξτε τότε ότι η  $G$  έχει το πολύ ένα σταθερό σημείο στον  $D_0$ . Μόνο όμως η (9), έστω και αν ισχύει ως αυστηρή ανισότητα, δεν είναι ικανή να εξασφαλίζει την ύπαρξη σταθερού σημείου:

εξετάστε την συνάρτηση  $g: \mathbb{R}^1 \rightarrow \mathbb{R}^1$

$$g(x) = \begin{cases} x + e^{-x/2}, & \text{αν } x \geq 0 \\ e^{x/2}, & \text{αν } x < 0. \end{cases}$$

για  $D_0 \subset [0, \infty)$  ή  $D_0 \subset (-\infty, 0]$

(β) θεωρείστε την συνάρτηση  $g(x) = -x$  με πεδίο ορισμού  $D_0 = [-1, 1]$ .

Δείξτε ότι δεν είναι διαστολή στο  $D_0$ , και επιπλέον ότι ικανοποιούνται οι υπόλοιπες συνθήκες του θεωρήματος που αναφέρεται στην παρατήρηση 1. Συνεπώς η ύπαρξη σταθερού σημείου ( $x^* = 0$ ) είναι εξασφαλισμένη στο  $D_0$ . Τί συμβαίνει αν θεωρήσουμε το πεδίο ορισμού  $D_0 = [-1, -1/2] \cup [1/2, 1]$ ;

6. Λέμε ότι η απεικόνιση  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι ομοιομορφισμός στον  $\mathbb{R}^n$  αν είναι 1 προς 1 και επί και αν οι  $F$  και  $F^{-1}$  είναι συνεχείς στον  $\mathbb{R}^n$ . Αποδείξτε το εξής μη γραμμικό ανάλογο του θεωρήματος του Neumann: Έστω  $F = I - G$  όπου  $G: \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συστολή στον  $\mathbb{R}^n$  και όπου  $I$  είναι η ταυτότητα στον  $\mathbb{R}^n$ . Δείξτε ότι η  $F$  είναι ομοιομορφισμός στον  $\mathbb{R}^n$ .

7. (Η άσκηση αυτή αποτελεί συνέχεια της Άσκησης 2.1.8)

(α) Η πιο προφανής μέθοδος για την επίλυση του συστήματος (2.1.19) είναι μία γενική επαναληπτική μέθοδος. Κατ' αρχήν δείξτε ότι ο  $A$  είναι θετικά ορισμένος. θεωρείστε την απεικόνισή  $P: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$(11) P(x) = -A^{-1}\phi(x), \quad x \in \mathbb{R}^n.$$

Προφανώς, κάθε λύση του (2.1.19) είναι σταθερό σημείο της  $P$  και αντίστροφα. Δείξτε ότι αν η  $g': \mathbb{R}^1 \rightarrow \mathbb{R}^1$  είναι συνεχής, τότε η  $P$  είναι συνεχώς παραγωγίσιμη (κατά Fréchet) στον  $\mathbb{R}^n$  και ότι η παράγωγός της δίδεται από  $P' = -A^{-1}\phi'$ . Υποθέστε επιπλέον ότι η  $g'$  ικανοποιεί τις συνθήκες

$$(12) 0 \leq g'(s) \leq b < \infty \quad \forall s \in \mathbb{R}^1.$$



Υπολογίζοντας ένα κατάλληλο άνω φράγμα της ευκλείδειας νόρμας  $\|P'(x)\|_2$  δείξτε ότι αν

$$(13) \quad b < n^2$$

τότε, για  $h$  αρκετά μικρό, η  $P$  είναι ευστολή (ως προς την νόρμα  $\|\cdot\|_2$ ) στον  $\mathbb{R}^n$ . Συνεπώς η ακολουθία που ορίζεται από την επανάληψη

$$(14) \quad x^{k+1} = P(x^k), \quad k \geq 0, \quad x^0 \in \mathbb{R}^n \text{ αυθαίρετο,}$$

ευκλίνει στο μοναδικό (κάτω από αυτές τις συνθήκες) σταθερό σημείο της  $P$  στον  $\mathbb{R}^n$ , δηλ. στη μοναδική λύση του (2.1.19). (Υπόδειξη: βρείτε τις ιδιοτιμές του  $A$  και υπολογίστε το  $\|A^{-1}\|_2$  συναρτήσει του  $h$ ).

(β) Η συνθήκη (13) και το γεγονός ότι το  $h$  πρέπει να είναι αρκετά μικρό (δηλ. το  $n$  αρκετά μεγάλο) για να ευκλίνει η (14) κάνει την επιλογή του  $P = -A^{-1}\Phi$  ως απεικόνιση επανάληψης όχι επιτυχή. Υπάρχει όμως η εξής εναλλακτική λύση: Για  $\gamma > 0$  σταθερό, θεωρείστε την απεικόνιση  $G: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$(15) \quad G(x) = (A + \gamma I)^{-1}(\gamma x - \Phi(x)), \quad x \in \mathbb{R}^n.$$

Προφανώς κάθε σταθερό σημείο της  $G$  είναι λύση του (2.1.19) και αντίστροφα. Υποθέστε ότι η  $g'$  είναι συνεχής στον  $\mathbb{R}^1$  και ικανοποιεί την (12) για οποιαδήποτε  $b > 0$ . Διαλέξτε  $\gamma = h^2 b / 2$ . Δείξτε τότε ότι η  $G$  είναι ευστολή (ως προς την νόρμα  $\|\cdot\|_2$ ) στον  $\mathbb{R}^n$ . Συνεπώς η επανάληψη  $x^{k+1} = G(x^k)$ ,  $k \geq 0$  ευκλίνει για κάθε  $x^0 \in \mathbb{R}^n$  στην μοναδική (κάτω από αυτές τις συνθήκες) λύση του (2.1.19).

(γ) Υποθέστε τώρα ότι αντί της (12) ισχύει απλώς ότι

$$(16) \quad 0 \leq g'(s), \quad \forall s \in \mathbb{R}^1.$$

Δείξτε τότε ότι για οποιαδήποτε τιμή του  $\gamma$  η  $G$  που ορίζεται από την (15) δεν είναι αναγκαστικά ευστολή στον  $\mathbb{R}^n$ .

### 2.3 ΜΕΘΟΔΟΣ ΤΟΥ ΝΕΥΤΩΝΑ: ΤΟΠΙΚΗ ΣΥΓΚΛΙΣΗ ΚΑΙ ΤΑΧΥΤΗΤΑ ΣΥΓΚΛΙΣΗΣ

Σ' αυτήν την παράγραφο θα μελετήσουμε εξειδικεύσεις της γενικής επαναληπτικής μεθόδου (2.2.2) για την λύση του ελατήματος (2.2.1) στην περίπτωση που η απεικόνιση  $G$  είναι της μορφής

$$(1) \quad G(x) = x - A(x)^{-1} F(x)$$

όπου υποθέτουμε ότι για  $x$  σε κατάλληλο υποέυναλο του  $\mathbb{R}^n$ , ο  $A(x) \in L(\mathbb{R}^n)$  είναι αντιστρέψιμος γραμμικός τελεστής. Θα μελετήσουμε δηλ. επαναλήψεις της μορφής

$$(2) \quad x^{k+1} = x^k - A(x^k)^{-1} F(x^k), \quad k \geq 0, \quad x^0 \in \mathbb{R}^n \text{ δεδομένο.}$$

Η σημαντικότερη ειδική περίπτωση των (1), (2) είναι η μέθοδος του Νεύτωνα για την οποία  $A(x) = F'(x)$ , όπου  $F'$  είναι η παράγωγος (Frechet) της  $F$ . Η μέθοδος αυτή γενικεύει την γνωστή μας από τον  $\mathbb{R}^1$  μέθοδο του Νεύτωνα ("μέθοδος των Newton-Raphson")  $x^{k+1} = x^k - f(x^k)/f'(x^k)$ .

Θα αποδείξουμε πρώτα ένα προκαταρκτικό αποτέλεσμα (θεώρημα 1) για επαναληπτικές μεθόδους της μορφής (2) το οποίο θα χρησιμοποιήσουμε για να δείξουμε εν συνεχεία ένα θεώρημα τοπικής εύγκλισης για την μέθοδο του Νεύτωνα. Πρώτα αναφέρουμε την εξής γενίκευση της πρότασης του Neumann (βλ. (1.2.3)) η οποία προκύπτει άμεσα από την (1.2.3) αν θεωρήσουμε αντί του  $A$ , του πίνακα  $-A^{-1}B$ :

**Λήμμα 1.** Έστω  $A \in L(\mathbb{R}^n)$  αντιστρέψιμος και έστω  $B \in L(\mathbb{R}^n)$  τέτοιος ώστε  $\|A^{-1}\| \|B\| < 1$ . Τότε ο  $A+B$  είναι αντιστρέψιμος' επίσης έχουμε

$$(3) \quad \|(A+B)^{-1}\| \leq \|A^{-1}\| / (1 - \|A^{-1}\| \|B\|). \quad \odot$$

**ΘΕΩΡΗΜΑ 1** Υποθέτουμε ότι η απεικόνιση  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη σ' ένα σημείο  $x^* \in \text{Int}(D)$  όπου  $F(x^*) = 0$ . Για  $x \in S_0$ , όπου

## 2.3.2

$S_0 \subset \Omega$  είναι μία (ανοιχτή) περιοχή του  $x^*$ , θεωρούμε μία οικογένεια γραμμικών τελεστών  $A(x) \in L(\mathbb{R}^n)$ . Υποθέτουμε ότι η  $x \mapsto A(x)$  είναι συνεχής στο  $x^*$  και ότι ο  $A(x^*)$  είναι αντιστρέψιμος. Τότε υπάρχει κλειστή μπάλα  $\bar{S} = \bar{S}(x^*, \delta) \subset S_0$ ,  $\delta > 0$ , στην οποία η απεικόνιση  $G: \bar{S} \rightarrow \mathbb{R}^n$ , που ορίζεται για  $x \in \bar{S}$  από την (1) είναι καλά ορισμένη· επιπλέον η  $G$  είναι παραγωγίσιμη στο  $x^*$  και η παράγωγός της εκεί είναι

$$(4) \quad G'(x^*) = I - A(x^*)^{-1} F'(x^*)$$

Απόδειξη: Η  $G$  θα είναι καλά ορισμένη στην  $\bar{S}$  αν δείξουμε ότι ο  $A(x)$  είναι αντιστρέψιμος για  $x \in \bar{S}$ . Θέτουμε  $\beta = \|A(x^*)^{-1}\|$ . Έστω  $\epsilon > 0$  τέτοιο ώστε  $0 < \epsilon < 1/2\beta$ . Από την συνέχεια της  $x \mapsto A(x)$  στο  $x^*$  συμπεραίνουμε ότι υπάρχει  $\delta = \delta(x^*, \epsilon) > 0$  τέτοιο ώστε  $\bar{S} = \bar{S}(x^*, \delta) \subset S_0$  και

$$(5) \quad \|A(x) - A(x^*)\| \leq \epsilon, \quad \forall x \in \bar{S}.$$

Χρησιμοποιούμε τώρα το Λήμμα 1 με  $A = A(x^*)$ ,  $B = A(x) - A(x^*)$ : επειδή υπάρχει ο  $A(x^*)^{-1}$  και  $\|A(x^*)^{-1}\| \|A(x) - A(x^*)\| \leq \epsilon\beta < 1/2$  για  $x \in \bar{S}$  συμπεραίνουμε ότι ο  $A(x)$  είναι αντιστρέψιμος για  $x \in \bar{S}$  και ότι

$$(6) \quad \|A(x)^{-1}\| = \|(A(x^*) + (A(x) - A(x^*)))^{-1}\| \leq \\ \leq \|A(x^*)^{-1}\| / (1 - \|A(x^*)^{-1}\| \|A(x) - A(x^*)\|) \leq \beta / (1 - 1/2) = 2\beta, \quad \forall x \in \bar{S}$$

Όπως λοιπόν η  $G$  είναι καλά ορισμένη για  $x \in \bar{S}$ . Για να δείξουμε τώρα ότι η  $G$  είναι παραγωγίσιμη στο  $x^*$  και ότι η παράγωγός της εκεί δίνεται από την (4), παρατηρούμε πρώτα ότι επειδή το  $x^*$  είναι λύση του  $F(x) = 0$ , θα είναι σταθερό της  $G$ , δηλ. θα ισχύει

$$(7) \quad G(x^*) = x^*.$$

Χρησιμοποιώντας τώρα την παραγωγισιμότητα της  $F$  στο  $x^*$  και υποθέτοντας ότι η ακτίνα  $\delta$  της μπάλας  $\bar{S} = \bar{S}(x^*, \delta)$  είναι αρκούντως μικρή έχουμε ότι

$$(8) \quad \|F(x) - F(x^*) - F'(x^*)(x-x^*)\| \leq \epsilon \|x-x^*\|, \quad \forall x \in \bar{S}$$

Συνοψώς, για  $x \in \bar{S}$

$$\begin{aligned} & \|G(x) - G(x^*) - (I - A(x^*)^{-1} F'(x^*)) (x-x^*)\| = (\text{από τις (1), (7)}) \\ & \|A(x^*)^{-1} F'(x^*)(x-x^*) - A(x)^{-1} F(x)\| = (\text{λόγω της } F(x^*)=0) \\ & \| -A(x)^{-1} (F(x) - F(x^*) - F'(x^*)(x-x^*)) - A(x)^{-1} (A(x^*) - A(x)) A(x^*)^{-1} \\ & F'(x^*) (x-x^*) \| \leq \|A(x)^{-1} [F(x) - F(x^*) - F'(x^*)(x-x^*)]\| \\ & + \|A(x)^{-1} (A(x^*) - A(x)) A(x^*)^{-1} F'(x^*)(x-x^*)\| \leq (\text{λόγω των (6), (8),} \\ & (5) \text{ και του οριεμού του } \beta) \leq 2\beta \|x-x^*\| + 2\beta^2 \epsilon \|F'(x^*)\| \|x-x^*\| \\ & \equiv \epsilon \gamma \|x-x^*\|, \text{ όπου } \gamma = 2\beta + 2\beta^2 \|F'(x^*)\|. \text{ Συμπεραίνουμε ότι η } G \text{ είναι} \\ & \text{παραγωγίσιμη (Frechet) στο } x^* \text{ και ότι όντως η παράγωγός της } G'(x^*) \\ & \text{δίνεται από την (4). } \odot \end{aligned}$$

Άμεως προκύπτει τώρα ένα αποτέλεσμα τοπικής εύκλισης για την επαναληπτική μέθοδο (2):

**Πόρισμα 1** Έστω ότι ισχύουν οι υποθέσεις του θεωρήματος 1. Επιπλέον υποθέτουμε ότι

$$(9) \quad \rho(G'(x^*)) \equiv \rho(I - A(x^*)^{-1} F'(x^*)) = \epsilon < 1.$$

Τότε το  $x^*$  είναι σημείο έλξης της (2).

Απόδειξη: Εφαρμόζουμε το θεώρημα 2.2.1 του Ostrowski, το θεώρημα 1 και την υπόθεσή μας (9).  $\odot$

θεωρούμε τώρα την μέθοδο του Νεύτωνα για την, οποία  $A(x) = F'(x)$ . Τότε έχουμε το εξής θεώρημα τοπικής εύκλισης:

**ΘΕΩΡΗΜΑ 2** Υποθέτουμε ότι η  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη σε μία (ανοιχτή) περιοχή  $S_0 \subset D$  ενός σημείου  $x^*$  όπου  $F(x^*) = 0$ . Υποθέτουμε επίσης ότι η  $F'$  είναι συνεχής στο  $x^*$  και ότι ο  $F'(x^*)$  είναι αντιστέψιμος. Τότε το  $x^*$  είναι σημείο έλξης της μεθόδου του Νεύτωνα:

$$(10) \quad x^{k+1} = x^k - F'(x^k)^{-1} F(x^k), \quad k \geq 0$$

Απόδειξη Χρησιμοποιώντας το θεώρημα 1 με  $A(x) = F'(x)$  για  $x \in S_0$  (ελέγχουμε εύκολα ότι ισχύουν όλες οι υποθέσεις του), συμπεραίνουμε ότι η απεικόνιση επαναλήψεως  $G(x) = x - F'(x)^{-1} F(x)$  της μεθόδου του Νεύτωνα είναι καλά ορισμένη σε μία μπάλα  $\bar{S} = \bar{S}(x^*, \delta) \subset S_0$ ,  $\delta > 0$ .

Επιπλέον τώρα  $G'(x^*) = I - F'(x^*)^{-1} F'(x^*) = 0$  συνεπώς  $\rho(G'(x^*)) = \rho = 0$  και το πρόβλημα 1 δίνει ότι το  $x^*$  είναι σημείο έλξης της (10), δηλ. ότι αν  $x^0 \in S$  τότε  $x^k \in S$  και  $x^k \rightarrow x^*$ ,  $k \rightarrow \infty$ . @

Μετά από αυτό το σημαντικό αποτέλεσμα θα διερευνήσουμε την ταχύτητα σύγκλισης της ακολουθίας  $x^k$  της μεθόδου του Νεύτωνα προς το σημείο έλξης  $x^*$ . Συγκεκριμένα κάτω από ορισμένες ευνοϊκές μπορούμε να δείξουμε ότι η σύγκλιση είναι "τετραγωνική".

**ΠΡΟΤΑΣΗ 1.** Έστω ότι ισχύουν όλες οι υποθέσεις του θεωρήματος 2. Τότε για το σημείο έλξης  $x^*$  ισχύει ότι

$$(11) \quad \lim_{k \rightarrow \infty} \|x^{k+1} - x^*\| / \|x^k - x^*\| = 0$$

Αν επιπλέον για κάποια σταθερά  $c_1$  έχουμε

$$(12) \quad \|F'(x) - F'(x^*)\| \leq c_1 \|x - x^*\|$$

για κάθε  $x$  σε μία (ανοιχτή) περιοχή του  $x^*$ , τότε υπάρχει φυσικός  $k_0$  και σταθερά  $c_2$  τέτοια ώστε

$$(13) \quad \|x^{k+1} - x^*\| \leq c_2 \|x^k - x^*\|^2 \quad \text{για } k \geq k_0,$$

δηλ. όπως λέμε, η σύγκλιση είναι "τετραγωνική" για  $k$  αρκετά μεγάλο.

Απόδειξη Για την μέθοδο του Νεύτωνα, θέτουμε  $G(x) = x - F'(x)^{-1}F(x)$ , έχουμε από το Θεώρημα 2 ότι η  $G$  είναι καλά ορισμένη σε μία μπάλα  $\bar{S} = \bar{S}(x^*, \delta)$ ,  $\delta > 0$  με κέντρο  $x^*$ , όπου  $G(x^*) = x^*$  και όπου η  $G$  είναι παραγωγίσιμη με  $G'(x^*) = 0$ . Χρησιμοποιώντας αυτά τα δεδομένα έχουμε για  $x^0 \in \bar{S}$  ότι  $x^{k+1} - x^* = G(x^k) - G(x^*) = G(x^k) - G(x^*) - G'(x^*)(x^k - x^*)$ . Άρα  $\lim_{k \rightarrow \infty} \|x^{k+1} - x^*\| / \|x^k - x^*\| = \lim_{k \rightarrow \infty} \|G(x^k) - G(x^*) - G'(x^*)(x^k - x^*)\| / \|x^k - x^*\| = 0$  από την παραγωγισιμότητα της  $G$  στο  $x^*$ . Συνεπώς αποδείχθη η (11). Έστω τώρα ότι για  $k \geq k_0$  οι όροι  $\{x^k\}$  περιέχονται στην περιοχή του  $x^*$  όπου ισχύει η (12). Για τέτοιο  $k$  θεωρείστε ως κυρτό εύνολο  $D_0$  το ευθύγραμμο τμήμα  $[x^k, x^*]$  και υποθέστε χωρίς περιορισμό της γενικότητας ότι η (12) - πιθανώς με διαφορετική σταθερά  $c_1$  - ισχύει στο  $D_0$ . Εφαρμόζοντας το αποτέλεσμα της πρότασης 2.1.4 και την (12) έχουμε τότε ότι για κάποια σταθερά  $c$  ανεξάρτητη του  $k$  ισχύει

$$(14) \quad \|F(x^k) - F(x^*) - F'(x^*)(x^k - x^*)\| \leq c \|x^k - x^*\|^2.$$

Άρα

$$\begin{aligned} \|x^{k+1} - x^*\| &= \|x^k - F'(x^k)^{-1}F(x^k) - x^*\| \leq \|F(x^*) - F(x^k) + F'(x^k)^{-1}F(x^k) - F'(x^k)^{-1}F(x^*)\| \\ &\leq \|F'(x^k)^{-1}\| \|F(x^k) - F(x^*) - F'(x^k)(x^k - x^*)\| \\ &\quad + \|F'(x^k)^{-1}\| \|F'(x^k) - F'(x^*)\| \|x^k - x^*\| \\ &\leq (\text{από τις (14), (12)}) \|F'(x^k)^{-1}\| c \|x^k - x^*\|^2 \\ &\quad + \|F'(x^k)^{-1}\| c_1 \|x^k - x^*\| = \|F'(x^k)^{-1}\| (c + c_1) \|x^k - x^*\|^2. \end{aligned}$$

Όπως στην (6) της απόδειξης του θεωρήματος 1 έχουμε τώρα, επειδή ο  $A = F'$  ικανοποιεί όλες τις συνθήκες του θεωρήματος 1, ότι  $\|F'(x^k)^{-1}\| \leq 2 \|F'(x^*)^{-1}\| = \text{εταθ}$ . Συμπεραίνουμε λοιπόν ότι  $\|x^{k+1} - x^*\| \leq c_2 \|x^k - x^*\|^2$  για  $k \geq k_0$  και για σταθερά  $c_2$  ανεξάρτητη του  $k$ . @

Η εφαρμογή στην πράξη της μεθόδου του Νεύτωνα (10) σε ένα εύνολο του  $\mathbb{R}^n$  όπου η  $F'(x^k)$  είναι αντιστρέψιμη γίνεται ως εξής:

Γιά δεδομένο  $x^k$ , υπολογίζουμε το διάνυσμα  $F(x^k)$  και την τρέχουσα τιμή του Ιακωβιανού πίνακα  $J(x^k)$  που παριετάνει την  $F'(x^k)$ , λύνουμε το σύστημα

$$(15) J(x^k) y^k = F(x^k)$$

και μετά υπολογίζουμε το  $x^{k+1}$  από την σχέση  $x^{k+1} = x^k - y^k$ .

### Παρατηρήσεις

1. Είδαμε ότι εάν συνέπεια του γεγονότος ότι  $G'(x^*)=0$ , (βλ. απόδειξη της Πρότασης 1) ισχύει η (11) για την μέθοδο του Νεύτωνα. Στην περίπτωση μιάς γενικής επαναληπτικής μεθόδου είδαμε (αν η  $G$  είναι ευστολή) ότι έχουμε απλώς  $\|x^{k+1}-x^*\|/\|x^k-x^*\| \leq a < 1$ ,  $k \geq 0$  (βλ. απόδειξη θεωρήματος 2.2.2). Αν ισχύει η (12) είδαμε επίσης ότι για  $k$  αρκετά μεγάλο η ακολουθία της μεθόδου του Νεύτωνα ικανοποιεί την  $\|x^{k+1}-x^*\| \leq c\|x^k-x^*\|^2$  ("τετραγωνική" σύγκλιση) ενώ για την γενική επαναληπτική μέθοδο έχουμε "γραμμική" σύγκλιση:  $\|x^{k+1}-x^*\| \leq a\|x^k-x^*\|$ . (Γενικά λέμε ότι μία μέθοδος έχει τάξη σύγκλισης  $p$  - σε κάποια ρίζα  $x^*$  - αν υπάρχει σταθερά  $c$  ανεξάρτητη του  $k$  και ακέραιος  $k_0 \geq 0$  τέτοιος ώστε για  $k \geq k_0$ ,  $\|x^{k+1}-x^*\| \leq c\|x^k-x^*\|^p$ ). Είναι φανερό ότι όσο μεγαλύτερη είναι η τάξη σύγκλισης τόσο πιο γρήγορη είναι η σύγκλιση της  $x^k$  στο  $x^*$  (για  $x^k$  κοντά στο  $x^*$ ). Ο προεδιορισμός καλών μέτρων σύγκλισης της ταχύτητας σύγκλισης διαφόρων μεθόδων είναι σημαντικό ζήτημα: βλέπε την ανάλυση και τους ορισμούς του κεφ. 9 του [2.3].

2. Είδαμε ότι η εφαρμογή της μεθόδου του Νεύτωνα στην πράξη απαιτεί για κάθε  $k$  τον υπολογισμό ενός  $n \times n$  πίνακα  $J(x^k)$ , ενός  $n \times 1$  διανύσματος  $F(x^k)$  και την επίλυση του  $n \times n$  γραμμικού συστήματος (15) με  $n^3/3 + O(n^2)$  γενικά πράξεις. Σε πολλές εφαρμογές (όπου μάλιστα μπορεί να απαιτείται η επίλυση πολλών μη γραμμικών συστημάτων) το κόστος του υπολογισμού των πινάκων  $J(x^k)$  και της λύσης των γραμμικών συστημάτων μπορεί να είναι απαγορευτικό. Καταφεύγουμε λοιπόν ευχιά σε απλοποιήσεις της μεθόδου του Νεύτωνα, δηλ. σε μεθόδους της

μορφής (2), όπου ο γραμμικός τελεστής  $A(x^k)$  αποτελεί γενικά προέκταση της  $F'(x^k)$ , οι οποίες έχουν μόνιμο μικρότερο κόστος ανά βήμα, ευκλίνουσι όμως πιά όχι "τετραγωνικά" στη  $x^*$ . Πολλές φορές το ευνολικό κόστος τέτοιων μεθόδων είναι μικρότερο από το ευνολικό κόστος της μεθόδου του Νεύτωνα (ιδίως αν δεν απαιτείται μία πολύ ακριβής λύση).

Μία προφανής τέτοια απλούστευση είναι η λεγόμενη "μέθοδος της χορδής":

$$(16) \quad x^{k+1} = x^k - F'(x^0)^{-1} F(x^k), \quad k \geq 0, \quad x^0 \text{ δεδομένο,}$$

η οποία σε κάθε βήμα απαιτεί την λύση γραμμικού συστήματος με ε σταθερό πίνακα  $J(x^0)$ , ο οποίος ευνεπώς μπορεί να υπολογισθεί και να αναλυθεί σε μορφή LU πριν αρχίσει η επανάληψη (16). Μπορεί να αποδειχθεί ότι γενικά αυτή η μέθοδος ευκλίνει γραμμικά.

Σε μία άλλη στρατηγική απλούστευσης, παρατηρούμε ότι ο τακτιανός πίνακας  $J(x)$  μπορεί πάντα να γραφεί ως άθροισμα

$$(17) \quad J(x) = D(x) + L(x) + U(x)$$

ενός διαγωνίου πίνακα  $D(x)$  με  $D_{ii} = J_{ii}$ , ενός αυστηρά κάτω τριγωνικού πίνακα  $L(x)$  με  $L_{ij} = J_{ij}$ , αν  $i > j$ ,  $L_{ij} = 0$  αν  $i \leq j$  και ενός αυστηρά άνω τριγωνικού πίνακα  $U(x)$  με  $U_{ij} = J_{ij}$ , αν  $i < j$ ,  $U_{ij} = 0$  αν  $i \geq j$ , σε αναλογία με ότι γίνεται για τις κλασικές επαναληπτικές μεθόδους, βλ. [5.4, Κεφ. 11-13]. Θεωρούμε τότε, (αν  $J_{ii} \neq 0$ ), για  $A(x)$  π.χ. τον πίνακα  $D(x) + L(x)$  οπότε προκύπτει η επαναληπτική μέθοδος

$$(18) \quad x^{k+1} = x^k - (D(x^k) + L(x^k))^{-1} F(x^k), \quad k \geq 0, \quad x^0 \text{ δεδομένο,}$$

που απαιτεί την λύση ενός τριγωνικού γραμμικού συστήματος για κάθε  $k$ . (Στην περίπτωση ενός γραμμικού συστήματος η μέθοδος αυτή συμπίπτει με την γνωστή μέθοδο των Gauss-Seidel· ανάλογα λοιπόν η (18) ονομάζεται μέθοδος Newton/Gauss-Seidel). Μία πιο ακραία επιλογή είναι  $A(x) = D(x)$  οπότε προκύπτει η λεγόμενη μέθοδος Newton/Jacobi:



$$(19) \quad x^{k+1} = x^k - D(x^k)^{-1} F(x^k), \quad k \geq 0, \quad x^0 \text{ δεδομένο.}$$

Γιά τις μεθόδους αυτές μπορούμε να αποδείξουμε τοπικά θεωρήματα εύγκλισης (βλ. Ασκήσεις 2,3 παρακάτω). Βέβαια η εύγκλισή τους είναι γενικά γραμμική.

3. Επιπλέον ετών πράξη πολλές φορές δεν είναι γνωστές (ή είναι δύσκολο να υπολογισθούν αναλυτικά) οι μερικές παράγωγοι  $\partial_j f_i$  των  $f_i$  αλλά μόνο οι  $f_i$ . Τότε μπορούμε να καταφύγουμε σε προερχίσεις του  $J_{ij}$  όπου αντί των  $\partial_j f_i(x^k)$  χρησιμοποιούμε πεπερασμένες διαφορές  $(f_i(x^k + h e^j) - f_i(x^k))/h$ . Αν το  $h$  είναι αρκετά μικρό (και μάλιστα αν πάρουμε ακολουθία όλο και μικρότερων  $h_k$ ) μπορούμε να δείξουμε ότι η προέχιστική αυτή τεχνική διατηρεί πολλά από τα χαρακτηριστικά της μεθόδου του Νεύτωνα. Πάλιτα μπορεί να αποδειχθεί (βλ. π.χ. [2.1, κελ. 94]) ότι αν η  $F$  και η λύση  $x^*$  ικανοποιούν τις συνθήκες του θεωρήματος 2 και την (12), τότε υπάρχει  $\epsilon > 0$  και περιοχή  $S$  του  $x^*$  τέτοια ώστε αν διαλέξουμε ακολουθία πραγματικών  $\{h_k\}$  με  $0 < |h_k| \leq \epsilon$ ,  $k=0,1,2,\dots$ , τότε η ακολουθία  $\{x^k\}$  που παράγει η μέθοδος (1), όπου ο  $A(x^k)$  παριστάεται από τον πίνακα

$$A_{ij}(x^k) = (f_i(x^k + h_k e^j) - f_i(x^k))/h_k, \quad 1 \leq i, j \leq n, \quad k \geq 0,$$

είναι καλά ορισμένη για  $x^0 \in S$  και έχει την λύση  $x^*$  ως σημείο έλξης. Επιπλέον η εύγκλισή της στο  $x^*$  είναι γενικά γραμμική. Αν διαλέξουμε ακολουθία  $\{h_k\}$  τέτοια ώστε  $\lim_{k \rightarrow \infty} h_k = 0$  τότε η εύγκλιση είναι "υπεργραμμική", δηλ. η μέθοδος έχει τάξη εύγκλισης  $1 < p < 2$ . Πάλιτα, αν η ακολουθία  $\{h_k\}$  τείνει αρκετά γρήγορα στο 0, έτσι ώστε π.χ.  $|h_k| \leq c_1 \|x_k - x^*\|$  ή  $|h_k| \leq c_2 \|F(x^k)\|$  για σταθερές  $c_1, c_2$ , τότε η εύγκλιση είναι τετραγωνική.

Σε μία διάσταση είναι γνωστή μας βέβαια και η λεγόμενη "μέθοδος της τέμνουσας" που αντικαθιστά το πηλίκο  $f(x^k)/f'(x^k)$  της μεθόδου του Νεύτωνα με το πηλίκο διαφορών  $f(x^k)(x^k - x^{k-1})/(f(x^k) - f(x^{k-1}))$ .

Η μέθοδος αυτή έχει υπερχρυσή σύγκλιση, βλ. [5.2]. Γενικεύσεις της μεθόδου αυτής σε ευετήματα καθώς και παρεμφερείς προεχίσεις  $H(x)$  της  $F'(x)$  οδηγούν σε μία μεγάλη κλάση μεθόδων, τις λεγόμενες μεθόδους του "τύπου του Νεύτωνα" (Quasi-Newton methods) που αποτελούν τα τελευταία χρόνια αντικείμενο ιδιαίτερου ενδιαφέροντος· βλ. το βιβλίο [2.1].

### Άσκησης 2.3

1. Θεωρείστε την μέθοδο του Νεύτωνα στον  $\mathbb{R}^1$  για τις συναρτήσεις  $f(x)=x^2$  και  $f(x)=x+x^{1+\alpha}$  όπου  $0 < \alpha < 1$ . Δείξτε απ' ευθείας ότι και στις δύο περιπτώσεις το σημείο  $x^*=0$  είναι σημείο έλξης της ακολουθίας της μεθόδου του Νεύτωνα αλλά ότι η σύγκλιση δεν είναι τετραγωνική.

2. Διατυπώστε και αποδείξτε χρησιμοποιώντας το θεώρημα 1 και το Πρόλημα 1 (όπως κάναμε δηλ. στο θεώρημα 2 για την μέθοδο του Νεύτωνα) θεώρημα τοπικής σύγκλισης για την μέθοδο της χορδής (16).

3. Ομοίως για την μέθοδο Newton/Gauss-Seidel (18) και για την μέθοδο Newton/Jacobi (19).

4. (Η άσκηση αυτή αποτελεί συνέχεια των ασκήσεων 2.1.8, 2.2.7).

(α) Υποθέστε ότι  $g'(s) \geq 0 \quad \forall s \in \mathbb{R}^1$ . Δείξτε ότι ο Ιακωβιανός πίνακας που περιέχει την παράγωγο  $F'(x)$  που δίνεται από την (2.1.20) είναι αντιστρέψιμος.

(β) Αποδείξτε ότι η μέθοδος του Νεύτωνα για το μη γραμμικό εύστημα (2.1.19) συγκλίνει τοπικά. Δείξτε δηλ. ότι μία λύση  $U^*$  του (2.1.19) είναι σημείο έλξης της ακολουθίας της μεθόδου του Νεύτωνα. Υποθέστε ότι η  $g$  είναι συνεχώς παραγωγίσιμη και ότι  $g'(s) \geq 0, \quad \forall s \in \mathbb{R}^1$ .

(γ) Θεωρείστε την μέθοδο Newton/Gauss-Seidel για την επίλυση του (2.1.19). Με τις υποθέσεις του (β) δείξτε ότι μία λύση  $U^*$  του (2.1.19) είναι σημείο έλξης για την μέθοδο αυτή.

## 2.4 ΤΟ ΘΕΩΡΗΜΑ ΤΟΥ ΚΑΝΤΟΡΟΥΙΧ ΓΙΑ ΤΗΝ ΣΥΓΚΛΙΣΗ ΤΗΣ ΜΕΘΟΔΟΥ ΤΟΥ ΝΕΥΤΩΝΑ

Σ' αυτήν την παράγραφο θα μελετήσουμε την απόδειξη του Καντορουίχ (1948) για την σύγκλιση της μεθόδου του Νεύτωνα. Το αποτέλεσμα αυτό είναι του τύπου "περιορισμένης σύγκλισης", όπως π.χ. ήταν και το θεώρημα 2.2 της εισαγωγής.

**ΘΕΩΡΗΜΑ 1** (Καντορουίχ). Έστω ότι η απεικόνιση  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη ε' ένα κυρτό εύνολο  $D_0 \subset D$  και έστω ότι για κάποια σταθερά  $\gamma$  ισχύει

$$(1) \|F'(x) - F'(y)\| \leq \gamma \|x - y\|, \quad x, y \in D_0.$$

Έστω ότι  $F'(x^0)$  είναι αντιστρέψιμος ε' ένα σημείο  $x^0 \in D_0$ . Έστω επίσης ότι για κάποιες σταθερές  $\beta, \zeta$  έχουμε εκεί ότι

$$(2) \|F'(x^0)^{-1}\| \leq \beta,$$

$$(3) \|F'(x^0)^{-1} F(x^0)\| \leq \zeta,$$

όπου

$$(4) \alpha = \beta \gamma \zeta \leq 1/2.$$

θέτουμε

$$(5) t^* = [1 - (1 - 2\alpha)^{1/2}] / \beta \gamma, \quad t^{**} = [1 + (1 - 2\alpha)^{1/2}] / \beta \gamma$$

και υποθέτουμε ότι  $\bar{S}(x^0, t^*) \subset D_0$ . Τότε οι όροι της ακολουθίας της μεθόδου του Νεύτωνα

$$(6) x^{k+1} = x^k - F'(x^k)^{-1} F(x^k), \quad k=0, 1, 2, \dots$$

είναι καλά ορισμένοι, παραμένουν μέσα στην  $\bar{S}(x^0, t^*)$  και συγκλίνουν

για  $k \rightarrow \infty$  σε μία λύση  $x^*$  του  $F(x)=0$  (στην μοναδική λύση του  $F(x)=0$  στο εύνολο  $\bar{S}(x^0, t^*)$ ). Επιπλέον έχουμε την εξής εκτίμηση του σφάλματος:

$$(7) \|x^k - x^*\| \leq (2\alpha)^2 / 2^k \beta \gamma, \quad k=0, 1, 2, \dots$$

Απόδειξη θα υποθέσουμε ότι η (4) ισχύει ως αμεταβλητή, δηλ. ότι  $\alpha < 1/2$ . Η περίπτωση  $\alpha = 1/2$  είναι γενικά απλούστερη αλλά απαιτεί ορισμένες απλές αλλαγές στις αποδείξεις (βλ. 1)

Κατ' αρχήν δείχνουμε ότι η απεικόνιση επανάληψης  $G = I - (F')^{-1}F$  της μεθόδου του Νεύτωνα είναι καλά ορισμένη στην  $\bar{S}(x^0, t^*)$ . Ορίζουμε  $D_1 = S(x^0, (\beta\gamma)^{-1}) \cap D_0$ . Τότε, για  $x \in D_1$  έχουμε απ' την (1) ότι  $\|F'(x) - F'(x^0)\| \leq \gamma \|x - x^0\| < \beta^{-1}$ . Άρα η (2) δίνει ότι  $\|F'(x^0)^{-1}\| \|F'(x) - F'(x^0)\| < 1$  για  $x \in D_1$ . Συμπεραίνουμε από το Λήμμα 2.3.1 ότι οι γραμμικοί τελεστές  $F'(x)$ ,  $x \in D_1$  είναι αντιστρέψιμοι. Επιπλέον από τις (2.3.3), (2) και τα παραπάνω έχουμε για  $x \in D_1$ :

$$(8) \|F'(x)^{-1}\| \leq \|F'(x^0)^{-1}\| / (1 - \|F'(x^0)^{-1}\| \|F'(x) - F'(x^0)\|) \\ \leq \beta / (1 - \beta\gamma \|x - x^0\|)$$

Τώρα, επειδή  $\alpha < 1/2$ , η (5) δίνει ότι  $t^* < (\beta\gamma)^{-1}$ , δηλ. ότι  $\bar{S}(x^0, t^*) \subset D_1$ , έτσι ώστε η  $G$  που δίνεται από

$$(9) G(x) = x - F'(x)^{-1}F(x)$$

είναι καλά ορισμένη στην  $\bar{S}(x^0, t^*)$ .

Σαν ένα προκαταρκτικό βήμα θα κάνουμε τώρα μία εκτίμηση για  $x$  και  $G(x) \in S(x^0, t^*)$  της ποσότητας  $\|G^2(x) - G(x)\|$ , όπου συμβολίζουμε  $G^2(x) = G(G(x))$ . Η (9) δίνει

## 2.4.3

$$\begin{aligned}
 (10) \quad \|G^2(x) - G(x)\| &= \|F'(G(x))^{-1} F(G(x))\| = \\
 &= \|F'(G(x))^{-1} [F(G(x)) - F(x) - F'(x)(G(x) - x)]\| \\
 &\leq \|F'(G(x))^{-1}\| \|F(G(x)) - F(x) - F'(x)(G(x) - x)\| \leq (\text{από την (8)}) \\
 &\leq \beta \|F(G(x)) - F(x) - F'(x)(G(x) - x)\| / (1 - \beta \gamma \|G(x) - x^0\|).
 \end{aligned}$$

Τώρα η (1) και η (2.1.10) δίνουν για  $x, G(x) \in S(x^0, t^*)$

$$(11) \quad \|F(G(x)) - F(x) - F'(x)(G(x) - x)\| \leq \gamma \|G(x) - x\|^2 / 2.$$

Συνεπώς από την (10) και (11) παίρνουμε για  $x, G(x) \in S(x^0, t^*)$  ότι

$$(12) \quad \|G^2(x) - G(x)\| \leq \beta \gamma \|G(x) - x\|^2 / 2(1 - \beta \gamma \|G(x) - x^0\|).$$

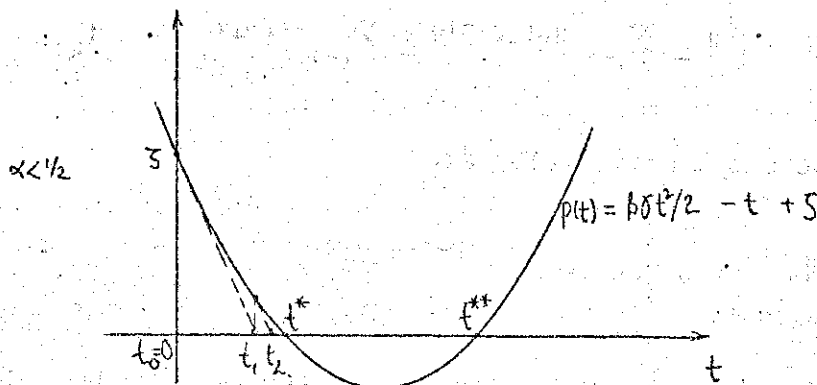
Θα επιτρέψουμε σε λίγο στην εκτίμηση (12). Προς το παρόν ορίζουμε την ακολουθία των πραγματικών αριθμών  $\{t_k\}$ ,  $k \geq 0$  από τις εκθέσεις

$$(13) \quad t_{k+1} = t_k - ((\beta \gamma t_k^2 / 2) - t_k + \zeta) / (\beta \gamma t_k - 1), \quad k \geq 0, \quad t_0 = 0.$$

Είναι εύκολο να διαπιστώσουμε ότι η ακολουθία  $\{t_k\}$ ,  $k \geq 0$  είναι καλά ορισμένη, μονοτονικά αύξουσα και συγκλίνει στον αριθμό  $t^*$ . Πράγματι, θεωρούμε το διωνυμο

$$(14) \quad p(t) = (\beta \gamma t^2 / 2) - t + \zeta.$$

Επειδή  $a = \beta \gamma \zeta < 1/2$  το  $p(t)$  έχει ρίζες τους αριθμούς  $t^*$  και  $t^{**}$  που δίνονται από την (5). Επιπλέον η (13)



δεν είναι τίποτε άλλο από την ακολουθία που παράγει η μέθοδος του Νεύτωνα για την προσέγγιση της ρίζας  $t^*$  της εξίσωσης  $p(t)=0$  με αρχική τιμή  $t_0=0$ . Η κυρτότητα και η ομαλότητα του  $p(t)$  δίνουν τώρα ότι  $0 < t_i < t_{i+1}$  για  $i \geq 1$  και ότι  $t_k \rightarrow t^*$ ,  $k \rightarrow \infty$ .

Ο σκοπός μας είναι να δείξουμε τώρα, χρησιμοποιώντας την (12) και τις ιδιότητες της ακολουθίας  $\{t_k\}$ ,  $k \geq 0$ , ότι οι όροι  $\{x^k\}$ ,  $k \geq 0$  της ακολουθίας (6) της μεθόδου του Νεύτωνα είναι καλά ορισμένοι, ανήκουν στην  $\bar{S}(x^0, t^*)$  και ικανοποιούν τις ανισότητες

$$(15) \quad \|x^k - x^{k-1}\| \leq t_k - t_{k-1}, \quad k \geq 1.$$

Είδαμε ήδη, επειδή η  $G(x)$  είναι καλά ορισμένη στην  $\bar{S}(x^0, t^*)$ , ότι το  $x^1 = G(x^0)$  υπάρχει. Επιπλέον, από τις (6), (3), (13)

$$\|x^1 - x^0\| = \|F'(x^0)^{-1}F(x^0)\| \leq \xi = t_1 = t_1 - t_0.$$

Συνεπώς ισχύει η (15) για  $k=1$ . Επιπλέον, επειδή  $\xi = t_1 < t^*$  έχουμε ότι  $x^1 \in S(x^0, t^*)$ . Υποθέτουμε τώρα ότι για κάποιο  $k \geq 1$  τα  $x^0, \dots, x^k$  υπάρχουν στην  $\bar{S}(x^0, t^*)$  και ικανοποιούν τις σχέσεις

$$(16) \quad \|x^i - x^{i-1}\| \leq t_i - t_{i-1}, \quad 1 \leq i \leq k.$$

Επειδή η  $G$  είναι καλά ορισμένη στην  $\bar{S}(x^0, t^*)$ , το  $x^{k+1} = G(x^k)$  υπάρχει. Επιπλέον επειδή η υπόθεση επαγωγής (16) συνεπάγεται ότι

$$(16') \quad \|x^k - x^0\| \leq \sum_{i=1}^k \|x^i - x^{i-1}\| \leq \sum_{i=1}^k (t_i - t_{i-1}) = t_k,$$

εμπεραίνουμε από τις (12), (16) ότι

$$(17) \quad \begin{aligned} \|x^{k+1} - x^k\| &= \|G(x^k) - G(x^{k-1})\| = \|G^2(x^{k-1}) - G(x^{k-1})\| \\ &\leq \beta\gamma \|x^k - x^{k-1}\|^2 / 2(1 - \beta\gamma \|x^k - x^0\|) \leq \beta\gamma (t_k - t_{k-1})^2 / 2(1 - \beta\gamma t_k) \\ &= (\text{κάντε τις πράξεις χρησιμοποιώντας την (13)}) = t_{k+1} - t_k. \end{aligned}$$

Συνεπώς ισχύει η (16) για  $i=k+1$ . Επιπλέον

$$\|x^{k+1}-x^0\| \leq \|x^k-x^0\| + \|x^{k+1}-x^k\| \leq t_k + t_{k+1}-t_k = t_{k+1} < t^*.$$

Συμπεραίνουμε ότι  $x^{k+1} \in \bar{S}(x^0, t^*)$ . Το επαγωγικό βήμα τελείωσε· έχουμε δείξει ότι όλοι οι όροι της ακολουθίας  $x^k$  υπάρχουν και παραμένουν στην  $\bar{S}(x^0, t^*)$ · επιπλέον δε ικανοποιούν για  $i \geq 1$  την (16).

Η τελευταία παρατήρηση μας δίνει βέβαια ότι για κάθε  $k \geq 0$  και  $m \geq 1$  έχουμε

$$(18) \quad \|x^{k+m}-x^k\| \leq \|x^{k+m}-x^{k+m-1}\| + \dots + \|x^{k+1}-x^k\| \leq t_{k+m}-t_k,$$

και επειδή  $t_k \rightarrow t^*$ ,  $k \rightarrow \infty$  συμπεραίνουμε ότι η  $\{x^k\}$ ,  $k \geq 0$ , είναι Cauchy στην  $\bar{S}(x^0, t^*)$ . Συνεπώς, υπάρχει  $x^* \in \bar{S}(x^0, t^*)$  τέτοιο ώστε

$$\lim_{k \rightarrow \infty} x_k = x^*.$$

Το  $x^*$  είναι άρα λύση του συστήματος  $F(x)=0$ : Πράγματι, η ύπαρξη της  $F'$  στο  $D_0$  συνεπάγεται την συνέχεια της  $F$  στο  $D_0$  (Πρόταση 2.1.1). Άρα επειδή

$$\begin{aligned} \|F(x^k)\| &= \|F'(x^k)(x^{k+1}-x^k)\| \leq \|F'(x^k)\| \|x^{k+1}-x^k\| \\ &\leq (\|F'(x^k)-F'(x^0)\| + \|F'(x^0)\|) \|x^{k+1}-x^k\| \leq (\text{από την (1)}) \\ &\leq (\gamma \|x^k-x^0\| + \|F'(x^0)\|) \|x^{k+1}-x^k\| \leq (\text{από την (16')}) \\ &\leq (\gamma t^* + \|F'(x^0)\|) \|x^{k+1}-x^k\| \rightarrow 0, \quad k \rightarrow \infty, \text{ συμπεραίνουμε ότι} \\ &F(x^*) = \lim_{k \rightarrow \infty} F(x^k) = 0. \end{aligned}$$

Για την μοναδικότητα της λύσης  $x^*$ , υποθέτουμε ότι υπάρχει και άλλη λύση  $y^*$ ,  $F(y^*)=0$ ,  $y^* \in \bar{S}(x^0, t^*)$ . Τότε ισχύει ότι

$$(19) \quad \|y^*-x^k\| \leq t^*-t_k, \quad k=0, 1, 2, \dots$$

Πράγματι η (19) ισχύει για  $k=0$  γιατί  $\|y^*-x^0\| \leq t^*=t^*-t_0$ . Υποθέτουμε ότι ισχύει για  $0 \leq k \leq i$ . Τότε, επειδή  $y^*=G(y^*)=G^2(y^*)$  έχουμε από τις (12), (16') ότι

$$\begin{aligned} \|y^*-x^{i+1}\| &= \|G(y^*)-G(x^i)\| = \|y^*-x^i + F'(x^i)^{-1} F(x^i)\| \\ &= \|F'(x^i)^{-1}(F(x^i)-F(y^*)-F'(x^i)(x^i-y^*))\| \leq \\ &\leq \beta\gamma \|y^*-x^i\|^2 / 2(1-\beta\gamma \|x^i-x^0\|) \leq \beta\gamma \|y^*-x^i\|^2 / 2(1-\beta\gamma t_i). \end{aligned}$$

'Αρα η (19) δίνει ότι

$$(19') \|y^*-x^{i+1}\| \leq \beta\gamma (t^*-t^i)^2 / 2(1-\beta\gamma t_i) = (\text{πράξεις!}) = t^*-t_{i+1}.$$

'Αρα ισχύει η (19) και για  $k=i+1$  και ευνεπώς για κάθε  $k$ .

Επειδή  $t_k \rightarrow t^*$ ,  $k \rightarrow \infty$  έχουμε ότι  $y^* = \lim_{k \rightarrow \infty} x^k = x^*$ , δηλ. ότι η λύση  $x^*$  είναι η μοναδική λύση του συστήματος  $F(x)=0$  στην μπάλα  $\bar{S}(x^0, t^*)$ .

Τέλος ερχόμαστε στην απόδειξη του φράγματος (7) για το εφάλμα της μεθόδου. Κατ' αρχήν δείχνουμε ότι

$$(20) t_{k+1} - t_k \leq \zeta 2^{-k}, \quad k=0, 1, 2, \dots$$

Πράγματι η (20) ισχύει για  $k=0$  επειδή  $t_1 - t_0 = t_1 = \zeta$ . Υποθέτοντας ότι ισχύει για  $0 \leq k \leq i$  έχουμε από την τελευταία ισότητα της (17) και την επαγωγική υπόθεση ότι

$$\begin{aligned} (21) \quad t_{i+2} - t_{i+1} &= \beta\gamma (t_{i+1} - t_i)^2 / 2(1-\beta\gamma t_{i+1}) \\ &\leq \beta\gamma \zeta^2 2^{-2i} / 2(1-\beta\gamma t_{i+1}). \end{aligned}$$

Εξ άλλου πάλι απ' την υπόθεση της επαγωγής έχουμε

$$t_{i+1} = \sum_{k=0}^i (t_{k+1} - t_k) \leq \zeta \sum_{k=0}^i 2^{-k} = 2\zeta(1-2^{-(i+1)}).$$

'Αρα στην (21), λόγω της (4)



$$(20') \quad 1 - \beta \gamma t_{i+1} = 1 - \alpha t_{i+1} / \zeta \geq 1 - 2\alpha(1 - 2^{-i+1}) \geq 2^{-i+1}.$$

Συνεπώς από τις (21), (4),  $t_{i+2} - t_{i+1} \leq \beta \gamma \zeta^2 2^{-i} = \zeta \alpha 2^{-i} \leq \zeta 2^{-i+1}$ .

δηλ. ισχύει η (20) για κάθε  $k \geq 0$ , πράγμα που συνεπάγεται την ισχύ της (20') για κάθε  $i \geq 0$ . Τώρα ισχυριζόμαστε ότι έχουμε επίσης

$$(22) \quad t^* - t_k \leq (\beta \gamma 2^k)^{-1} (2\alpha)^2, \quad k=0,1,2,\dots$$

Πράγματι, για  $k=0$  έχουμε ότι  $t^* - t_0 = t^* = (1 - (1 - 2\alpha)^{1/2}) / \beta \gamma \leq 2\alpha / \beta \gamma$  γιατί  $0 \leq \alpha < 1/2$ . Αν η (21) ισχύει για  $0 \leq k \leq i$ , έχουμε από την τελευταία ιδιότητα της (19'), την (20') και την επαγωγική υπόθεση ότι

$$\begin{aligned} t^* - t_{i+1} &= \beta \gamma (t^* - t_i)^2 / 2(1 - \beta \gamma t_i) \leq \beta \gamma (2\alpha)^2 \quad 2^{i-1} / \beta^2 \gamma^2 2^{2i} \\ &= (2\alpha)^2 / \beta \gamma 2^{i+1}. \end{aligned}$$

Συνεπώς ισχύει η (22) για κάθε  $k \geq 0$ . Η (18) για  $m \rightarrow \infty$  δίνει τώρα  $\|x^* - x^k\| \leq t^* - t^k$ . Συνεπώς από την (22) έπεται το ζητούμενο φράγμα (7).@

### Παρατηρήσεις

1. Το θεώρημα του Kantorovich (και η απόδειξή του ε' αυτήν την παράγραφο) ισχύει και για γενικότερες απεικονίσεις μεταξύ χώρων Banach. Ένα σημαντικό του πλεονέκτημα είναι ότι δεν υποθέτει εκ των προτέρων την ύπαρξη μίας λύσης  $x^*$  της εξίσωσης  $F(x)=0$ . Χρησιμεύει λοιπόν, και σαν εργαλείο για απόδειξη θεωρημάτων ύπαρξης και μοναδικότητας λύσεων ευστημάτων μη γραμμικών εξισώσεων αλλά και συναρτησιακών εξισώσεων όπως π.χ. ολοκληρωτικών και διαφορικών εξισώσεων. Ας σημειωθεί ότι μπορεί να αποδειχθεί (βλ. [2.3, Παρ. 12.5.4]) ότι η λύση  $x^*$  είναι η μοναδική λύση του ευστήματος  $F(x)=0$  όχι μόνο στην  $\bar{S}(x^0, t^*)$  αλλά και στο μεγαλύτερο, αν  $\alpha < 1/2$ , εύνολο  $\bar{S}(x^0, t^{**}) \cap D_0$ .

2. Για την λύση  $x^*$  όπως δείχνει η (19) ισχύει η εκτίμηση

$$\|x^* - x^k\| \leq t^* - t_k, \quad k=0,1,2,\dots$$

Συνεπώς για  $k=1$  έχουμε χρησιμοποιώντας την συνθήκη  $p(t^*)=0$ ,

$$\|x^* - x^1\| \leq t^* - t^1 = t^* - \zeta = \beta\gamma(t^*)^2/2.$$

Η (5) δίνει τώρα ότι  $t^* = [1 - (1 - 2\alpha)^{1/2}]/\beta\gamma \leq 2\alpha/\beta\gamma$ . Άρα έχουμε την

$$\|x^* - x^1\| \leq 2\alpha^2/\beta\gamma = 2\beta\gamma\zeta^2,$$

εχέση η οποία ισχύει για κάθε  $\zeta \geq \|F'(x^0)^{-1} F(x^0)\|$ , για το οποίο ικανοποιείται η (4). Παίρνοντας  $\zeta = \|F'(x^0)^{-1} F(x^0)\|$ , συμπεραίνουμε λοιπόν ότι οι υποθέσεις του θεωρήματος 1 δίνουν επίσης ότι

$$(23) \quad \|x^* - x^1\| \leq 2\beta\gamma \|x^1 - x^0\|^2.$$

Η ανισότητα (23) είναι πολύ χρήσιμη για την a posteriori εκτίμηση του εφάλματος της μεθόδου του Νεύτωνα: Υποθέστε ότι υπολογίζουμε τους όρους της ακολουθίας της μεθόδου  $\{x^k\}$ ,  $k=0,1,2,\dots$  μέχρις ότου ισχύει για πρώτη φορά το "κριτήριο τερματισμού"  $\|x^j - x^{j-1}\| \leq \epsilon$  όπου  $\epsilon > 0$  δεδομένο. Για να βρούμε τώρα μία εκτίμηση του εφάλματος της  $x^j$ , εφαρμόζουμε το θεώρημα 1 με  $x^{j-1}$  αντί του  $x^0$ , δηλ. με  $\beta \geq \|F'(x^{j-1})^{-1}\|$  (ποσότητα που μπορεί να εκτιμηθεί στην πράξη με την επίλυση δύο-τριών συστημάτων με πίνακα  $F'(x^{j-1})$  κατ' αναλογία με την εκτίμηση στην πράξη του δείκτη κατάταξη ενός πίνακα) και με  $\zeta = \|F'(x^{j-1})^{-1} F(x^{j-1})\| = \|x^j - x^{j-1}\| \leq \epsilon$ , δηλ. με  $\alpha = \beta\gamma\zeta \leq \beta\gamma\epsilon$ . Αν λοιπόν  $\beta\gamma\epsilon \leq 1/2$  και  $\bar{S}(x^0, t^*) \subset D_0$  - και τα δύο αυτά εξασφαλίζονται αν το  $\epsilon$  είναι αρκετά μικρό - η (23) ερμηνεύεται σαν

$$\|x^* - x^j\| \leq 2\beta\gamma\epsilon^2,$$

που είναι η ζητούμενη a posteriori εκτίμηση του εφάλματος  $\|x^* - x^j\|$ . ο υπολογισμός του φράγματος απαιτεί βέβαια και κάποια εκτίμηση του  $\gamma$ , δηλ. ενός φράγματος της "δεύτερης" παραγώγου του  $F$ .

3. Οι υποθέσεις του θεωρήματος του Καντοροίτς επισημαίνουν ένα σημαντικό πρόβλημα που εμφανίζεται και στην εφαρμογή της μεθόδου του Νεύτωνα στην πράξη: δηλ. ότι πρέπει να επιλέξουμε τον πρώτο όρο  $x^0$  της ακολουθίας αρκετά κοντά στην (άγνωστη) λύση  $x^*$  για να έχουμε εύγκλιση της μεθόδου. Ζητάμε να βρούμε λοιπόν τεχνικές για τον εντοπισμό καλών  $x^0$ , ή παραλλαγές της μεθόδου του Νεύτωνα που έχουν μεγαλύτερες "περιοχές εύγκλισης", δηλ. μεγαλύτερα εύνολα μέσα στα οποία οποιαδήποτε επιλογή του  $x^0$  θα οδηγήσει σε εύγκλιση της  $\{x^k\}$ .

Μία κατηγορία τεχνικών που χρησιμοποιούνται στην πράξη είναι οι λεγόμενες μέθοδοι "ευνέχισης" (continuation) που μπορούν να περιγραφούν ως εξής: θεωρούμε το σύστημα  $F(x)=0$  που πρέπει να λύσουμε σαν το τελευταίο σε μία συνεχή μονοπαραμετρική οικογένεια συστημάτων

$$H(x,t)=0,$$

όπου  $x \in \mathbb{R}^n$  και  $t \in [0,1]$  παράμετρος. Υποθέτουμε δηλ. ότι  $H(x,1)=F(x)$ , και ότι η απεικόνιση  $H(x,0)$  είναι τέτοια ώστε το σύστημα  $H(x,0)=0$  να λύνεται εύκολα. Π.χ. ευθείες επιλογές είναι, για δεδομένο  $a \in \mathbb{R}^n$ ,

$$H(x,t) = tF(x) + (1-t)(F(x) - F(a))$$

ή

$$H(x,t) = tF(x) + (1-t)(x-a).$$

Βρίσκουμε πρώτα την λύση  $\xi^{(0)}$  του συστήματος  $H(x,0)=0$ . Κατόπιν, για ένα διαμερισμό  $0=t_0 < t_1 < \dots < t_N=1$  του  $[0,1]$ , λύνουμε τα συστήματα  $H(x,t_i)=0$  υπολογίζοντας τις λύσεις τους  $\xi^{(i)}$  με την μέθοδο του Νεύτωνα και χρησιμοποιώντας ως αρχική τιμή για κάθε  $i$  π.χ. την λύση  $\xi^{(i-1)}$  του προηγούμενου συστήματος, βασιζόμενοι στο ότι τα προβλήματα  $H(x,t^{i-1})=0$  και  $H(x,t^i)=0$  είναι αρκετά "κοντά" μεταξύ τους έτσι ώστε η λύση  $\xi^{(i-1)}$  του πρώτου να αποτελεί μία καλή πρώτη προσέγγιση στη λύση του δεύτερου.

Πιο ευχάριστα όμως στην πράξη χρησιμοποιούνται κάποιες μορφές συνδυασμοί της μεθόδου του Νεύτωνα (ή των απλουστευμένων της ή

μεθόδων του "τύπου-του Νεύτωνα" βλ. παρατηρήσεις 2.3.2 και 2.3.3) και μιάς μεθόδου ελαχιστοποίησης ενός καταλλήλου ευαρτησιακού. Η γενική δομή μιάς τέτοιας μεθόδου είναι η εξής: δεδομένου του  $x^k$  υπολογίζουμε με την μέθοδο του Νεύτωνα (ή κάποια παραλλαγή της) την διαφορά  $y^k = x^{k+1} - x^k$ . Κατόπιν, είτε αποδεχόμαστε το  $x^{k+1} = x^k + y^k$  σαν επόμενο όρο της ακολουθίας, είτε το απορρίπτουμε με κριτήριο το αν ελαττώνει ή όχι κάποιο "φυσικό" μη αρνητικό ευαρτησιακό του προβλήματος το οποίο μηδενίζει η λύση, όπως, π.χ. την συνάρτηση  $f(x) = \|F(x)\|_2^2$ . Αν το  $x^{k+1}$  απορριφθεί, τότε υπολογίζουμε ένα νέο  $x^{k+1}$  είτε της μορφής  $x^k + \alpha y^k$  (δηλ. πάνω στην κατεύθυνση  $y^k$  από το  $x^k$ , μιά τέτοια κατεύθυνση είναι κατεύθυνση τοπικής ελαχιστοποίησης του  $f(x)$  στο  $x^k$ , βλ. Θεωρ. 4), είτε εκτελώντας ένα βήμα κάποιου αλγόριθμου ελαχιστοποίησης, π.χ. μιάς μεθόδου καθόδου μεγίστης κλίσεως ή μιάς μεθόδου συζυγών κλίσεων για ευαρτησιακά που προέρχονται από μη γραμμικά προβλήματα. Βλέπε το κεφ. 6 του βιβλίου [2.1] για μία καλή εισαγωγή στο θέμα.

4. Η γεωμετρική ερμηνεία της μεθόδου του Νεύτωνα στον  $\mathbb{R}^1$  μας πείθει αμέσως ότι αν η  $f: \mathbb{R}^1 \rightarrow \mathbb{R}^1$  είναι κυρτή (ή κοίλη), αυστηρά αύξουσα (ή φθίνουσα) και έχει ρίζα  $x^*$ , τότε η ακολουθία  $\{x^k\}$  της μεθόδου του Νεύτωνα συγκλίνει στην  $x^*$  για οποιοδήποτε  $x^0 \in \mathbb{R}^1$ . Κάτι ανάλογο συμβαίνει και στον  $\mathbb{R}^n$ . Λέμε ότι μία απεικόνιση  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι κυρτή στο κυρτό εύνολο  $D_0 \subset D$  αν  $F(\lambda x + (1-\lambda)y) \leq \lambda F(x) + (1-\lambda)F(y)$ , για κάθε  $x, y \in D_0$  και  $\lambda \in [0, 1]$ , όπου για  $x, y \in \mathbb{R}^n$   $x \leq y \Leftrightarrow x_i \leq y_i, \quad 1 \leq i \leq n$  (και ανάλογα για  $A, B \in L(\mathbb{R}^n)$ ,  $A \leq B \Leftrightarrow A_{ij} \leq B_{ij}, \quad 1 \leq i, j \leq n$ ). Ισχύει το εξής θεώρημα του Balmer: 'Εστω ότι η  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συνεχώς παραγωγίσιμη και κυρτή σ' όλο του  $\mathbb{R}^n$ . Υποθέτουμε επίσης ότι η  $F'(x)^{-1}$  υπάρχει και μάλιστα ότι  $F'(x)^{-1} \geq 0, \quad \forall x \in \mathbb{R}^n$ . 'Εστω ότι υπάρχει ρίζα  $x^*$  του συστήματος  $F(x) = 0$ . Τότε η  $x^*$  είναι μοναδική και η ακολουθία  $\{x^k\}$  της μεθόδου του Νεύτωνα συγκλίνει στην  $x^*$  για οποιοδήποτε  $x^0 \in \mathbb{R}^n$ . Επιπλέον  $\forall k \geq 1, \quad x^* \leq x^{k+1} \leq x^k$ . (βλ. Θεωρ. 5).

#### Θεωρήματα 2.4

1. Αποδείξτε το θεώρημα 1 του Kantorovich για  $\alpha = 1/2$ .

2. Αποδείξτε ότι αν  $t_0 < (\beta\gamma)^{-1}$ , τότε η ακολουθία  $\{t_k\}$  που ορίζεται από την σχέση (13) συγκλίνει, μονοτονικά αυξανόμενη (για  $k \geq 1$ ), στην  $t^*$ .

3. Υποθέστε ότι η  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι συνεχώς παραγωγίσιμη σε μία περιοχή μιάς λύσης  $x^*$  του  $F(x)=0$  και ότι υπάρχει η  $F'(x^*)^{-1}$ . Δείξτε ότι η απεικόνιση επανάληψης της μεθόδου του Νεύτωνα  $G(x) = x - F'(x)^{-1}F(x)$  είναι καλά ορισμένη και είναι ευστολή σε μία περιοχή του  $x^*$ .

4. Θεωρείστε το συναρτησιακό  $f(x) = \|F(x)\|_2^2$ , όπου  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  και το διάνυσμα  $y^k = x^{k+1} - x^k$ , όπου  $\{x^k\}$  η ακολουθία της μεθόδου του Νεύτωνα για το σύστημα  $F(x)=0$ , την οποία υποθέτουμε καλά ορισμένη. Δείξτε ότι το  $f(x)$  ελαττώνεται τοπικά κατά μήκος της ακτίνας με αρχή  $x^k$  και κατεύθυνση  $y^k$ .

5. (Η άσκηση αυτή σκοπό έχει να αποδείξει το θεώρημα του Valueen της Παρατήρησης 4).

(α) Έστω ότι η  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  είναι παραγωγίσιμη στο κυρτό σύνολο  $D_0 \subset D$ . Τότε η  $F$  είναι κυρτή στο  $D_0$  (βλ. Παρατήρηση 4) αν και μόνο αν

$$F(y) - F(x) \geq F'(x)(y-x), \quad \forall x, y \in D_0.$$

(β) Χρησιμοποιώντας το (α) δείξτε με επαγωγή ότι οι υποθέσεις του θεωρήματος του Valueen συνεπάγονται ότι  $x^* \leq x^k$  και  $F(x^k) \geq 0$  για  $k=1, 2, \dots$ . Συμπέρασμα:  $x^* \leq x^{k+1} \leq x^k$  για  $k \geq 1$ , από τον ορισμό του  $x^{k+1}$  και την υπόθεση ότι  $F'(x)^{-1} \geq 0$ .

(γ) Δείξτε τώρα ότι η ακολουθία  $\{x^k\}$  συγκλίνει όταν  $k \rightarrow \infty$  ε' ένα διάνυσμα  $y$  το οποίο είναι λύση του  $F(x)=0$ . (Υπόδειξη: θεωρείστε για  $1 \leq i \leq n$  την ακολουθία  $\{x^k\}$ ,  $k \geq 1$  και χρησιμοποιείστε το (β)).

(δ) Χρησιμοποιώντας το (α) δείξτε ότι η ρίζα  $x^*$  του  $F(x)=0$  είναι μοναδική στον  $\mathbb{R}^n$  και συνεπώς ότι  $x^k \rightarrow x^*$ ,  $k \rightarrow \infty$ .

(ε) Η υπόθεση της κυρτότητας της  $F$  είναι απαραίτητη στην υπόθεση του θεωρήματος του Βαλιέν: θεωρείτε τις απεικονίσεις  $f: \mathbb{R}^1 \rightarrow \mathbb{R}^1$ ,  $f(x) = \arctan x$  και  $f(x) = 2x + \sin x$ . Δείξτε ότι ικανοποιούν όλες τις άλλες συνθήκες του θεωρήματος με εξαίρεση την κυρτότητα της  $f$  και ότι η μέθοδος του Νεύτωνα δεν συγκλίνει για κάθε  $x^0 \in \mathbb{R}^1$  στις περιπτώσεις τους.

6. (Η άσκηση αυτή αποτελεί συνέχεια - και τέλος (!) - των ασκήσεων 2.1.8, 2.2.7 και 2.3.4).

(α) Διατυπώστε κατάλληλες υποθέσεις (η ισχύς των οποίων να είναι δυνατόν να ελέγχεται εύκολα, αλλά τέτοιες ώστε να μην οδηγούν σε τετριμμένα  $g$  και  $U^0$ !) πάνω στα  $g, U^0, h$  για να ισχύει ένα θεώρημα του Kantorovich για την ακολουθία  $\{U^k\}$ ,  $k \geq 0$  της μεθόδου του Νεύτωνα για το σύστημα (2.1.19). Διατυπώστε επίσης τις υποθέσεις και τα συμπεράσματα του θεωρήματός σας.

(β) Έστω ότι η  $g: \mathbb{R}^1 \rightarrow \mathbb{R}^1$  είναι κυρτή, (βλ. Παρατήρηση 4). Δείξτε ότι οι απεικονίσεις  $\Phi$  και  $F$ , που ορίζονται αντίστοιχα από τις (2.1.18) και (2.1.19) είναι τότε κυρτές στον  $\mathbb{R}^n$ . Υποθέστε επιπλέον ότι η  $g$  είναι συνεχώς παραγωγίσιμη με μη αρνητική παράγωγο στον  $\mathbb{R}^1$  και ότι το σύστημα (2.1.19) έχει λύση  $U^*$ . Δείξτε ότι η  $U^*$  είναι η μοναδική λύση του (2.1.19) στην οποία συγκλίνει η μέθοδος του Νεύτωνα για κάθε  $U^0 \in \mathbb{R}^n$ .

**ΕΘΝΙΚΟ ΚΑΙ ΚΑΠΟΔΙΣΤΡΙΑΚΟ  
ΠΑΝΕΠΙΣΤΗΜΙΟ ΑΘΗΝΩΝ**

**ΑΡΙΘΜΗΤΙΚΗ ΑΝΑΛΥΣΗ**

**(ΣΗΜΕΙΩΣΕΙΣ ΜΕΤΑΠΤΥΧΙΑΚΟΥ ΜΑΘΗΜΑΤΟΣ)  
Β' ΜΕΡΟΣ (ΚΕΦ. 3,4)**

**ΒΑΣΙΛΕΙΟΣ Α. ΔΟΥΓΑΛΗΣ**

**Μαθηματικό Τμήμα Πανεπιστημίου Αθηνών**

**ΑΘΗΝΑ 1996**

### 3. ΑΡΙΘΜΗΤΙΚΗ ΛΥΣΗ ΣΥΝΗΘΩΝ ΔΙΑΦΟΡΙΚΩΝ ΕΞΙΣΩΣΕΩΝ

#### 3.1 ΠΡΟΒΛΗΜΑ ΑΡΧΙΚΩΝ ΤΙΜΩΝ. Η ΜΕΘΟΔΟΣ ΤΟΥ EULER

Στο κεφάλαιο αυτό θα ασχοληθούμε με αριθμητικές μεθόδους για την προσεγγιστική λύση του προβλήματος αρχικών τιμών για ευετήματα πρώτης τάξης Συνήθων Διαφορικών Εξισώσεων. (Σ.Δ.Ε) Έστω  $-\infty < a < b < +\infty$ . Ζητάμε μία συνάρτηση  $y: [a, b] \rightarrow \mathbb{R}^m$ , παραγωγίσιμη στο  $[a, b]$ , που να ικανοποιεί το πρόβλημα αρχικών τιμών

$$(1) \quad \begin{cases} y'(t) = f(t, y(t)), & a \leq t \leq b, \\ y(a) = y_0, \end{cases}$$

όπου  $f$  δεδομένη συνάρτηση,  $f: [a, b] \times \mathbb{R}^m \rightarrow \mathbb{R}^m$  και  $y_0 \in \mathbb{R}^m$  δεδομένη αρχική τιμή. Γράφοντας  $y = (y_1, \dots, y_m)^T$  και  $f = (f_1, \dots, f_m)^T$ , το πρόβλημα (1) διατυπώνεται και ως εξής: Ζητάμε  $y_i: [a, b] \rightarrow \mathbb{R}^1$ ,  $1 \leq i \leq m$ , τέτοια ώστε

$$(1') \quad \begin{cases} y_i'(t) = f_i(t, y_1, \dots, y_m), & 1 \leq i \leq m, \quad a \leq t \leq b, \\ y_i(a) = y_{0,i}, & 1 \leq i \leq m \end{cases}$$

όπου  $f_i: [a, b] \times \mathbb{R}^m \rightarrow \mathbb{R}^1$  δεδομένες συναρτήσεις  $m+1$  μεταβλητών.

Από την θεωρία ύπαρξης-μοναδικότητας λύσεων των συνηθών διαφορικών εξισώσεων\* μας είναι γνωστό ότι για γενικά δεδομένα  $a, b, y_0, f$ , το πρόβλημα (1) μπορεί και να μην έχει λύση ή να έχει

---

\* Βλ. π.χ. τα βιβλία G. Birkhoff and G.-C. Rota, "Ordinary differential equations", 2<sup>nd</sup> ed., Wiley, New York 1969 ή E.A. Coddington and N. Levinson, "Theory of ordinary differential equations", McGraw-Hill, New York 1955 ή C. Corduneanu, "Principles of differential and integral equations", 2<sup>nd</sup> ed., Chelsea, New York 1977, κ.ά.



πολλές λύσεις' βλ. τις Παρατηρήσεις και Ηεκλήσεις αυτής της παραγράφου για μία εύστοχη επισκόπηση. Το παρακάτω θεώρημα δίνει ικανές συνθήκες για ύπαρξη και μοναδικότητα λύσεων για κάθε  $y_0 \in \mathbb{R}^m$ .

**ΘΕΩΡΗΜΑ 1.** Θεωρούμε το πρόβλημα αρχικών τιμών (1) και υποθέτουμε ότι η συνάρτηση  $f$  έχει τις εξής ιδιότητες:

(i) Η  $f(t, y)$  είναι συνεχής για  $(t, y) \in [a, b] \times \mathbb{R}^m$ .

(ii) Η  $f(t, y)$  ικανοποιεί μία συνθήκη Lipschitz ως προς  $y$  για  $y \in \mathbb{R}^m$ , αμοιόμορφα ως προς  $t \in [a, b]$  δηλ. υπάρχει νόρμα  $\|\cdot\|$  του  $\mathbb{R}^m$  και σταθερά  $L \geq 0$  (η "σταθερά Lipschitz" ως προς  $\|\cdot\|$ ) τέτοιες ώστε

$$(2) \|f(t, y) - f(t, y^*)\| \leq L \|y - y^*\|, \quad \forall y, y^* \in \mathbb{R}^m, t \in [a, b].$$

Τότε, για κάθε  $y_0 \in \mathbb{R}^m$ , το πρόβλημα (1) έχει μοναδική λύση' υπάρχει δηλ. μοναδική συνάρτηση  $y: [a, b] \rightarrow \mathbb{R}^m$ , παραγωγίσιμη στο  $[a, b]$  (μάλιιστα έχει συνεχή παράγωγο  $y'(t)$  στο  $[a, b]$ ) που ικανοποιεί το σύστημα των ΣΔΕ  $y'(t) = f(t, y(t))$  για  $a \leq t \leq b$  και την αρχική συνθήκη  $y(a) = y_0$ . @

Μία ικανή συνθήκη για την ισχύ της (2) είναι να είναι οι παράγωγοι  $\partial f_i / \partial y_j$ ,  $1 \leq i, j \leq m$  συνεχείς και φραγμένες συναρτήσεις για  $(t, y) \in [a, b] \times \mathbb{R}^m$ . Τότε, από το θεώρημα μέσης τιμής για συναρτησιακά  $f_i: [a, b] \times \mathbb{R}^m \rightarrow \mathbb{R}^1$  έχουμε ότι για κάθε  $i$ ,  $1 \leq i \leq m$ ,  $y, y^* \in \mathbb{R}^m$  και  $t \in [a, b]$  υπάρχει  $\xi = \xi(i, t, y, y^*) \in \mathbb{R}^m$ , πάνω στο ευθύγραμμο τμήμα  $[y, y^*]$ , τέτοιο ώστε

$$f_i(t, y) - f_i(t, y^*) = \sum_{j=1}^m \frac{\partial f_i}{\partial y_j}(t, \xi)(y_j - y_j^*), \quad 1 \leq i \leq m.$$

Έπεται ότι ισχύει η (2) για κάποια σταθερά  $L = L(M, \|\cdot\|)$  όπου  $M = \max_{i, j} (\sup_{(t, y) \in [a, b] \times \mathbb{R}^m} |\partial f_i / \partial y_j|) < \infty$ .

Οι υποθέσεις του θεωρήματος 1 είναι αρκετά περιοριστικές, ιδιαίτερα βέβαια η (2) που υποθέσαμε ότι ισχύει  $\forall y, y^* \in \mathbb{R}^m$ . Σάν επακόλουθο όμως παίρνουμε ύπαρξη και μοναδικότητα της λύσης  $y(t)$

ε' όλο το διάστημα  $[a, b]$  και για κάθε αρχική τιμή  $y_0 \in \mathbb{R}^m$ . Μία τοπική συνθήκη Lipschitz (δηλ. μία συνθήκη της μορφής (2) για  $(t, y) \in [a, b] \times \bar{S}$  όπου  $\bar{S} = \bar{S}(y_0, r)$  π.χ.) και η συνέχεια της  $f$  στο  $[a, b] \times \bar{S}$  εχθύνονται ύπαρξη-μοναδικότητα της λύσης  $y(t)$  μόνο τοπικά, π.χ. ε' ένα διάστημα της μορφής  $[a, \delta)$ , δ>a - βλ. τις Παρατηρήσεις και Ηεκρήσεις αυτής της παραγράφου. Είναι δυνατόν να αναλύσουμε την εύκλιση προεχχιστικών μεθόδων για την λύση του (1) ακόμα και όταν αντί της (2) πληρούται μόνο μία τοπική συνθήκη Lipschitz (βλ. π.χ. Ηεκ. 2). Για να αποφύγουμε όμως τεχνικές πολυπλοκότητες θα υποθέσουμε από δω κι εμπρός εισηηρά ότι ισχύουν οι υποθέσεις του θεωρήματος 1 και (ευσενώς) ότι υπάρχει μοναδική λύση  $y(t)$  του (1) με συνεχή παράγωγο στο  $[a, b]$ . Συνήθως θα υποθέτουμε επιηλέον ότι η λύση είναι αρκετά ομαλή (δηλ. έχει ένα απαιτούμενο αριθμό υψηλοτέρων συνεχών παραγώγων  $y^{(j)}(t) = d^j y(t)/dt^j$ ) ώστε να παίρνουμε την επιθυμητή τάξη εύκλισης των μεθόδων.

Γιά γενικό δεύτερο μέλος  $f(t, y)$  (ακόμα και όταν έχουμε ένα γενικό γραμμικό ομογενές σύστημα, δηλ. με  $f(t, y) = A(t)y$  όπου  $A: [a, b] \rightarrow \mathbb{R}^{m \times m}$ ) είναι γνωστό ότι δεν μπορούμε να λύσουμε το (1) αναλυτικά. Καταφεύγουμε λοιπόν σε αριθμητικές μεθόδους για την προεχχιστική του λύση. Έστω ο διαμερισμός  $a = t^0 < t^1 < \dots < t^N = b$  του  $[a, b]$ , τον οποίο υποθέτουμε για απλούστευση ομοιόμορφο, δηλ. ότι  $t^n = a + nh$ ,  $0 \leq n \leq N$ , όπου  $h = (b-a)/N$  είναι το (εταθερό) βήμα του διαμερισμού. Μία αριθμητική μέθοδος για την λύση του (1) κατασκευάζει διανύσματα  $y^n \in \mathbb{R}^m$ ,  $0 \leq n \leq N$ , τέτοια ώστε  $y^n \approx y(t^n)$ ,  $0 \leq n \leq N$  όπου  $y(t)$  η λύση του (1).

Η απλούστατη μέθοδος για την προεχχιστική λύση του (1) είναι ίσως η μέθοδος του Euler (ή "πολυχωνική μέθοδος του Cauchy"). Αν το  $h = t^{n+1} - t^n$  είναι αρκετά μικρό, τότε η  $y'(t^n)$  προσεχχίζεται ικανοποιητικά από το ηηλίο διαφορών  $(y(t^{n+1}) - y(t^n))/h$  και η Σ.Δ.Ε. στην (1) δίνει για  $t = t^n$  την προεχχιστική εξέση  $y(t^{n+1}) \approx y(t^n) + h f(t^n, y(t^n))$ . Η μέθοδος του Euler στηρίζεται ακριβώς ε' αυτήν την παράσταση: Κατασκευάζει  $y^n \in \mathbb{R}^m$ ,  $0 \leq n \leq N$ , προεχχίσεις των τιμών  $y(t^n)$ , από τις εξέσεις

$$(3) \begin{cases} y^0 = y_0 \\ y^{n+1} = y^n + hf(t^n, y^n), n=0, 1, \dots, N-1. \end{cases}$$

Η μέθοδος του Euler (3) είναι προφανώς πολύ απλή στην υλοποίησή της. Απαιτεί για κάθε βήμα  $h$  τον υπολογισμό της (διανυσματικής) συνάρτησης  $f(t^n, y^n)$ ,  $m$  πολλαπλασιασμούς και  $m$  προσθέσεις. (Στην περιοχή αυτή της αριθμητικής ανάλυσης έχει επικρατήσει να εκτιμάται η πολυπλοκότητα των προβλημάτων από τον αριθμό των υπολογισμών της  $f(t, y)$  που απαιτούνται ανά βήμα και τον αριθμό των βηματιών - γραμμικών ή μη - αλγεβρικών εξισώσεων που τυχόν απαιτούνται ανά βήμα). Συνεπώς η μέθοδος του Euler έχει πολύ χαμηλό κόστος (ανά βήμα). Ας προχωρήσουμε να εξετάσουμε την ακρίβειά της, δηλ. ας προσπαθήσουμε να εκτιμήσουμε τα εφάλματα  $\|y^n - y(t^n)\|$ ,  $0 \leq n \leq N$ .

**ΛΗΜΜΑ 1** Έστω ότι υπάρχουν σταθερές  $\delta > 0$  και  $K \geq 0$ , τέτοιες ώστε η ακολουθία των μη αρνητικών αριθμών  $d_0, d_1, \dots$  να ικανοποιεί τις σχέσεις

$$(4) \quad d_{j+1} \leq d_j(1+\delta) + K, \quad j=0, 1, 2, \dots$$

Τότε ισχύει για κάθε  $n \geq 0$

$$(5) \quad d_n \leq d_0 e^{n\delta} + K(e^{n\delta} - 1)/\delta.$$

Απόδειξη: Η (4) δίνει  $d_n \leq (1+\delta)^n d_0 + K[1 + (1+\delta) + \dots + (1+\delta)^{n-1}] = (1+\delta)^n d_0 + K[(1+\delta)^n - 1]/\delta$  απ' την οποία έπεται η (5) λόγω της  $1+\delta \leq e^\delta$ . @

**ΘΕΩΡΗΜΑ 2** Έστω  $y(t)$  η λύση του (1). Υποθέτουμε (επιπλέον των υποθέσεων του θεωρήματος 1) ότι η  $y^{(2)}(t)$  είναι συνεχής στο  $[a, b]$ . Έστω  $\{y^n\}$   $0 \leq n \leq N$  η λύση της μεθόδου του Euler (3). Τότε υπάρχει σταθερά  $C_1$  (είναι η σταθερά σύγκρισης  $\|x\| \leq C_1 \|x\|_\infty$  στον  $\mathbb{R}^m$ ) τέτοια ώστε για  $0 \leq n \leq N$

$$(6) \|y^n - y(t^n)\| \leq h \{C_1 (e^{L(t^n - a)} - 1) / 2L\} \max_{t \in [a, t^n]} \|y^{(2)}(t)\|_\infty,$$

όπου η  $L$  η σταθερά Lipschitz της (2)

Απόδειξη: Έστω  $e^n = y^n - y(t^n)$ ,  $0 \leq n \leq N$  το σφάλμα της  $y^n$ . Εκ κατασκευής  $e^0 = 0$ . Για  $1 \leq n \leq N$  από το θεώρημα του Taylor έχουμε ότι

$$(7) y(t^{j+1}) = y(t^j) + hy^{(1)}(t^j) + p^j, \quad 0 \leq j \leq n-1$$

όπου  $p^j = \frac{h^2}{2} y^{(2)}(\xi_j)$

$$p_i = h^2 y^{(2)}(\xi_i) / 2, \quad 1 \leq i \leq m$$

για κάποια  $\xi_i \in [t^j, t^{j+1}]$ . Συνεπώς

$$(8) \|p^j\| \leq C \max_i \|p_i\| \leq K_j = C_1 h^2 \max_{t \in [a, t^{j+1}]} \|y^{(2)}(t)\|_\infty / 2.$$

Αφαιρούμε τώρα κατά μέλη την (7) από την αναδρομική σχέση της (3). Έχουμε, χρησιμοποιώντας την ΣΔΕ (1)

$$e^{j+1} = e^j + h [f(t^j, y^j) - f(t^j, y(t^j))] - p^j.$$

Συνεπώς, χρησιμοποιώντας τις (2), (8) έχουμε για  $0 \leq j \leq n-1$

$$\|e^{j+1}\| \leq \|e^j\| + hL \|y^j - y(t^j)\| + \|p^j\| \leq (1+hL)\|e^j\| + K_j \leq (1+hL)\|e^j\| + K_n$$

Εφαρμόζοντας τώρα το Λήμμα 1 με  $d_j = \|e^j\|$ ,  $d_0 = \|e^0\| = 0$ ,  $\delta = hL$ ,  $K = K_n$  παίρνουμε την (6) χρησιμοποιώντας το γεγονός ότι  $nh = t^n - a$ . (Σίωπηρά υποθέσαμε ότι  $L > 0$ . Αν  $L = 0$  τότε στην (6), αντί του όρου  $(e^{L(t^n - a)} - 1) / L$  βέτουμε  $(t^n - a)$  @

Το θεώρημα 2 μας δίνει συνεπώς μία εκτίμηση του σφάλματος της μορφής

$$(9) \quad \max_{0 \leq n \leq N} \|y^n - y(t^n)\| \leq Ch,$$

όπου  $C$  μία σταθερά ανεξάρτητη του  $h$  (ή του  $N$ ) που εξαρτάται όμως από τα δεδομένα  $a, b, f$  και την λύση  $y(t)$  του (1). Παρατηρούμε ότι το φράγμα στην (9) είναι γραμμικό ως προς  $h$ . Η δύναμη αυτή του  $h$  δεν μπορεί να αυξηθεί, όσο ομαλή και αν είναι η λύση (υποθέτουμε ότι  $y''(t) \neq 0$ ): Πράγματι, θεωρείστε το πρόβλημα αρχικών τιμών (στον  $\mathbb{R}^1$ ):  $y' = 2t$ ,  $0 \leq t \leq 1$ ,  $y(0) = 0$  το οποίο φυσικά έχει την μοναδική λύση  $y(t) = t^2$ . Εύκολα βλέπουμε ότι η μέθοδος του Euler με ομοιόμορφο διαμερισμό στο  $[0, 1]$  δίνει  $y^n = n(n-1)h^2$ . Συνεπώς για  $Nh = 1$ , π.χ., έχουμε το εφάλμα  $y(1) - y^N = h$ . Η (9) εκφράζει το ότι αν η  $y(t)$  είναι αρκετά ομαλή στο  $[a, b]$  (τουλάχιστον  $C^2$ ) τότε η μέθοδος του Euler (3) έχει τάξη ακρίβειας ίση με 1. Αντίθετα, αν  $y(t)$  είναι λιγώτερο ομαλή (αλλά τουλάχιστον  $C^1$ ), τότε θα έχουμε γενικά μικρότερη δύναμη του  $h$  στο δεύτερο μέλος της (9), που εξαρτάται από κατάλληλη εκτίμηση του "μέτρου συνέχειας" της  $y'(t)$ .

Τέλος, η ανισότητα (6) ή (9) ευνεπάγεται, υπό τις προϋποθέσεις π.χ. του θεωρήματος 2, ότι η μέθοδος του Euler "ευχκλίνει". Μ' αυτό εννοούμε ακριβώς ότι αν  $n \rightarrow \infty$  και  $h \rightarrow 0$  έτσι ώστε  $t^n = a + nh \rightarrow t$  για κάθε σταθερό  $t \in (a, b]$  (π.χ. αν  $n \rightarrow \infty$ ,  $h \rightarrow 0$  έτσι ώστε  $nh = t - a$ ), τότε  $y^n \rightarrow y(t)$ , όπου, για δεδομένο  $h$ , η  $y^n$  είναι προσέγγιση που παίρνουμε μετά από  $n$  βήματα της μεθόδου του Euler.

#### Παρατηρήσεις

1. Σε μαθήματα Συνήθων Διαφορικών Εξισώσεων συνήθως αποδεικνύονται και τοπικές μορφές του θεωρήματος 1 ύπαρξης και μοναδικότητας λύσεων του προβλήματος (1). Π.χ. ισχύει το εξής:

**ΘΕΩΡΗΜΑ 1'.** Υποθέστε ότι

(1) Η  $f(t, y)$  είναι συνεχής στο σύνολο  $D = \{(t, y) \in \mathbb{R}^{m+1} : |t - a| \leq A, \|y - y_0\| \leq B\}$ . Έστω  $M = \max_{(t, y) \in D} \|f(t, y)\|$ .

(ii) Η  $f$  ικανοποιεί την εξής συνθήκη Lipschitz: υπάρχει σταθερά  $L \geq 0$  τέτοια ώστε για  $(t, y), (t, y^*) \in D$

$$(2') \|f(t, y) - f(t, y^*)\| \leq L \|y - y^*\|.$$

Τότε, υπάρχει μοναδική λύση του προβλήματος αρχικών τιμών (1) τοπικά, δηλ. υπάρχει μοναδική συνάρτηση  $y(t)$  ορισμένη για  $|t - a| \leq \delta = \min\{A, B/M\}$ , που ικανοποιεί την  $y' = f(t, y)$  για  $|t - a| \leq \delta$  και την αρχική συνθήκη  $y(a) = y_0$ . @

Σημειώτεον ότι η συνθήκη (i) εγχυάται μόνη της την ύπαρξη (αλλά όχι την μοναδικότητα) λύσεων του προβλήματος για  $|t - a| \leq \min\{A, B/M\}$ . (Αυτό είναι το λεγόμενο θεώρημα του Peano). Για μοναδικότητα χρειαζόμαστε κάτι παραπάνω από συνέχεια της  $f'$  η συνθήκη Lipschitz (2') είναι "πολύ κούτα" στο να είναι αναγκαία για μοναδικότητα (βλ. θεώρ. 1). Γενικεύοντας λίγο (θεώρημα του Osgood) μπορούμε να αποδείξουμε το εξής: Στο θεώρημα 1' υποθέτουμε ότι ισχύει η (i) και αντικαθιστούμε την (ii) με το αίτημα να ισχύει στο  $D$  η

$$(2'') \|f(t, y) - f(t, y^*)\| \leq \Phi(\|y - y^*\|),$$

όπου  $\Phi(u)$  μη αρνητική συνεχής αύξουσα συνάρτηση ορισμένη στο διάστημα  $[0, 2B]$ , τέτοια ώστε  $\Phi(0) = 0$  και

$$\int_0^{2B} \frac{du}{\Phi(u)} = +\infty$$

(Η συνθήκη Lipschitz αντιστοιχεί σε  $\Phi(u) = Lu$ . Θα είχαμε επίσης μοναδικότητα αν π.χ.  $\Phi(u) = Lu|\ln u|$ ). Τότε ισχύει το συμπέρασμα του θεωρήματος 1'.

2. Λόγω της (6), δηλ. λόγω μιάς εκτίμησης του εφάλματος της μορφής (9), πρέπει να υπολογίζουμε με πολύ μικρό βήμα  $h$  στην μέθοδο του Euler για να πάρουμε μικρό εφάλμα  $\max_n \|y^n - y(t^n)\|$ . Αυτό έχει ως αποτέλεσμα μεγάλο αριθμό πράξεων, δηλ. την αύξηση των εφαλμάτων

ετροχχύλευσης. Η ακριβής ανάλυση των αποτελεσμάτων της αριθμητικής πεπερασμένης ακρίβειας στις πράξεις δεν είναι δυνατόν να χίνει στο γενικό πρόβλημα (1) λόγω του ότι δεν είναι γνωστό το πώς υπολογίζεται η  $f(t, y)$ . Μπορούμε όμως να κάνουμε μιά σειρά από "λογικές" παραδοχές. Π.χ. έστω ότι αντί των ακριβών τιμών  $y^n$  υπολογίζουμε τις τιμές  $\tilde{y}^n$  όπου υποθέτουμε ότι

$$(i) \tilde{y}^0 = y^0 + \delta^0, \quad \|\delta^0\| \leq \delta, \quad \delta \text{ "μικρό"},$$

ότι αντί του  $f(t^n, \tilde{y}^n)$  υπολογίζουμε στην πραγματικότητα το διάυσεμα  $\tilde{f}(t^n, \tilde{y}^n)$  όπου

$$(ii) \tilde{f}(t^n, \tilde{y}^n) = f(t^n, \tilde{y}^n) + \epsilon^n, \quad \|\epsilon^n\| \leq \epsilon, \quad \epsilon \text{ "μικρό"}$$

και τελικά ότι έχουμε και εφάλμα ετροχχύλευσης κατά τον υπολογισμό του  $\tilde{y}^{n+1}$ , δηλ. ότι

$$(iii) \tilde{y}^{n+1} = \tilde{y}^n + hf(t^n, \tilde{y}^n) + \rho^n, \quad \|\rho^n\| \leq \rho, \quad \rho \text{ "μικρό"}.$$

Μπορεί να δείχθει τότε (βλ. Ρσκ. 4) ότι υπό τις προϋποθέσεις του θεωρήματος 2 υπάρχουν σταθερές  $C, C_1, C_2$  ανεξάρτητες των  $h, h$  τέτοιες ώστε

$$(10) \max_n \|\tilde{y}^n - y(t^n)\| \leq C h + C_1 \delta + C_2 (\epsilon + \rho h^{-1}),$$

(όπου η  $C$  είναι η ίδια με την  $C$  της (9), δηλ. όπου  $C = \max_{t \in [a, b]} \|y''(t)\| (e^{L(b-a)} - 1)/2L$ ). Η εκτίμηση αυτή δείχνει ότι υπάρχει μία κρίσιμη τιμή  $h_0$  του βήματος (που εξαρτάται από το συγκεκριμένο πρόβλημα και την ακρίβεια της αριθμητικής) τέτοια ώστε για  $h < h_0$  το εφάλμα όχι μόνο δεν ελαττώνεται αλλά αυξάνεται!

### Βακίσεις 3.1

1. (α) θεωρείστε για  $\epsilon > 0$  σταθερό το πρόβλημα (στον  $\mathbb{R}^1$ )

$$\begin{cases} y' = |y|^{1+\epsilon}, & 0 \leq t \leq b, \\ y(0) = 1 \end{cases}$$

Δείξτε ότι το πρόβλημα δεν έχει λύση αν  $b \geq e^{-1}$ . Γιατί δεν ισχύει το θεώρημα 1; Ποιό είναι το μέγιστο διάστημα  $[0, b]$  ύπαρξης λύσης που δίνει το θεώρημα 1', π.χ. για  $\epsilon = 1$ ; είναι η λύση μοναδική εκεί;

(β) θεωρείστε για  $0 < \epsilon < 1$  σταθερό το πρόβλημα

$$\begin{cases} y' = |y|^{1-\epsilon}, & 0 \leq t \leq b, \\ y(0) = 0 \end{cases}$$

Δείξτε ότι το πρόβλημα δεν έχει μοναδική λύση σε οποιοδήποτε διάστημα  $[0, b]$ . Γιατί δεν ισχύει το θεώρημα 1'; Ποιά λύση προερχίσει πάντα η μέθοδος του Euler για το πρόβλημα; (Δίδαγμα: η συνθήκη Lipschitz  $\|f(t, y) - f(t, y^*)\| \leq L \|y - y^*\|$  δεν μπορεί να αντικατασταθεί με συνθήκη της μορφής  $\|f(t, y) - f(t, y^*)\| \leq L \|y - y^*\|^a$ ,  $a \neq 1$ ).

2. Υποθέστε ότι το πρόβλημα (1) στον  $\mathbb{R}^1$  έχει μοναδική λύση  $y(t) \in C^2[a, b]$  και υποθέστε ότι αντί της (2) ισχύει η συνθήκη Lipschitz

$$|f(t, y) - f(t, y^*)| \leq L |y - y^*| \quad \forall t \in [a, b], y, y^* \in M_\delta,$$

όπου  $\delta > 0$  και όπου  $M_\delta = [m_1 - \delta, m_2 + \delta]$  με  $m_1 \leq y(t) \leq m_2$  για  $t \in [a, b]$ , δηλ. ότι έχουμε συνθήκη Lipschitz μόνο σε μία περιοχή  $M_\delta$  του πεδίου τιμών  $M_0 = [m_1, m_2]$  της λύσης  $y(t)$  του (1) (μία πολύ χρήσιμη και ρεαλιστική συνθήκη). Δείξτε ότι υπάρχει  $h_0 > 0$  τέτοιο ώστε για  $0 < h \leq h_0$ , η μέθοδος του Euler για το πρόβλημα δίνει προεχίσεις  $\{y^n\}$  για τις οποίες ισχύει το ανάλογο της εκτίμησης (6).



3. θεωρείστε ένα μη ομοιόμορφο διαμερισμό  $a=t_0 < t_1 < \dots < t_N=b$  και υποθέστε ότι αν  $h_n = t^{n+1} - t^n$ ,  $0 \leq n \leq N-1$ , είναι το μεταβλητό βήμα, τότε  $\min_n h_n \geq \lambda \max_n h_n$  για κάποια σταθερά  $\lambda > 0$  ανεξάρτητη του  $n$ . Δείξτε ένα φράγμα του εφάλματος της μεθόδου του Euler ανάλογο του (6) όπου  $h = \max_n h_n$ .

4. Αποδείξτε την (10) υπο τις προϋποθέσεις της Παρατήρησης 2.

5. θεωρείστε το 2x2 σύστημα

$$\begin{cases} x' = -y \\ y' = x & t \geq 0 \\ x(0)=1, y(0)=0, \end{cases}$$

που έχει προφανώς την μοναδική λύση  $x(t)=\cos t$ ,  $y(t)=\sin t$ , που ικανοποιεί την σχέση ("νόμο διατήρησης")

$$(*) \quad x^2(t) + y^2(t) = 1, \quad t \geq 0.$$

(α) θεωρείστε την μέθοδο του Euler για το σύστημα (με σταθερό βήμα  $h > 0$ ). Δείξτε ότι δεν ικανοποιεί το διακριτό ανάλογο της (\*), δηλ. ότι  $(x^n)^2 + (y^n)^2 \neq 1$  αν  $n > 1$ . (Μάλιστα δείξτε ότι, για σταθερό  $h$ ,  $\lim_{n \rightarrow \infty} [(x^n)^2 + (y^n)^2] = \infty$ ). Σχεδιάστε στο επίπεδο  $(x, y)$  τα σημεία  $(x^n, y^n)$  για μερικές τιμές του  $n$ .

(β) διερευνήστε αν υπάρχει διακριτό ανάλογο της (\*) για την μέθοδο

$$\begin{cases} (x^{n+1} - x^n)/h = -y^n, & n \geq 0 \\ (y^{n+1} - y^n)/h = x^{n+1}, & n \geq 0 \\ x^0 = 1, y^0 = 0. \end{cases}$$

(γ) 'ίδιο ερώτημα για την μέθοδο ("πεπλεγμένη Euler")

$$\begin{cases} (x^{n+1}-x^n)/h = -y^{n+1}, & n \geq 0 \\ (y^{n+1}-y^n)/h = x^{n+1}, & n \geq 0 \\ x^0=1, & y^0=0 \end{cases}$$

(δ) 'Ιδιο ερώτημα για την μέθοδο ("τραπεζίου")

$$\begin{cases} (x^{n+1}-x^n)/h = -(y^n+y^{n+1})/2, & n \geq 0 \\ (y^{n+1}-y^n)/h = (x^n+x^{n+1})/2, & n \geq 0 \\ x^0=1, & y^0=0 \end{cases}$$

(Παρατηρείστε ότι οι μέθοδοι (γ) και (δ) απαιτούν την επίλυση γραμμικού συστήματος για κάθε βήμα  $n$  - είναι παραδείγματα "πεπλεγμένων" μεθόδων. Σ' όλα τα ερωτήματα θεωρείστε το  $h > 0$  σταθερό και διερευνήστε την συμπεριφορά του  $(x^n)^2 + (y^n)^2$  καθώς αυξάνεται το  $n$ ).

## 3.2-ΜΕΘΟΔΟΙ RUNGE-KUTTA

Στο κεφάλαιο αυτό θα ασχοληθούμε με μία κατηγορία αριθμητικών μεθόδων για το πρόβλημα αρχικών τιμών (3.1.1), με τις λεγόμενες μεθόδους Runge-Kutta (RK)\*. Σκοπός μας θα είναι να κατασκευάσουμε και να αναλύσουμε μεθόδους ανώτερης τάξης ακρίβειας, δηλ. μεθόδους για τις οποίες

$$(1) \max_{0 \leq n \leq N} \|y^n - y(t^n)\| = O(h^p), \quad p > 1.$$

Ο αριθμός  $p$  λέγεται τάξη ακρίβειας της μεθόδου. Για την απόδειξη εκτιμήσεων όπως η (1) υποθέτουμε ότι η λύση  $y(t)$  του (3.1.1) είναι αρκετά ομαλή και ότι η μέθοδός μας υπολογίζει προσεγγίσεις  $y^n$  των τιμών  $y(t^n)$  στους κόμβους  $t^n = a + nh$ ,  $0 \leq n \leq N$  ενός ομοιόμορφου διαμερισμού με βήμα  $h = (b-a)/N$ . Στην (1) ο συμβολισμός  $O(\cdot)$  ερμηνεύεται ως εξής. Έστω  $e: [0, h_0] \rightarrow \mathbb{R}^m$ . Λέμε ότι  $e(h) = O(h^q)$  (όταν  $h \rightarrow 0$ ) αν υπάρχει σταθερά  $C$  ανεξάρτητη του  $h$  τέτοια ώστε

$$(2) \|e(h)\| \leq Ch^q \text{ για } 0 \leq h \leq h_0.$$

Παραδείγματος χάριν για το σφάλμα της μεθόδου του Euler δείξαμε ότι ισχύει η (1) με  $p=1$ , αν  $y \in C^2[a, b]$  (με  $C^k[a, b]$  θα συμβολίζουμε τον χώρο των ευκαμπτήσεων  $y: [a, b] \rightarrow \mathbb{R}^m$  με  $k$  συνεχείς παραχώχους στο  $[a, b]$ ) και ότι η τάξη ακρίβειας δεν αυξάνει αν η  $y$  είναι περισσότερο ομαλή.

Οι μέθοδοι Runge-Kutta προκύπτουν κατά συστηματικό τρόπο από κατάλληλη εφαρμογή κανόνων αριθμητικής ολοκλήρωσης στην Δ.Ε. του (3.1.1) στο διάστημα  $[t^n, t^{n+1}]$ . Ολοκληρώνοντας π.χ. την Δ.Ε. του (3.1.1) ως προς  $t$  έχουμε την ακριβή σχέση:

---

\* Στους Runge (1885) και Kutta (1901) οφείλεται η λεγόμενη "κλασική" μέθοδος Runge-Kutta (21a). Σήμερα η κατηγορία των μεθόδων που ονομάζουμε RK έχει επεκταθεί σημαντικά.

$$(3) \quad y(t^{n+1}) - y(t^n) = \int_{t^n}^{t^{n+1}} f(t, y(t)) dt$$

Η αντικατάσταση του ολοκληρώματος στο δεύτερο μέλος της (3) από το "εμβαδόν του ορθογωνίου" πλάτους  $h = t^{n+1} - t^n$  και "ύψους"  $f(t^n, y(t^n))$  δίνει την προσεχιστική σχέση  $y(t^{n+1}) - y(t^n) \approx hf(t^n, y(t^n))$  που οδηγεί στην γνωστή μας μέθοδο του Euler. Αν ως "ύψος" του ορθογωνίου χρησιμοποιηθεί η τιμή  $f(t^{n+1}, y(t^{n+1}))$  παίρνουμε την λεγόμενη "πεπλεγμένη" (ή "οπισθοδρομική") μέθοδο του Euler

$$(4) \quad y^{n+1} = y^n + hf(t^{n+1}, y^{n+1}), \quad n \geq 0,$$

(θα παίρνουμε πάντα  $y^0 = y(a) = y_0$ ) που σε κάθε βήμα της απαιτεί την λύση ενός πχμ μη γραμμικού συστήματος για το άγνωστο διάνυσμα  $y^{n+1}$ . Το σύστημα αυτό, για  $h$  αρκετά μικρό, έχει πάντα μοναδική λύση (βλ. Παρατήρηση 1 και Πρόταση 1 παρακάτω). Εφαρμόζοντας τώρα τον κανόνα του τραπεζίου στο β' μέλος της (3) παίρνουμε την (επίσης πεπλεγμένη) μέθοδο του τραπεζίου

$$(5) \quad y^{n+1} = y^n + h[f(t^{n+1}, y^{n+1}) + f(t^n, y^n)]/2, \quad n \geq 0.$$

Εξ άλλου η εφαρμογή του κανόνα "του μέσου" δίνει στην (3) την προσεχιστική σχέση

$$y(t^{n+1}) - y(t^n) \approx hf(t^{n+h/2}, y(t^{n+h/2}))$$

όπου η τιμή  $y(t^{n+h/2})$  μπορεί πάλι με κανόνα ορθογωνίου (ή μέθοδο Euler!) στο διάστημα  $[t^n, t^{n+h/2}]$  να προσεχισθεί ως  $y(t^{n+h/2}) \approx y(t^n) + hf(t^n, y(t^n))/2$ . Συνεπώς η εφαρμογή δύο κανόνων αριθμητικής ολοκλήρωσης, του ενός στο  $[t^n, t^{n+1}]$  και του άλλου στο  $[t^n, t^{n+h/2}]$  οδηγεί στην λεγόμενη "βελτιωμένη μέθοδο του Euler" ή "άμεση μέθοδο του μέσου"

$$(6) \quad \begin{cases} \bar{y}^n = y^n + hf(t^n, y^n)/2 \\ y^{n+1} = y^n + hf(t^n + h/2, \bar{y}^n) \end{cases}$$

που απαιτεί τον υπολογισμό της ενδιάμεσης ποσότητας  $\bar{y}^n$  και δύο υπολογισμούς της  $f$  ανά βήμα. Η (6) είναι παράδειγμα μίας άμεσης μεθόδου, δηλ. μίας μεθόδου που υπολογίζει το  $y^{n+1}$  από το  $y^n$  με αντικατάσταση, δηλ. χωρίς να απαιτεί την επίλυση μη γραμμικών ελασμάτων.

Δεν είναι δύσκολο να αποδειχθεί (βλ. Ασκ. 2) ότι η (6) έχει τάξη ακρίβειας  $p=2$ , δηλ. ότι ισχύει  $\max_n \|y^n - y(t^n)\| = O(h^2)$ . Θα δούμε αργότερα ότι η ηηπλεχμένη μέθοδος του Euler είναι πρώτης τάξης ενώ η μέθοδος του τραπεζίου δεύτερης. Προς το παρόν προχωρούμε στον ορισμό της γενικής μεθόδου RK.

Έστω  $\tau_i$ ,  $1 \leq i \leq q$  αριθμοί (ευνήθως  $0 \leq \tau_i \leq 1$ ) και έστω ότι οι κάμβοι  $\tau_i$  και τα βάρη (ευντελεστές)  $a_{ij}$ ,  $1 \leq j \leq q$ ,  $b_i$ ,  $1 \leq i \leq q$ , ορίζουν τους  $q+1$  κανόνες αριθμητικής ολοκλήρωσης:

$$(7) \quad \int_0^{\tau_i} \psi(s) ds \approx \sum_{j=1}^q a_{ij} \psi(\tau_j), \quad 1 \leq i \leq q,$$

$$(8) \quad \int_0^1 \psi(s) ds \approx \sum_{j=1}^q b_j \psi(\tau_j).$$

Ορίζουμε τώρα για το πρόβλημα (3.1.1)  $t^n = a + nh$ , κατά τα γνωστά, και τα σημεία (για  $0 \leq n \leq N-1$ )

$$(9) \quad t^{n,i} = t^n + \tau_i h, \quad 1 \leq i \leq q.$$

Ολοκληρώνοντας την Δ.Ε. του (3.1.1) ως προς  $t$  από  $t^n$  μέχρι  $t^{n,i}$  έχουμε για  $1 \leq i \leq q$ :

$$(10) \quad y(t^{n,i}) - y(t^n) = \int_{t^n}^{t^{n,i}} f(t, y(t)) dt = (\text{με αλλαγή μεταβλητών} \\ t \rightarrow s, t = t^n + sh, 0 \leq s \leq \tau_i) = h \int_0^{\tau_i} f(t^n + sh, y(t^n + sh)) ds \approx (\text{με εφαρμογή} \\ \text{του κανόνα ολοκλήρωσης (7) και χρήση των (9)})$$

$\approx h \sum_{j=1}^q a_{ij} f(t^{n,j}, y(t^{n,j}))$ . Ευτελώς ανάλογα, ολοκληρώνοντας την δ.ε. του (3.1.1) από  $t^n$  μέχρι  $t^{n+1}$  έχουμε (χρησιμοποιώντας τις (8) και (9)) ότι

$$(11) \quad y(t^{n+1}) - y(t^n) \approx h \sum_{j=1}^q b_j f(t^{n,j}, y(t^{n,j})).$$

Οι προεχθιστικές σχέσεις (10), (11) ορίζουν την λεγόμενη γενική μέθοδο Runge-Kutta με  $q$  (ενδιάμεσα) στάδια για τον υπολογισμό της  $y^{n+1}$  αν είναι γνωστή η  $y^n$ ,  $0 \leq n \leq N-1$ , (με  $y^0 = y_0$ ):

$$(12\alpha) \quad y^{n,i} = y^n + h \sum_{j=1}^q a_{ij} f(t^{n,j}, y^{n,j}), \quad 1 \leq i \leq q,$$

$$(12\beta) \quad y^{n+1} = y^n + h \sum_{j=1}^q b_j f(t^{n,j}, y^{n,j}),$$

όπου τα  $t^{n,i}$  έχουν οριστεί από τις (9). Τα (ενδιάμεσα) "στάδια" είναι τα  $q$  διανύσματα  $y^{n,i} \in \mathbb{R}^m$ ,  $1 \leq i \leq q$ , τα οποία υπολογίζονται ως λύση του (μη γραμμικού γενικά)  $m_q \times m_q$  συστήματος που παριστάνουν οι σχέσεις (12α). Στην πράξη συνήθως η μέθοδος γράφεται ισοδύναμα στην εξής μορφή: θέτουμε  $k^{n,i} = f(t^{n,i}, y^{n,i})$ ,  $1 \leq i \leq q$  οπότε οι (12α-β) γίνονται

$$(13) \quad \begin{aligned} k^{n,i} &= f(t^{n,i}, y^n + h \sum_{j=1}^q a_{ij} k^{n,j}), \quad 1 \leq i \leq q, \\ y^{n+1} &= y^n + h \sum_{j=1}^q b_j k^{n,j}. \end{aligned}$$

### 3.2.5

Συνοπώς η γενική μέθοδος RK με  $q$  στάδια ορίζεται μέσω των  $q^2+2q$  σταθερών  $a_{ij}, b_j, \tau_i$  που συνήθως διατάσσουμε στο μητρώο\*

$$(14) \quad \begin{array}{ccc|c} a_{11} & \dots & a_{1q} & \tau_1 \\ \vdots & & \vdots & \vdots \\ a_{q1} & \dots & a_{qq} & \tau_q \\ \hline b_1 & \dots & b_q & \end{array}$$

ή σε συντομογραφία στο μητρώο

$$(14') \quad \begin{array}{c|c} A & \tau \\ \hline b & \end{array}$$

όπου  $A=(a_{ij}) \in \mathbb{R}^{q \times q}$ ,  $\tau=(\tau_i)$ ,  $b=(b_j) \in \mathbb{R}^q$ . Αν ο πίνακας  $A$  είναι (πιθανώς μετά από αναδιάταξη γραμμών) αυστηρά κάτω τριγωνικός, δηλ. αν  $a_{ij}=0$  για  $i \leq j$ , τότε η μέθοδος είναι άμεση, δηλ. δεν χρειάζεται την λύση μη γραμμικού συστήματος αφού τότε τα ενδιάμεσα στάδια υπολογίζονται με απλή αντικατάσταση:

$$\begin{aligned} y^{n,1} &= y^n \\ y^{n,2} &= y^n + h a_{21} f(t^{n,1}, y^{n,1}) \\ &\vdots \\ y^{n,q} &= y^n + h \sum_{j=1}^{q-1} a_{qj} f(t^{n,j}, y^{n,j}) \end{aligned}$$

Σε κάθε άλλη περίπτωση η μέθοδος είναι πεπλεγμένη και ο υπολογισμός των  $y^{n,i}$  απαιτεί επίλυση μη γραμμικών συστημάτων. Σημαντικές μεταξύ των πεπλεγμένων είναι οι λεγόμενες ημιπεπλεγμένες μέθοδοι για τις

\* Ο φορμαλισμός της μεθόδου οφείλεται στον J.C. Butcher: (Math. Comp. 18(1964), 50-64 και 233-244, ibid. 26 (1972), 79-106 κ.ά.)

οποίες ο πίνακας  $A$  είναι κάτω τριγωνικός, δηλ. έχει  $a_{ij}=0, i < j$ . Τότε τα στάδια υπολογίζονται από τις εκθέσεις

$$y^{n,1} = y^n + h a_{11} f(t^{n,1}, y^{n,1})$$

$$y^{n,2} = y^n + h a_{21} f(t^{n,1}, y^{n,1}) + h a_{22} f(t^{n,2}, y^{n,2})$$

$$y^{n,q} = y^n + h \sum_{j=1}^q a_{qj} f(t^{n,j}, y^{n,j}).$$

Τα μη γραμμικά συστήματα είναι τώρα αποευνδεδεμένα: Η πρώτη εξίσωση δίνει το  $y^{n,1}$  ως λύση ενός  $m \times m$  μη γραμμικού συστήματος και γενικά η  $i$ -ετή δίνει το  $y^{n,i}$  ως λύση του  $m \times m$  μη γραμμικού συστήματος (αν  $a_{ii} \neq 0$ )

$$y^{n,i} = h a_{ii} f(t^{n,i}, y^{n,i}) + g^{n,i}$$

όπου  $g^{n,i}$  ήδη γνωστό διάνυσμα. (Η επίλυση  $q$   $m \times m$  συστημάτων είναι προτιμότερη από την επίλυση ενός  $qm \times qm$ )

Ας δούμε τώρα παραδείγματα μεθόδων RK. Κατ' αρχήν το μητρώο (για  $q=1$  στάδιο)

$$\begin{array}{c|c} 0 & 0 \\ \hline 1 & 0 \end{array}, \quad \begin{array}{c|c} 1 & 1 \\ \hline 1 & 1 \end{array}$$

περιγράφουν, αντίστοιχα, την μέθοδο του Euler (3.1.3) και την πεπλεγμένη μέθοδο του Euler (4). Το μητρώο ( $q=1$ )

$$(15) \quad \begin{array}{c|c} 1/2 & 1/2 \\ \hline 1 & 1 \end{array}$$

περιγράφει στην μορφή (12) την (πεπλεγμένη) μέθοδο του μέσου

$$(15') \quad y^{n+1} = y^n + hf(t^n+h/2, (y^n+y^{n+1})/2),$$

(της οποίας η άμεση μέθοδος του μέσου (6) αποτελεί "γραμμικοποίηση"). Η τάξη ακρίβειας της (15) είναι  $p=2$



## 3.2.7

Γιά  $q=2$  τώρα, παρατηρούμε ότι η μέθοδος του τραπέζιου (5) περιγράφεται από το μητρώο

$$(16) \begin{array}{cc|c} 0 & 0 & 0 \\ 1/2 & 1/2 & 1 \\ \hline 1/2 & 1/2 & \end{array}$$

ενώ η άμεση μέθοδος του μέσου (6) από το

$$(17) \begin{array}{cc|c} 0 & 0 & 0 \\ 1/2 & 0 & 1/2 \\ \hline 0 & 1 & \end{array}$$

Μία ενδιαφέρουσα οικογένεια ημιπελεγμένων μεθόδων με  $q=2$  στάδια δίνεται από το μονοπαραμετρικό μητρώο για  $\lambda \in \mathbb{R}$

$$(18) \begin{array}{cc|c} \lambda & 0 & \lambda \\ 1-2\lambda & \lambda & 1-\lambda \\ \hline 1/2 & 1/2 & \end{array}$$

Οι μέθοδοι που περιγράφονται από το (18) είναι γενικά τάξης  $p=2$ . Οι τιμές  $\lambda=(1 \pm 3^{-1/2})/2$  δίνουν τις πολύ ενδιαφέρουσες μεθόδους "(2,3) PIARK", τάξης  $p=3$ . Ιδιαίτερα ενδιαφέρουσα είναι η μέθοδος με  $q=2$  στάδια που δίνεται από το μητρώο

$$(19) \begin{array}{cc|c} 1/4 & 1/4-\mu & 1/2-\mu \\ 1/4+\mu & 1/4 & 1/2+\mu \\ \hline 1/2 & 1/2 & \end{array}$$

όπου  $\mu=1/2\sqrt{3}$  είναι η λεγόμενη μέθοδος "Gauss-Legendre δύο σημείων" (Τα  $\tau_i=1/2 \pm \mu$  της (19) είναι οι κάμφοι - και τα  $b_i=1/2$  οι αντίστοιχοι ευτελεστές - του κανόνα ολοκλήρωσης Gauss με 2 σημεία με συνάρτηση βάρους  $w(x)=1$  στο  $[0,1]$ ). Η (19) είναι η μόνη μέθοδος με  $q=2$  στάδια που έχει τάξη ακρίβειας  $p=4$ , έχει επίσης πολλές άλλες ενδιαφέρουσες ιδιότητες.

Μερικές κλασικές άμεσες μέθοδοι RK δίνονται από το μητρώο

$$(20\alpha) \begin{array}{ccc|c} 0 & 0 & 0 & 0 \\ 1/2 & 0 & 0 & 1/2 \\ -1 & 2 & 0 & 1 \\ \hline 1/6 & 2/3 & 1/6 & \end{array}, (20\beta) \begin{array}{ccc|c} 0 & 0 & 0 & 0 \\ 1/3 & 0 & 0 & 1/3 \\ 0 & 2/3 & 0 & 2/3 \\ \hline 1/4 & 0 & 3/4 & \end{array}$$

(και λέγονται, αντίστοιχα, μέθοδοι "Kutta τρίτης τάξης" και "Heun τρίτης τάξης" - και οι δύο με  $q=p=3$ ). Η κλασική μέθοδος Runge-Kutta είναι η άμεση μέθοδος με  $q=p=4$  που δίνεται από το μητρώο (21α) ενώ η μέθοδος (21β) έχει επίσης  $q=p=4$ :

$$(21\alpha) \begin{array}{cccc|cccc} 0 & & & & 0 & & & & 0 \\ 1/2 & 0 & & & 1/2 & & & & 1/3 \\ 0 & 1/2 & 0 & & 1/2 & & & & 2/3 \\ 0 & 0 & 1 & 0 & 1 & & & & 1 \\ \hline 1/6 & 1/3 & 1/3 & 1/6 & & & & & \end{array} (21\beta) \begin{array}{cccc|cccc} 0 & & & & 0 & & & & 0 \\ 1/3 & 0 & & & 1/3 & 0 & & & 1/3 \\ -1/3 & 1 & 0 & & -1/3 & 1 & 0 & & 2/3 \\ 1 & -1 & 1 & 0 & 1 & -1 & 1 & 0 & 1 \\ \hline 1/8 & 3/8 & 3/8 & 1/8 & & & & & \end{array}$$

Προφανώς μας ενδιαφέρει με όσο μικρότερο δυνατό αριθμό βημάτων  $q$  να επιτυχάνουμε όσο το δυνατό μεγαλύτερη τάξη ακρίβειας  $p$ . Όπως θα δούμε υπάρχει περιορισμός: η μέγιστη τάξη ακρίβειας μίας μεθόδου  $q$  σταδίων είναι  $p=2q$ . (Οι μέθοδοι που έχουν αυτήν την ιδιότητα είναι (πλήρως) πεπλεγμένες και γενικεύουν για  $q>2$  την (19)). Για άλλα παραδείγματα μεθόδων RK βλ. τα βιβλία [3.1]-[3.4], [3.6], τις εργασίες του Butcher (op.cit) και του συ συνεργατού του, τις εργασίες του Crouzeix\*, και την βιβλιογραφία του [3.1].

Προχωρούμε τώρα στην ανάλυση των μεθόδων RK. Αρχίζουμε με την απόδειξη ότι για πεπλεγμένες μεθόδους υπάρχει λύση του μη γραμμικού συστήματος που ορίζει τα  $y^{n,i}$  (και συνεπώς και το  $y^{n+1}$ ) συναρτήσει του  $y^n$ . Για ό,τι επακολουθεί υποθέτουμε ειληπηρά ότι ισχύουν οι υποθέσεις του θερήματος 3.1.1 ύπαρξης-μοναδικότητας λύσεων του προβλήματος αρχικών τιμών (3.1.1)

\* Ιδιαίτερα βλ. M. Crouzeix, Thèse, Paris VI, 1975 και Num. Math. 32(1979), 75-82.

**ΠΡΟΤΑΣΗ 1** θεωρούμε μία πεπλεγμένη μέθοδο RK της μορφής (12).  
 Έστω  $|A|$  ο πίνακας  $(|a_{ij}|)$  των απολύτων τιμών των στοιχείων  $a_{ij}$  του  
 (14) και έστω  $\rho(|A|)$  η φασματική του ακτίνα. Τότε για κάθε  $h$  τέτοιο  
 ώστε

$$(22) \quad h \leq h_0 < (L \rho(|A|))^{-1}$$

το σύστημα (12α) έχει μοναδική λύση  $(y^{n,i})$ ,  $1 \leq i \leq q$ .

Απόδειξη: θεωρείτε την απεικόνιση  $F: (\mathbb{R}^m)^q \rightarrow (\mathbb{R}^m)^q$ , που  
 ορίζεται για  $\zeta = (\zeta_1, \dots, \zeta_q)^T \in (\mathbb{R}^m)^q$  (δηλ.  $\zeta_j \in \mathbb{R}^m$ ,  $1 \leq j \leq q$ ) από

$$\zeta = \begin{pmatrix} \zeta_1 \\ \zeta_2 \\ \vdots \\ \zeta_q \end{pmatrix} \mapsto F(\zeta) = \begin{pmatrix} y^n + h \sum_{j=1}^q a_{1j} f(t^{n,i}, \zeta_j) \\ \vdots \\ y^n + h \sum_{j=1}^q a_{qj} f(t^{n,i}, \zeta_j) \end{pmatrix}$$

Για την  $i$ -ετη "συνιστώσα" (διάνυσμα στον  $\mathbb{R}^m$ ) του  $F(\zeta)$  έχουμε λοιπόν  
 λόγω της συνθήκης Lipschitz στην  $f$  ότι για κάθε  $\zeta, \zeta^* \in (\mathbb{R}^m)^q$

$$\|(F(\zeta))_i - (F(\zeta^*))_i\| \leq Lh \sum_{j=1}^q |a_{ij}| \|\zeta_j - \zeta_j^*\|, \quad 1 \leq i \leq q.$$

Για  $x \in (\mathbb{R}^m)^q$  έστω  $[x] \in \mathbb{R}^q$  το διάνυσμα με συνιστώσες  $(\|x_1\|, \dots, \|x_q\|)$   
 όπου  $x = (x_1, \dots, x_q)^T$ ,  $x_i \in \mathbb{R}^m$ . Τότε έχουμε (για διανύσματα  $u \leq v \Leftrightarrow$   
 $u_i \leq v_i \quad \forall i$ )

$$\|F(\zeta) - F(\zeta^*)\| \leq Lh |A| [\zeta - \zeta^*] \quad \forall \zeta, \zeta^* \in (\mathbb{R}^m)^q.$$

Έστω τώρα  $F^2(\zeta) = F(F(\zeta))$  και γενικά  $F^v(\zeta) = F(F^{v-1}(\zeta))$ . Τότε η  
 παραπάνω σχέση δίνει εφαρμοζόμενη κατ' επανάληψιν ότι

$$(23) [F^v(\zeta) - F^v(\zeta^*)] \leq (Lh|A|)^v [\zeta - \zeta^*] \quad \forall \zeta, \zeta^* \in (\mathbb{R}^m)^q$$

Η συνθήκη (22) λέει τώρα ότι,  $hLp(|A|) = p(hL|A|) < 1$ . Συνεπώς η ακολουθία πινάκων  $(Lh|A|)^v$  τείνει στο 0 (γιατί;) καθώς  $v \rightarrow \infty$ . Άρα η (23) δίνει ότι για κάποιο  $v_0$  αρκετά μεγάλο όλα τα στοιχεία του πίνακα  $(Lh|A|)^{v_0}$  θα γίνουν μικρότερα από  $a/q$  όπου  $a = a(v_0) < 1$ . Συνεπώς η (23) δίνει ότι  $\exists a < 1$ ,  $v_0 \in \mathbb{N}$  τέτοια ώστε

$$\| (F^v(\zeta))_i - (F^v(\zeta^*))_i \| \leq a/q \sum_{i=1}^q \|\zeta_i - \zeta_i^*\| \Rightarrow$$

$$\sum_{i=1}^q \| (F^v(\zeta))_i - (F^v(\zeta^*))_i \| \leq a \sum_{i=1}^q \|\zeta_i - \zeta_i^*\|$$

δηλ. ότι η απεικόνιση  $F$  είναι ευσταθής, ως προς την νόρμα

$\| \cdot \|$  :  $\| \zeta \| = \sum_{i=1}^q \|\zeta_i\|$  του  $(\mathbb{R}^m)^q$  ε' όλο του  $(\mathbb{R}^m)^q$ . Έπεται ότι η απεικόνιση  $\zeta \mapsto F(\zeta)$  έχει μοναδικό σταθερό σημείο στον  $(\mathbb{R}^m)^q$ , (γιατί;), δηλ. ότι το σύστημα (12α) έχει μοναδική λύση  $\{y^{n,i}\}$ ,  $1 \leq i \leq q$ . @

Θα αποδείξουμε τώρα ένα βασικό αποτέλεσμα ευστάθειας των μεθόδων RK:

**ΠΡΟΤΑΣΗ 2.** Έστω ότι ισχύουν οι υποθέσεις της πρότασης 1. Με  $y^0 \in \mathbb{R}^m$  δεδομένο θεωρούμε τις προσεγγίσεις  $\{y^n\}$ ,  $0 \leq n \leq N$  που παράγει η μέθοδος RK (12α-β) για  $n=0,1,2,\dots,N-1$ . Θεωρούμε επίσης τα διανύσματα  $z^{n,i}$ ,  $0 \leq n \leq N-1$ ,  $1 \leq i \leq q$  και  $z^n$ ,  $0 \leq n \leq N$  του  $\mathbb{R}^m$  που ορίζουν οι εξισώσεις:

$z^0 \in \mathbb{R}^m$  δεδομένο. Για  $n=0,1,\dots,N-1$ :

$$(24a) \quad z^{n,i} = z^{n+h} \sum_{j=1}^q a_{ij} f(t^{n,j}, z^{n,j}), \quad 1 \leq i \leq q,$$

$$(24\beta) \quad z^{n+1} = z^n + h \sum_{j=1}^q b_j f(t^{n,j}, z^{n,j}) + e^n,$$

που προφανώς έχουν επίσης μοναδική λύση για  $0 \leq n \leq R-1$ . (Το σύστημα (24α,β) είναι μία "διαταραχή" του συστήματος (12α,β)). Τότε ισχύει ότι

$$(25) \quad \max_{0 \leq n \leq N} \|y^n - z^n\| \leq C_1 \|y^0 - z^0\| + C_2 h^{-1} \max_{0 \leq n \leq N} \|e^n\|$$

όπου οι σταθερές  $C_1$  και  $C_2$  είναι ανεξάρτητες του  $h$ .

Απόδειξη: Χρησιμοποιούμε τον υπολογισμό της απόδειξης της προηγούμενης πρότασης. Αφαιρώντας κατά μέλη τις (24α), (12α) και χρησιμοποιώντας την συνθήκη Lipschitz για την  $f$  έχουμε

$$\|y^{n,i} - z^{n,i}\| \leq \|y^n - z^n\| + hL \sum_{j=1}^q |a_{ij}| \|y^{n,j} - z^{n,j}\|, \quad 1 \leq i \leq q.$$

Έστω  $Y^n = (y^{n,1}, \dots, y^{n,q})^T$ ,  $Z^n = (z^{n,1}, \dots, z^{n,q})^T$  διανύσματα του  $(\mathbb{R}^m)^q$  με "ευτεταχμένες"  $y^{n,i}, z^{n,i} \in \mathbb{R}^m$  αντίστοιχα. Η παραπάνω σχέση (θυμόμαστε ότι το  $[Y^n]$  είναι το διάνυσμα  $(\|y^{n,1}\|, \dots, \|y^{n,q}\|)^T$  του  $\mathbb{R}^q$ ) γράφεται ως:

$$(26) \quad [Y^n - Z^n] \leq \|y^n - z^n\| u + hL |A| [Y^n - Z^n],$$

όπου  $u = (1, 1, \dots, 1)^T \in \mathbb{R}^q$ . Η (αυτή) αναδρομική ανισότητα (26) δίνει τώρα εφαρμοζόμενη κατ' επανάληψιν για κάθε  $v \in \mathbb{N}$ , (όπου  $1_q$  είναι η ταυτότητα στον  $\mathbb{R}^q$ ):

$$[Y^n - Z^n] \leq \|y^n - z^n\| (1_q + hL|A| + \dots + (hL)^v |A|^v) u + (hL)^{v+1} |A|^{v+1} [Y^n - Z^n]$$

απ' την οποία, χρησιμοποιώντας την  $h \leq h_0$  και παίρνοντας το όριο για  $v \rightarrow \infty$  έχουμε (γιατί:)

$$[Y^n - Z^n] \leq \|y^n - z^n\| (1_q - h_0 L |A|)^{-1} u$$

Συμπεραίνουμε ότι υπάρχει σταθερά  $C$  ανεξάρτητη των  $h, n$  τέτοια ώστε

$$(27) \|y^{n,i} - z^{n,i}\| \leq C \|y^n - z^n\|, \quad 1 \leq i \leq q, \quad 0 \leq n \leq N-1.$$

Τώρα, αφαιρώντας τις (12β), (24β) κατά μέλη έχουμε για  $0 \leq n \leq N-1$

$$(28) \|y^{n+1} - z^{n+1}\| \leq \|y^n - z^n\| + Lh \sum_{j=1}^q |b_j| \|y^{n,i} - z^{n,i}\| + \|e^n\|$$

Αντικαθιστώντας την (27) στην (28) έχουμε με  $C' = LC \sum_{j=1}^q |b_j|$

$$\|y^{n+1} - z^{n+1}\| \leq (1 + C'h) \|y^n - z^n\| + \|e^n\|, \quad 0 \leq n \leq N-1,$$

από την οποία και το Λήμμα 3.1.1 παίρνουμε ότι

$$\max_n \|y^n - z^n\| \leq e^{C'Nh} \|y^0 - z^0\| + \max_n \|e^n\| (e^{C'Nh} - 1) / C'h,$$

απ' όπου προκύπτει η (25) με  $C_1 = e^{C'(b-a)}$ ,  $C_2 = (e^{C'(b-a)} - 1) / c'$  @.

Η Πρόταση 2 ("ευεταθεια") είναι ένα από τα δύο κύρια αποτελέσματα που χρειαζόμαστε για την απόδειξη ενός θεωρήματος "εύγκλισης" (και εκτίμησης του σφάλματος) για τις μεθόδους RK. Το άλλο αφορά την τάξη (ακρίβειας) της μεθόδου.

Υποθέτουμε ότι το δεύτερο μέλος  $f(t, y)$ , και συνεπώς και η λύση  $y(t)$  του προβλήματος αρχικών τιμών (3.1.1), είναι αρκετά ομαλά. Θεωρούμε την μέθοδο RK (12). Για να ορίσουμε την τάξη (ακρίβειας) της, ορίζουμε πρώτα για  $0 \leq n \leq N-1$  το  $\delta^n \in \mathbb{R}^m$  από τις εξισώσεις:

$$(29\alpha) \zeta^{n,i} = y(t^n) + h \sum_{j=1}^q a_{ij} f(t^{n,j}, \zeta^{n,j}), \quad 1 \leq i \leq q,$$

$$(29\beta) \delta^n = y(t^{n+1}) - y(t^n) - h \sum_{j=1}^q b_j f(t^{n,j}, \zeta^{n,j}),$$

όπου υποθέτουμε βέβαια ότι  $h \leq h_0 < (L_p(|A|))^{-1}$  έτσι ώστε (βλ. Πρόταση 1) το σύστημα (29α) να έχει μοναδική λύση  $(\zeta^{n,i})$ ,  $1 \leq i \leq q$ . Λέμε ότι η μέθοδος (12) έχει τάξη (ακρίβειας)  $p$  αν  $\varepsilon^n = O(h^{p+1})$ , δηλ. αν υπάρχει σταθερά  $D$ , ανεξάρτητη του  $h$  ή  $N$  (που μπορεί να εξαρτάται όμως από τα δεδομένα και την λύση του προβλήματος (3.1.1) και την συγκεκριμένη μέθοδο RK) τέτοια ώστε για  $h$  αρκετά μικρό να ισχύει

$$(30) \max_{0 \leq n \leq N} \|\varepsilon^n\| \leq D h^{p+1}.$$

Λέμε ότι η μέθοδος (12) είναι ευνεπής αν η τάξη της  $p$  είναι τουλάχιστον 1. Τονίζουμε το γεγονός ότι οι ποσότητες  $y(t^n)$ ,  $y(t^{n+1})$  στις (29α-β) είναι τιμές της λύσης του προβλήματος (3.1.1). Η τάξη λοιπόν εκφράζει το κατά πόσον ένα βήμα της μεθόδου RK (12) με αρχική συνθήκη  $y(t^n)$  δίνει μία καλή προσέγγιση της τιμής  $y(t^{n+1})$ .

Η γνώση της τάξης μίας μεθόδου (δηλ. μίας εκτίμησης του τύπου (30)) και η "ευστάθεια" της που αποδείχθη στην Πρόταση 2 μάς εξασφαλίζει ένα φράγμα  $O(h^p)$  στο εφάλμα της μεθόδου, όπως μπορούμε εύκολα να διαπιστώσουμε:

**ΘΕΩΡΗΜΑ 1** Έστω ότι ισχύουν οι υποθέσεις της Πρότασης 1. Με  $y^0 = y_0 \in \mathbb{R}^m$  θεωρούμε την μέθοδο RK (12) για την οποία υποθέτουμε ότι ισχύει η (30), δηλ. η (12) έχει τάξη  $p$ . Τότε, αν η λύση  $y(t)$  του (3.1.1) είναι αρκετά ομαλή, έχουμε

$$(31) \max_{0 \leq n \leq N} \|y^n - y(t^n)\| \leq C h^p,$$

όπου  $C = O(e^{c'(b-a)} - 1) / C'$ , η σταθερά  $C'$  ορίστηκε στην απόδειξη της Πρότασης 2.

Απόδειξη: Οι εξισώσεις (29α,β) γράφονται, για  $0 \leq n \leq N-1$ :

$$\zeta^{n,i} = y(t^n) + h \sum_{j=1}^q a_{ij} f(t^{n,j}, \zeta^{n,j}), \quad 1 \leq i \leq q$$

$$y(t^{n+1}) = y(t^n) + h \sum_{j=1}^q b_j f(t^{n,j}, \zeta^{n,j}) + \delta^n,$$

δηλ. είναι της μορφής (24 $\alpha, \beta$ ) με  $z^n \equiv y(t^n)$ ,  $z^{n,i} \equiv \zeta^{n,i}$ ,  $e^n \equiv \delta^n$ . Συνεπώς η (25) και η υπόθεση ότι  $y^0 = y(t^0) \equiv z^0$  δίνουν

$$\max_{0 \leq n \leq N} \|y^n - y(t^n)\| \leq C_2 h^{-1} \max_{0 \leq n \leq N} \|\delta^n\|.$$

Η (31) προκύπτει τώρα από την παραπάνω ανισότητα, την μορφή της  $C_2$  από την απόδειξη της Πρότασης 2, και την (30). Σημειώνουμε ότι η σταθερά  $C'$  εξαρτάται από την μέθοδο RK (δηλ. τις σταθερές  $a_{ij}, b_j$ ), το  $h_0$  και την σταθερά Lipschitz  $L$  της  $f$ , ενώ η σταθερά  $D$  εξαρτάται από την μέθοδο RK και νόρμες παραχώχου της  $y$  και της  $f$ . @

Το θεώρημα 1 εκφράζει ποσοτικά την γενική αρχή ότι "ευστάθεια" + "ευέπεια"  $\Rightarrow$  "εύγκλιση". Η προσοχή μας στρέφεται λοιπόν τώρα στην διερεύνηση της τάξης ακρίβειας μίας δεδομένης μεθόδου, δηλ. στην απόδειξη μίας ανισότητας της μορφής (30).

Αν η μέθοδος είναι άμεση, τότε, με απλή αντικατάσταση, μπορούμε να απαλείψουμε τα  $\zeta^{n,i}$  από τις (29 $\alpha, \beta$ ) και να εκφράσουμε το  $\delta^n$  από την (29 $\beta$ ) ως

$$(32) \quad \delta^n = y(t^{n+1}) - y(t^n) - h \Phi(t^n, y(t^n), h),$$

όπου η  $\Phi$  θα είναι μιά (πολύπλοκη γενικά) συνάρτηση. Αναπτύσσοντας τώρα κατά Taylor την  $y(t^{n+1}) - y(t^n)$  γύρω από το σημείο  $t^n$ , παίρνουμε μιά σειρά δυνάμεων του  $h$  επί παραχώχους  $y^{(j)}(t^n)$ . Οι όροι αυτοί (μέχρι και τάξεως  $h^p$ ) θα πρέπει να απαλειφθούν από ανάλογους όρους του αναπτύχματος της συνάρτησης  $\Phi(t^n, y(t^n), h)$  γύρω από το σημείο  $(t^n, y(t^n))$ .



Ένα τέτοιο ανάπτυγμα της  $\Phi$  δίνει σειρά δυνάμεων του  $h$  οι συντελεστές των οποίων εξαρτώνται από την  $f$  και τις μερικές παραγώγους της. Οι συντελεστές αυτοί θα πρέπει να εκφραστούν ως παράγωγοι της  $y$  στο  $t^n$  μέσω των σχέσεων

$$y'' = f_x, \quad y''' = 0, \quad f_x = f_{x_1} + \sum_{i=1}^m f_{y_i} f_{y_i}, \quad \text{κλπ.}$$

βλ. π.χ. τις Αεκήσεις 2-5. Για πεπλεγμένες μεθόδους και  $h$  αρκετά μικρό πάντα μπορούμε βέβαια να γράφουμε την (29β) στην μορφή (32)· ο προσδιορισμός όμως της συνάρτησης  $\Phi$  απαιτεί την απαλοιφή των  $\zeta^{n,i}$  από τις (29α,β) δηλ. την αναλυτική επίλυση του μη γραμμικού συστήματός (29α) ως προς τους αγνώστους  $\zeta^{n,i}$ . Αυτό βέβαια είναι αδύνατο γενικά· στην πραγματικότητα όμως δεν χρειαζόμαστε ένα τύπο για την  $\Phi(t,y,h)$  αλλά μία δυναμοσειρά της ως προς  $h$ , την οποία μπορούμε να υπολογίσουμε (η δυσκολία αυξάνει εκθετικά με το  $q$ !) με διαδοχικές "ευθθείες" δυναμοσειρών ενδιάμεσων ποσοτήτων (ευνήθως των  $k^{n,i}$ , βλ. (13)). βλ. π.χ. τις Αεκήσεις 6,7 για τις περιπτώσεις απλών Δ.Ε.

Το πρόβλημα της κατασκευής όλων των μεθόδων RK με ορισμένο αριθμό σταδίων  $q$  που έχουν δεδομένη τάξη  $p$  γίνεται εξαιρετικά πολύπλοκο και δύσκολο όσο αυξάνει το  $q$ . Στην Αεκηση 5 βρίσκουμε όλες τις άμεσες μεθόδους (για απλές Δ.Ε.) με  $q \leq 3$ · τέτοιες μέθοδοι έχουν μέγιστη τάξη ακρίβειας  $p=3$ . Στην Αεκηση 7 λύσαμε το ανάλογο πρόβλημα για όλες τις μεθόδους RK με  $q=2$  βρίσκοντας διαδοχικά όλο και περιεωτέρα περιοριστικές συνθήκες ώστε να έχουν τάξη  $p=1,2,3,4$ , αντίστοιχα. Γενικά, για οποιοδήποτε  $q$  είναι πολύ δύσκολο να διατυπωθούν εύκολα ελέγξιμες ικανές και αναγκαίες συνθήκες πάνω στους συντελεστές  $a_{ij}, b_i, \tau_i$  έτσι ώστε η μέθοδος να έχει δεδομένη τάξη ακρίβειας  $p$  όταν εφαρμόζεται στο σύστημα (3.1.1) (βλ. όμως τις εργασίες του Butcher). Θα αποδείξουμε όμως τώρα ένα θεώρημα το οποίο δίνει ικανές συνθήκες, πάνω στους συντελεστές  $a_{ij}, b_i, \tau_i$  της μεθόδου, οι οποίες εξασφαλίζουν ορισμένη τάξη ακρίβειας.

**ΘΕΩΡΗΜΑ 2** (Butcher-Crouzeix). Υποθέτουμε ότι υπάρχουν ακέραιοι  $p, s, r \geq 0$  τέτοιοι ώστε

$$(33) \quad \sum_{i=1}^a b_i \tau_i^k = 1/(k+1), \text{ για } 0 \leq k \leq p-1$$

$$(34) \quad \sum_{j=1}^a a_{ij} \tau_j^k = \tau_i^{k+1} / (k+1), \quad 1 \leq i \leq a, \text{ για } 0 \leq k \leq s-1.$$

$$(35) \quad \sum_{i=1}^a b_i \tau_i^k a_{ij} = b_j (1 - \tau_j^{k+1}) / (k+1), \quad 1 \leq j \leq a, \text{ για } 0 \leq k \leq r-1.$$

$$(36) \quad p \leq r+s+1 \text{ και } p \leq 2s+2.$$

Τότε η μέθοδος RK (12) έχει τάξη  $p$  όταν εφαρμοσθεί στο σύστημα (3.1.1) (στο οποίο υποθέτουμε ότι οι  $y(t), f(t, y)$  είναι αρκετά ομαλές).

Απόδειξη: Συμβολίζουμε για  $k=0, 1, 2, \dots, 0 \leq n \leq N-1$

$$(37) \quad E_n(k) = \sum_{i=1}^a b_i (t^{n,i} - t^n)^k (f(t^{n,i}, y(t^{n,i})) - f(t^{n,i}, \zeta^{n,i})),$$

όπου  $\{\zeta^{n,i}\}_{1 \leq i \leq a}$  είναι η λύση του μη γραμμικού συστήματος (29α) - υποθέτουμε ότι υπάρχει και είναι μοναδική -. Από τον ορισμό (29β) του  $\delta^n$  έχουμε λοιπόν ότι

$$(38) \quad \begin{aligned} \delta^n &= y(t^{n+1}) - y(t^n) - h \sum_{j=1}^a b_j f(t^{n,j}, y(t^{n,j})) + h E_n(0) = \\ &= y(t^{n+1}) - y(t^n) - h \sum_{j=1}^a b_j y'(t^{n,j}) + h E_n(0). \end{aligned}$$

Η υποθέση (33) είναι ισοδύναμη με την δήλωση ότι ο κανόνας ολοκλήρωσης (8) ισχύει ως ιδιότητα για τις συναρτήσεις  $\psi(t) = t^k$ ,  $0 \leq k \leq p-1$ , δηλ. είναι ακριβής για πολυώνυμα βαθμού  $\leq p-1$ . Από την θεωρία αριθμητικής ολοκλήρωσης (βλ. π.χ. [5.2]) έχουμε τότε ότι

$$y(t^{n+1}) - y(t^n) - h \sum_{j=1}^a b_j y'(t^{n,j}) = \int_{t^n}^{t^{n+1}} y'(\tau) d\tau - h \sum_{j=1}^a b_j y'(t^{n,j}) = O(h^{p+1}).$$

Άρα η (38) δίνει ότι

$$\delta^n = hE_n(0) + O(h^{p+1}).$$

Γιά να αποδείξουμε λοιπόν το θεώρημά μας αρκεί να δείξουμε ότι  $E_n(0) = O(h^p)$ . Το αποτέλεσμα αυτό είναι ειδική περίπτωση της σχέσης

$$(39) \quad E_n(k) = O(h^p) \text{ για } k \geq 0$$

την οποία αποδεικνύουμε παρακάτω.

Γιά  $k \geq r$ , η (39) αποδεικνύεται ως εξής: Παρατηρούμε ότι η (34) ισοδυναμεί με την συνθήκη ότι οι κανόνες ολοκλήρωσης (7) είναι όλοι ακριβείς για πολυώνυμα βαθμού  $\leq s-1$ . Άρα έχουμε για  $1 \leq i \leq q$

$$y(t^{n,i}) - y(t^n) = \int_{t^n}^{t^{n,i}} y'(\tau) d\tau = h \sum_{j=1}^a a_{ij} y'(t^{n,j}) + O(h^{s+1}),$$

δηλ. ότι

$$y(t^{n,i}) = y(t^n) + h \sum_{j=1}^a a_{ij} f(t^{n,j}, y(t^{n,j})) + O(h^{s+1}), \quad 1 \leq i \leq q$$

Συνεπώς, η (29α) και η παραπάνω σχέση δίνουν

$$y(t^{n,i}) - \zeta^{n,i} = h \sum_{j=1}^a a_{ij} [f(t^{n,j}, y(t^{n,j})) - f(t^{n,j}, \zeta^{n,j})] + O(h^{s+1}), \quad 1 \leq i \leq q$$

Παίρνοντας υόρμες, χρησιμοποιώντας την συνθήκη Lipschitz στην  $f$  και χρησιμοποιώντας ανάλογους συλλογισμούς μ' αυτούς που οδήγησαν στην ανισότητα (27) της απόδειξης της Πρότασης 2 (όπου  $y^{n,i}$  θέτουμε  $y(t^{n,i})$  και όπου  $y^n - z^n$  θέτουμε  $O(h^{s+1})$ ) παίρνουμε την σχέση

$$(40) \quad y(t^{n,i}) - \zeta^{n,i} = O(h^{s+1}), \quad 1 \leq i \leq q,$$

η οποία, λόγω της (37), δίνει

$$(41) \quad E_n(k) = O(h^{k+s+1}), \quad k \geq 0.$$

Άρα ισχύει η (39) για  $k \geq r$  γιατί υποθέσαμε ((36)) ότι  $p \leq r+s+1$ . Έστω τώρα  $0 \leq k \leq r$ . Παρατηρούμε ότι αν  $(\partial_y f)_{ij} = \partial f_i / \partial y_j$ ,  $1 \leq i, j \leq m$ , τότε

$$\begin{aligned} f(t^{n,j}, y(t^{n,j})) - f(t^{n,j}, \zeta^{n,j}) &= \partial_y f(t^{n,j}, y(t^{n,j})) (y(t^{n,j}) - \zeta^{n,j}) \\ &= O(\|y(t^{n,j}) - \zeta^{n,j}\|^2) = O(h^{2s+2}). \end{aligned}$$

(Υποθέσαμε ότι η  $f$  είναι αρκετά ομαλή - και παραγωγίσιμη άρα - και ότι η παράγωγός της είναι Lipschitz - κάτι που εξασφαλίζεται αν η "δεύτερη" παράγωγος π.χ. είναι συνεχής. Χρησιμοποιήσαμε μετά την Πρόταση 2.1.4 με  $p=1$  και τέλος την (40)). Συμπεραίνουμε ότι

$$E_n(k) = \sum_{i=1}^q b_i (t^{n,i} - t^n)^k \partial_y f(t^{n,i}, y(t^{n,i})) (y(t^{n,i}) - \zeta^{n,i}) + O(h^{2s+2+k}).$$

Έστω τώρα  $\varphi(t) = \partial_y f(t, y(t)) \in \mathbb{R}^{m \times m}$  για  $a \leq t \leq b$  και έστω

$$C(n, \lambda) = \varphi^{(\lambda)}(t^n) / \lambda!, \quad \text{όπου } \varphi^{(\lambda)}(t) = D_t^\lambda \varphi(t).$$

Ο τύπος του Taylor δίνει λοιπόν αν αναπτύξουμε γύρω από το  $t^n$ :  
( $r > k \geq 0$ )

$$\partial_y f(t^{n,i}, y(t^{n,i})) = \sum_{\lambda=0}^{r-k-1} (t^{n,i} - t^n)^\lambda C(n, \lambda) + O(h^{r-k})$$

Άρα, για  $0 \leq k \leq r-1$  αντικαθιστώντας στον παρακάτω τύπο για το  $E_n(k)$  έχουμε με χρήση της (40) ότι

$$\begin{aligned}
 (42) \quad E_n(k) &= \sum_{i=1}^a b_i (t^{n,i} - t^n)^k \left( \sum_{\lambda=0}^{r-k-1} (t^{n,i} - t^n)^\lambda C(n, \lambda) \right) (y(t^{n,i}) - \zeta^{n,i}) \\
 &\quad + O(h^{r+s+1}) + O(h^{2s+2+k}) = \\
 &= \sum_{\lambda=k}^{r-1} C(n, \lambda-k) \left[ \sum_{i=1}^a b_i (t^{n,i} - t^n)^\lambda (y(t^{n,i}) - \zeta^{n,i}) \right] \\
 &\quad + O(h^{r+s+1}) + O(h^{2s+2+k})
 \end{aligned}$$

Τώρα, χρησιμοποιώντας τις εκθέσεις (29α), (35) έχουμε για  $0 \leq \lambda \leq r-1$

$$\begin{aligned}
 (43) \quad &\sum_{i=1}^a b_i (t^{n,i} - t^n)^\lambda (\zeta^{n,i} - y(t^n)) \\
 &= h^{\lambda+1} \sum_{i=1}^a b_i \tau_i \left( \sum_{j=1}^a a_{ij} f(t^{n,i}, \zeta^{n,j}) \right) \\
 &= h^{\lambda+1} \sum_{j=1}^a f(t^{n,j}, \zeta^{n,j}) \left( \sum_{i=1}^a b_i \tau_i a_{ij} \right) \\
 &= h^{\lambda+1} (\lambda+1)^{-1} \sum_{j=1}^a b_j (1 - \tau_j)^{\lambda+1} f(t^{n,j}, \zeta^{n,j})
 \end{aligned}$$

Χρησιμοποιώντας εξ' άλλου πάλι ότι η (33) σημαίνει ότι ο κανόνας ολοκλήρωσης (8) είναι ακριβής για πολυώνυμα βαθμού  $\leq r-1$ , έχουμε

$$\begin{aligned}
 (44) \quad &\sum_{i=1}^a b_i (t^{n,i} - t^n)^\lambda (y(t^{n,i}) - y(t^n)) \\
 &= h^{-1} \int_{t^n}^{t^{n+1}} (t - t_n)^\lambda (y(t) - y(t^n)) dt + O(h^p)
 \end{aligned}$$

$$= h^{-1} \int_{t^n}^{t^{n+1}} (t-t_n)^\lambda \left( \int_{t^n}^t f(x, y(x)) dx \right) dt + O(h^p)$$

= (με ολοκλήρωση κατά μέρη)

$$= (h(\lambda+1))^{-1} \int_{t^n}^{t^{n+1}} (h^{\lambda+1} - (t-t^n)^{\lambda+1}) f(t, y(t)) dt + O(h^p).$$

= (χρησιμοποιώντας τον κανόνα ολοκλήρωσης (8))

$$= (\lambda+1)^{-1} \sum_{j=1}^q b_j (h^{\lambda+1} - \tau_j^{\lambda+1}) f(t^n, j, y(t^n, j)) + O(h^p)$$

$$= h^{\lambda+1} (\lambda+1)^{-1} \sum_{j=1}^q b_j (1 - \tau_j^{\lambda+1}) f(t^n, j, y(t^n, j)) + O(h^p).$$

Αφαιρώντας κατά μέλη την (43) από την (44) έχουμε λοιπόν για  $0 \leq \lambda \leq n-1$

$$\begin{aligned} \sum_{i=1}^q b_i (t^{n,i} - t^n)^\lambda (y(t^{n,i}) - \zeta^{n,i}) &= \\ &= h^{\lambda+1} (\lambda+1)^{-1} \sum_{j=1}^q b_j (1 - \tau_j^{\lambda+1}) (f(t^n, j, y(t^n, j)) - f(t^n, j, \zeta^{n,i})) + O(h^p). \end{aligned}$$

Αντικαθιστώντας στην (42) παίρνουμε για  $0 \leq k \leq n-1$ , επειδή  $p \leq r+s+1$  και  $p \leq 2s+2$  λόγω των υποθέσεων μας

$$E_n(k) = \sum_{\lambda=k}^{r-1} h^{\lambda+1} (\lambda+1)^{-1} C(n, \lambda-k) \left\{ \sum_{i=1}^q b_i (1 - \tau_i^{\lambda+1}) [f(t^n, i, y(t^n, i)) - f(t^n, i, \zeta^{n,i})] \right\} + O(h^p).$$

Ανακαλώντας τώρα τον ορισμό του  $E_n(k)$  από την (37) έχουμε για  $0 \leq k \leq n-1$

$$E_n(k) = \sum_{\lambda=k}^{r-1} h^{\lambda+1} (\lambda+1)^{-1} C(n, \lambda-k) E_n(0) \\ - \sum_{\lambda=k}^{r-1} (\lambda+1)^{-1} C(n, \lambda-k) E_n(\lambda+1) + O(h^p)$$

και επειδή ήδη ξέρουμε ότι  $E_n(r) = O(h^p)$  (γιατί ήδη αποδείξαμε ότι ισχύει η (39) για  $k \geq r$ ), η παραπάνω σχέση γράφεται

$$(45) \quad E_n(k) = \sum_{\lambda=k}^{r-1} h^{\lambda+1} (\lambda+1)^{-1} C(n, \lambda-k) E_n(0) \\ - \sum_{\lambda=k}^{r-2} (\lambda+1)^{-1} C(n, \lambda-k) E_n(\lambda+1) + O(h^p), \quad 0 \leq k \leq r-1.$$

Ας θυμηθούμε ότι η (41) ισχύει για κάθε  $k \geq 0$ . Αντικαθιστώντας την στην (45) έχουμε για  $0 \leq k \leq r-1$

$$(46) \quad E_n(k) = O(h^{s+2k+2}) + O(h^{s+k+2}) + O(h^p) = O(h^{s+k+2}) + O(h^p).$$

Αντικαθιστώντας την σχέση αυτή πάλι στην (45) παίρνουμε για  $0 \leq k \leq r-1$

$$(47) \quad E_n(k) = O(h^{s+2k+3}) + O(h^{s+k+3}) + O(h^p) = O(h^{s+k+3}) + O(h^p).$$

Παρατηρούμε δηλ. ότι με κάθε νέα αντικατάσταση η (45) δίνει μία αυξημένη κατά 1 τάξη στον πρώτο όρο του δεύτερου μέλους. Επαναλαμβάνοντας αυτήν την διαδικασία όσο χρειάζεται καταρθώνουμε τελικά να δείξουμε ότι για κάθε  $0 \leq k \leq r-1$ :

$$(48) \quad E_n(k) = O(h^{s+r+1}) + O(h^p) = (\text{λόγω της υποθέσεως (36)}) = O(h^p)$$

που είναι το ζητούμενο αποτέλεσμα (39) και για  $0 \leq k \leq r-1$ . @

Είδαμε ότι οι πρώτες δύο συνθήκες του θεωρήματος 2 εκφράζουν, αντίστοιχα, ότι οι κανόνες αριθμητικής ολοκλήρωσης (8), αντίστ. (7), είναι ακριβείς για πολύνομα βαθμού  $\leq r-1$ , αντίστ.  $\leq s-1$ . Δηλ. είναι

όπως λέμε τάξης  $p-1$ , αντίστοιχα  $s-1$ . Είναι χρήσιμο να δούμε ένα παράδειγμα του θεωρήματος 2 που βρίσκει κανόνες ευθυθείας (για να έχει η μέθοδος RK(12) τάξη  $p$ ) που διατυπώνονται μόνο μέσω της τάξης των κανόνων ολοκλήρωσης (8) και (7), δηλ. χωρίς χρήση της ευθυθείας (35). Για (μερική) απόδειξη βλ. τις Ασκήσεις 8,9,10.

### ΠΟΡΙΣΜΑ 1 (Butcher-Crouzeix)

(α) Αν ο κανόνας ολοκλήρωσης (8) είναι τάξης  $p-1$  και όλοι οι κανόνες (7) είναι τάξης  $p-2$ , τότε η μέθοδος (12) έχει τάξη  $p$ .

(β) Έστω  $q'$  ο αριθμός των  $\{\tau_i\}$   $1 \leq i \leq q$  που είναι διάφορα μεταξύ τους. Αν όλοι οι κανόνες (7) είναι τάξης  $q'-1$  και ο κανόνας (8) είναι τάξης  $p-1$ , τότε η μέθοδος (12) έχει τάξη  $p$ .

(γ) Υπάρχει μόνο μία μέθοδος RK με  $q$  στάδια που έχει τάξη  $2q$ . Είναι η μέθοδος για την οποία τα  $\tau_i$  και  $b_i$  είναι, αντίστοιχα, οι κόμβοι και οι συντελεστές της ολοκλήρωσης Gauss-Legendre (δηλ. της ολοκλήρωσης Gauss με βάρος  $w(x)=1$ ) στο διάστημα  $[0,1]$ . (Τα  $a_{ij}$  κατασκευάζονται έτσι ώστε οι τύποι ολοκλήρωσης (7) να είναι ακριβείς για πολυώνυμα βαθμού  $\leq q-1$ . Αυτό δίνει ένα  $q^2 \times q^2$  ((34) με  $s=q$ ) γραμμικό αντιστέψιμο σύστημα για τα  $a_{ij}$ ). @

Κλείνουμε αυτήν την παράγραφο αναφέροντας μία αξιοσημείωτη οικογένεια μεθόδων, που δεν ανήκουν στην κατηγορία των μεθόδων RK (12) αλλά αποτελούν κατά κάποιο τρόπο ένα ενδιάμεσο βήμα μεταξύ των άμεσων και των πεπλεγμένων μεθόδων RK. Θεωρούμε για απλοστευση (βλ. Παρατήρηση 2) το αυτόνομο σύστημα  $y' = f(y)$  και υποθέτουμε ότι για  $y \in \mathbb{R}^m$  είναι γνωστός ο  $m \times m$  Ιακωβιανός πίνακας  $J(y) = (\partial_y f)(y)$ . Οι λεγόμενες μέθοδοι Rosenbrock είναι της μορφής:

$$y^{n,1} = hf(y^n) + h\gamma_1 J(y^n) y^{n,1}$$

$$y^{n,2} = hf(y^n + a_{21} y^{n,1}) + h\gamma_2 J(y^n + c_{21} y^{n,1}) y^{n,2}$$

(49α)

$$y^{n,q} = hf(y^n + \sum_{j=1}^{q-1} a_{qj} y^{n,j}) + h\gamma_q J(y^n + \sum_{j=1}^{q-1} c_{qj} y^{n,j}) y^{n,q}$$



$$(49\beta) \quad y^{n+1} = y^n + \sum_{i=1}^q b_i y^{n,i},$$

όπου  $a_{ij}, c_{ij}, 1 \leq j \leq i-1 \leq q-1, b_i, \gamma_i, 1 \leq i \leq q$  δεδομένες σταθερές. Οι μέθοδοι αυτής της μορφής, όπως φαίνεται από την (49α), απαιτούν την λύση  $q$  γραμμικών συστημάτων για τον υπολογισμό των ενδιαμέσων τιμών  $y^{n,i}, 1 \leq i \leq q$  και, όπως θα δούμε στην παράγραφο 3.4, έχουν μερικές από τις επιθυμητές ιδιότητες "απόλυτης ευστάθειας" των πεπλεγμένων μεθόδων RK. Ένα πολύ γνωστό παράδειγμα μεθόδων Rosenbrock είναι η λεγόμενη μέθοδος του Calahan:

$$(50) \quad \begin{cases} y^{n,1} = hf(y^n) + h\gamma J(y^n)y^{n,1} \\ y^{n,2} = hf(y^n + a_{21}y^{n,1}) + h\gamma J(y^n)y^{n,2} \\ y^{n+1} = y^n + b_1 y^{n,1} + b_2 y^{n,2} \end{cases}$$

όπου  $a_{21} = -2/3^{1/2}, \gamma = (1+3^{-1/2})/2, b_1 = 3/4, b_2 = 1/4$  που έχει τάξη ακρίβειας  $p=3$  και απαιτεί την λύση δύο γραμμικών συστημάτων (με τον ίδιο πίνακα  $1-h\gamma J(y^n)$ ) σε κάθε βήμα  $n$  για τον προσδιορισμό των  $y^{n,1}, y^{n,2}$ .

### Παρατηρήσεις

Εν Στήν Πρόταση 1 είδαμε ότι το μη γραμμικό σύστημα (12α) έχει μοναδική λύση  $\{y^{n,i}\}, 1 \leq i \leq q$ , δηλ. το μοναδικό σταθερό σημείο στον  $(\mathbb{R}^m)^q$  της απεικόνισης  $F$  (βλ. απόδειξη της Πρότασης 1). Η λύση  $Y_n = (y^{n,i}) \in (\mathbb{R}^m)^q$  του συστήματος  $Y_n = F(Y_n)$  μπορεί π.χ. να προεγγιεθεί με την απλή επαναληπτική μέθοδο  $Y_n^{[j+1]} = F(Y_n^{[j]})$ ,  $j=0,1,2,\dots$ , ή με την μέθοδο του Νεύτωνα ή (ευνηθέετα) με μία απλουστευμένη μέθοδο του τύπου του Νεύτωνα - π.χ. με την μέθοδο της χορδής -. Ως αρχική τιμή  $Y_n^{[0]}$  της ακολουθίας - ανακαλώντας ότι η λύση  $\{y^{n,i}\}$  του συστήματος αποτελεί προσέγγιση των τιμών  $\{y(t^{n,i})\}$  - μπορούμε π.χ.

να πάρουμε το διάνυσμα  $(y^n, \dots, y^n)^T$  ή ένα διάνυσμα γραμμικών συνδυασμών τιμών  $y^k$ ,  $k \leq n$ , για τις οποίες οι αντίστοιχοι γραμμικοί συνδυασμοί  $y(t^k)$  είναι καλές προσεγγίσεις των  $y(t^{n,i})$ . Τονίζουμε ότι δεν μας ενδιαφέρει η ακριβής λύση των ευστημάτων (49) - μιά και οι τιμές  $y^n$ , είναι τελικά προσεγγίσεις των  $y(t^n)$  - αλλά μιά προσεγγιστική τους λύση που να κατασκευάζεται σχετικά εύκολα (π.χ. με 2 ή 3 ανακυκλώσεις  $j$  το πολύ στην πράξη για κάθε  $n$ ) και που να δίνει εφάλμα αρκετά μικρό έτσι ώστε το συνολικό τοπικό εφάλμα να εξακολουθεί να είναι της τάξης  $h^{p+1}$ , όπου  $p$  η τάξη της αριθμητικής μεθόδου για την λύση της διαφορικής εξίσωσης.

Όπως αναφέραμε ήδη ε' αυτήν την παράγραφο, οι ημιπεπλεγμένες μέθοδοι RK είναι ιδιαίτερα ενδιαφέρουσες από την άποψη της ευκολίας επίλυσης των μη γραμμικών ευστημάτων τους. Ένα ενδιαφέρον υποσύνολο των ημιπεπλεγμένων μεθόδων αποτελούν οι λεγόμενες διαχώνια πεπλεγμένες (DIRK) για τις οποίες τα διαχώνια στοιχεία του πίνακα  $A$  είναι όλα ίσα. Π.χ. οι διαχώνια πεπλεγμένες μέθοδοι με  $q=2$  δίνονται από το μητρώο (18) και είναι γενικά τάξης  $p=2$  εκτός αν το  $\lambda$  είναι μία από τις ρίζες του διωνύμου  $\lambda^2 - \lambda + 1/6$  οπότε  $p=3$ . Αξιοσημείωτες είναι οι διαχώνια πεπλεγμένες μέθοδοι με  $q=3$  στάδια που δίνονται από το μητρώο

$$(51) \quad \begin{array}{ccc|c} \beta & 0 & 0 & \beta \\ 1/2 - \beta & \beta & 0 & 1/2 \\ 2\beta & 1 - 4\beta & \beta & 1 - \beta \\ \hline b_1 & b_2 & b_3 & \end{array}$$

όπου  $b_1 = b_3 = [6(2\beta - 1)^2]^{-1}$  και  $b_2 = 1 - 2b_1$ . Αν το  $\beta$  είναι μία από τις τρεις πραγματικές ρίζες του πολυωνύμου  $\beta^3 - 3\beta^2/2 + \beta/2 - 1/24$ , η μέθοδος (51) έχει τάξη  $p=4$ .

Ο λόγος για τον οποίο οι διαχώνια πεπλεγμένες μέθοδοι είναι ιδιαίτερα ενδιαφέρουσες στην πράξη είναι (εκτός από τις καλές ιδιότητες ευστάθειας τους, βλ. Παρ. 3.4) ο εξής: θεωρούμε, αντί του (3.1.1) το αυτόνομο ευστημα

$$(52) \quad y' = f(y),$$

δηλ. την περίπτωση που η  $f$  δεν εξαρτάται άμεσα από το  $t$ . (Σημειώστε ότι το σύστημα  $y' = f(t, y)$ ,  $y \in \mathbb{R}^m$  μετατρέπεται εύκολα σε αυτόνομο με  $m+1$  - εξισώσεις ως προς την μεταβλητή  $z = (y, t)^T \in \mathbb{R}^{m+1}$ : το σύστημα  $z' = f(z)$  είναι το προηγούμενο με την προεθήκη της εξίσωσης  $t' = 1$ ,  $t(\alpha) = \alpha$ , για την μεταβλητή  $z_{m+1} = t$ ). Για το (50) η επίλυση του μη γραμμικού συστήματος (12α) για μία διαχώνια πεπλεγμένη μέθοδο με  $a_{ii} = \lambda \neq 0$  ανάγεται στη επίλυση  $q$  μη γραμμικών συστημάτων της μορφής

$$(53) \quad y^{n,i} = h\lambda f(y^{n,i}) + z^{n,i}, \quad 1 \leq i \leq q,$$

όπου  $z^{n,i}$  γνωστά διανύσματα. Αν λύσουμε το (51) π.χ. με την μέθοδο της χορδής θα πρέπει για κάθε  $i$  να υπολογίσουμε και να αναλύσουμε σε μορφή LU μόνο μία φορά του Ιακωβιανό πίνακα  $I - h\lambda \partial_y f$ . (Συνηθως χρησιμοποιούμε τον ίδιο πίνακα για όλα τα στάδια  $q$  και του μεταβάλλουμε κάθε 10 ή 20 π.χ. χρονικά βήματα  $n$ ). Σημειώνουμε τέλος ότι οι διαχώνια πεπλεγμένες μέθοδοι είναι υποεύνολο των πεπλεγμένων μεθόδων για τις οποίες ο πίνακας  $H = (a_{ij})$  έχει μία μόνο ιδιοτιμή  $\lambda$  πολλαπλότητας  $q$ : η κατηγορία αυτή των μεθόδων είναι επίσης πολύ αποτελεσματική στην πράξη. (Χρειάζεται κατάλληλη αλλαγή μεταβλητών και χρήση της μορφής Jordan του  $H$ ).

2. Στην πράξη οι μέθοδοι RK χρησιμοποιούνται με μεταβλητό γενικά βήμα  $t^{n+1} - t^n = h_n$  το οποίο είναι επιθυμητό να αυξομειώνεται χωρίς την παρέμβασή μας ανάλογα με το αν η λύση αλλάζει πολύ από βήμα σε βήμα ή αν μεταβάλλεται αρχά και αμαλά. (Μας ενδιαφέρει να κρατάμε το εφάλμα μικρό: θυμηθείτε ότι τα φράγματα των εφαλμάτων είναι της μορφής  $C(y)h^p$  όπου  $C(y)$  κάποια συνάρτηση (ημι)συνάρτησης υψηλών παραγώγων της λύσης  $y(t)$ ). Η τεχνική που χρησιμοποιείται για την "αυτόματη" μεταβολή του βήματος στηρίζεται σε μία εκτίμηση κάποιου τοπικού εφαλματος, το οποίο συνήθως ορίζεται ως η διαφορά  $u(t^{n+1}) - y^{n+1}$ , όπου  $\{y^n\}$  η αριθμητική λύση και  $u(t)$  είναι η λύση, για  $t^n \leq t \leq t^{n+1}$ , του προβλήματος  $u'(t) = f(t, u(t))$ ,  $t^n \leq t \leq t^{n+1}$ ,  $u(t^n) = y^n$ . Αν κατορθώσουμε να ελέγξουμε το τοπικό εφάλμα, τότε το ολικό εφάλμα  $y(t^{n+1}) - y^{n+1} = (y(t^{n+1}) - u(t^{n+1})) + (u(t^{n+1}) - y^{n+1})$  θα είναι επίσης μικρό,

υπό την προϋπόθεση ότι η τιμή  $y^n$  ήταν κοντά στην πραγματική τιμή  $y(t^n)$ . Δεν είναι δύσκολο να δει κανείς (βλ. π.χ. [5.3, παρ. Β4]) ότι το τοπικό σφάλμα προεχχίζεται αρκετά καλά από μία διαφορά της μορφής  $\tilde{y}^{n+1} - y^{n+1}$ , όπου  $\tilde{y}^{n+1}$  μία προσέγγιση της  $y(t^{n+1})$  μεγαλύτερης τάξης ακρίβειας από την  $y^{n+1}$ . Αν  $|\tilde{y}^{n+1} - y^{n+1}| \leq \epsilon h_n$ , όπου  $\epsilon$  κάποιο ανεκτό επίπεδο "σφάλματος ανά βήμα", αποδεχόμαστε την προσέγγιση  $y^{n+1}$  και συνεχίζουμε τον υπολογισμό με το ίδιο  $\Delta t = h_n$ . Αν  $|\tilde{y}^{n+1} - y^{n+1}| > \epsilon h_n$ , την απορρίπτουμε, μειώνουμε το βήμα  $h_n$ , βρίσκουμε νέα  $y^{n+1}$ ,  $\tilde{y}^{n+1}$  και επαναλαμβάνουμε τον έλεγχο κ.ο.κ.

Μας ενδιαφέρει συνεπώς να μπορούμε, χωρίς σημαντική αύξηση του αριθμού των πράξεων ανά βήμα, να υπολογίζουμε και ένα  $\tilde{y}^{n+1}$  με μία ευνοϊκό-μέθοδο μεγαλύτερης τάξης ακρίβειας από την βασική μας μέθοδο που παράγει το  $y^{n+1}$ . Τέτοια ζεύγη μεθόδων είναι π.χ. οι μέθοδοι RKF (Runge-Kutta-Fehlberg) τάξεων (4,5) ή (5,6) κ.λ.π. που χρησιμοποιούνται ευρύτατα στην πράξη για την κατασκευή αλγορίθμων με "αυτόματη" επιλογή βήματος και εκτίμηση του τοπικού σφάλματος για συστήματα Σ.Δ.Ε. που δεν είναι άκαμπτα. Βλ. [5.3, παρΒ4], το βιβλίο [3.7] καθώς και τα άρθρα των Hull et al. στο SIAM J. Num. Anal., 9(1972), 603-637, και των Shampine et al. στο SIAM Review, 18(1976), 376-411, μεταξύ άλλων, για περισσότερες λεπτομέρειες.

### Ασκήσεις 3.2

1. Απαντήστε ετά δύο "γιατί" της απόδειξης της Πρότασης 1 και στο "γιατί" της απόδειξης της Πρότασης 2 διατυπώνοντας και αποδεικνύοντας ανάλογα θεωρήματα.

2(α) θεωρείστε την "άμεση μέθοδο του μέσου" (6) (0 παρακάτω προσδιορισμός της τάξης  $\epsilon$  ακρίβειάς της αποτελεί υπόδειγμα που μπορεί κανείς ν' ακολουθήσει κανείς για άμεσες μεθόδους. Μη χρησιμοποιείτε το θεώρημα 2 ή την Άσκηση 5). Δείξτε ότι η μέθοδος είναι της μορφής

$$y^{n+1} = y^n + h\phi(t^n, y^n; h), \quad n \geq 0$$

όπου

$$\varphi(t, y; h) = f(t+h/2, y+hf(t, y)/2).$$

Ορίστε την ποσότητα  $\delta^n$  όπως στην (32'), δηλ. ως

$$\delta^n = y(t^{n+1}) - y(t^n) - f(t^n+h/2, y(t^n)+hf(t^n, y(t^n))/2), \quad 0 \leq n \leq N-1$$

(β) θεωρείτε το πρόβλημα (3.1.1) στον  $\mathbb{R}^1$ , δηλ. για μία απλή Δ.Ε.. Αναπτύξτε την συνάρτηση  $\delta^n = \delta^n(t^n, y(t^n), h)$  σε δυναμοσειρά του  $h$  αναπτύσσοντας κατά Taylor την  $\varphi(t^n, y(t^n), h)$  ως συνάρτηση δύο μεταβλητών  $f(t^n+h/2, y(t^n)+hk)$  γύρω απ' το σημείο  $(t^n, y(t^n))$  και την διαφορά  $y(t^{n+1}) - y(t^n)$  γύρω από το  $t^n$ . Δείξτε ότι αν οι συναρτήσεις  $f, f_t, f_y, f_{tt}, f_{yy}, f_{ty}$  είναι φραγμένες και συνεχείς για  $t \in [a, b]$ , τότε

$$\max_{0 \leq n \leq N} |\delta^n| \leq Dh^3,$$

και ότι η μέθοδος έχει τάξη ακρίβειας  $p=2$  για απλές Δ.Ε. (Δηλ. ότι η δύναμις 3 στο φράγμα δεν μπορεί να αυξηθεί για γενικά  $f$  και  $y$ ).

(γ) Επεκτείνετε την παρακάτω ανάλυση στον  $\mathbb{R}^m$  και δείξτε ότι για  $y(t), f(t, y)$  ομαλές υπάρχει σταθερά  $D=D(a, b, y, f)$  τ.ψ.

$$\max_n \|\delta^n\| \leq Dh^3$$

δηλ. ότι η μέθοδος είναι τάξης  $p=2$  και για ευστήματα Δ.Ε.

3. (Για όσους αγαπούν τις πράξεις) Ακολουθώντας τα βήματα της Άσκησης 1 (δηλ. χωρίς χρήση του θεωρήματος 2 ή της Άσκησης 5) δείξτε, για απλές Δ.Ε (δηλ. στον  $\mathbb{R}^1$ ), ότι η άμεση μέθοδος RK (20β) έχει τάξη ακρίβειας  $p=3$ . (Το ίδιο ισχύει και για ευστήματα αλλά η λογιστική των αναπτυχμάτων Taylor αυξάνει).

4. (Για όσους αγαπούν πολύ τις πράξεις). Όπως στις Άσκησης 1 και 2 δείξτε ότι η τάξη της κλασικής μεθόδου RK (21α) είναι  $p=4$  για απλές Δ.Ε. (Το ίδιο ισχύει και για ευστήματα).

5. (Βραβείο Runge-Kutta) Σ' αυτήν την άσκηση θα κατασκευάσουμε όλες τις άμεσες μεθόδους RK με αριθμό σταδίων  $q \leq 3$  που έχουν την μεγαλύτερη δυνατή τάξη (για απλές Δ.Ε). Η γενική άμεση μέθοδος με 3 στάδια γράφεται - στη μορφή (13) με  $k^{n,i} = k^i$  - ως:

$$(i) \quad y^{n+1} = y^n + h\phi(t^n, y^n, h), \text{ όπου}$$

$$(ii) \quad \phi(t, y, h) = \sum_{j=1}^3 b_j k^j \quad \text{και όπου}$$

$$(iii) \quad k^1 = f(t, y)$$

$$k^2 = f(t+h\tau_2, y+h\alpha_{21}k^1)$$

$$k^3 = f(t+h\tau_3, y+h(\alpha_{31}k^1 + \alpha_{32}k^2))$$

(Παρατηρούμε ότι η γενική άμεση μέθοδος με  $q=1$ , αντίστοιχα με  $q=2$ , είναι της μορφής (i)-(iii) αν υποθέσουμε ότι  $b_2=b_3=0$ , αντίστοιχα  $b_3=0$ ).

(α) Δείξτε ότι για να έχει η μέθοδος τάξη ακρίβειας  $p \geq 3$  είναι αναγκαίο να έχει η  $\Phi$  ανάπτυγμα της μορφής

$$\Phi_*(t, y, h) = f + hF/2 + h^2(Ff_y + G)/6 + O(h^3)$$

όπου  $F = f_t + ff_y$ ,  $G = f_{tt} + 2ff_{ty} + f^2f_{yy}$ . (Η  $f$  και οι μερικές της παράγωγοι υπολογίζονται στο σημείο  $(t, y)$ . Φυσικά για  $p \geq 1$  είναι αναγκαίο  $\Phi_* = f + O(h)$ . Για  $p \geq 2$ ,  $\Phi_* = f + hF/2 + O(h^2)$ ).

(β) Αναπτύξτε την  $k^2$  σε σειρά Taylor περί το σημείο  $(t, y)$ , χρησιμοποιώντας την σχέση  $k^1 = f$ : υπολογίστε τους όρους τάξης μέχρι και  $O(h^2)$ . Κατόπιν, αντικαθιστώντας στην  $k^3$  το ανάπτυγμα της  $k^2$ , βρείτε παρόμοιο για την  $k^3$ . Αντικαθιστώντας αυτά τα αναπτύγματα στο δεύτερο μέλος της (ii) πάρτε τελικά ένα ανάπτυγμα της μορφής

$$\Phi(t, y, h) = A + Bh + Ch^2 + O(h^3)$$

(γ) Εξισώνοντας όρους του δευτέρου μέλος της  $\Phi_*$  με αντίστοιχους όρους του παραπάνω αναπτύγματος της  $F$  δείξτε ότι:

(γ<sub>1</sub>) Για  $q=1$  (δηλ. με  $b_2=b_3=0$ ), η μόνη άμεση μέθοδος με τάξη ακρίβειας  $p \geq 1$  είναι η μέθοδος του Euler, δηλ. η μέθοδος με  $b_1=1$ ,  $a_{11}=\tau_1=0$ , που έχει  $p=1$ .

(γ<sub>2</sub>) Για  $q=2$  (δηλ. με  $b_3=0$ ) δείξτε ότι υπάρχει μονοπαραμετρική οικογένεια μεθόδων με  $p \geq 2$ , οι σταθερές τους  $a_{21}, \tau_1, b_1, b_2$  ικανοποιούν τις σχέσεις:

$$b_1 + b_2 = 1$$

$$a_{21} = \tau_2$$

$$b_2 \tau_2 = 1/2$$

Δείξτε ότι όλες αυτές οι μέθοδοι έχουν τάξη  $p=2$ , (δηλ. ότι δεν υπάρχει άμεση μέθοδος με  $q=2$ ,  $p > 2$ ). Σε ποιές τιμές των σταθερών αντιστοιχεί η μέθοδος του μέσου; Ξαν τι μπορούν να ερμηνευθούν οι μέθοδοι με  $b_2=1$  και με  $b_2=1/2$ ;

(γ<sub>3</sub>) Για  $q=3$  προκύπτουν (αν θέλουμε  $p \geq 3$ ) οι ευσθήκες:  $b_1 + b_2 + b_3 = 1$ ,  $a_{21} = \tau_2$ ,  $a_{31} + a_{32} = \tau_3$ ,  $b_2 \tau_2 + b_3 \tau_3 = 1/2$ ,  $b_2 \tau_2^2 + b_3 \tau_3^2 = 1/3$ ,  $b_3 \tau_2 a_{32} = 1/6$ . Συνεπώς, υπάρχει διπαραμετρική οικογένεια άμεσων μεθόδων με  $q=3$ ,  $p \geq 3$ . (Μπορεί να δείχθει - με υπολογισμό του όρου  $O(h^3)$  - ότι καμιά από αυτές τις μεθόδους δεν έχει  $p > 3$ , δηλ. ότι για όλες  $p=3$ ). Βεβαιωθείτε ότι οι (20α), (20β) ανήκουν ε' αυτήν την κλάση.

6. (Διάλειμμα!) Δείξτε ότι η μόνη μέθοδος RK με 1 στάδιο και τάξη ακρίβειας 2 είναι η μέθοδος του μέσου (15).

7. (Μέγα βραβείο Runge-Kutta). Θεωρούμε την γενική (πεπλεγμένη) μέθοδο RK με  $q=2$  την οποία γράφουμε στην μορφή (13) - βλ. και άσκηση 5 - δηλ. ως

$$(i) \quad y^{n+1} - y^n = h\phi(t^n, y^n, h), \text{ όπου}$$

$$(ii) \quad \phi(t, y, h) = b_1 k^1 + b_2 k^2, \text{ όπου}$$

$$(iii) \quad \begin{cases} k^1 = f(t+h\tau_1, y+h(a_{11}k^1+a_{12}k^2)) \\ k^2 = f(t+h\tau_2, y+h(a_{21}k^1+a_{22}k^2)) \end{cases}$$

θα προσδιορίσουμε συνθήκες πάνω στα  $a_{ij}, b_i, \tau_i$  ώστε η μέθοδος να έχει τάξη ακρίβειας αντίστοιχα  $p \geq 1, 2, 3, 4$

(α) Δείξτε ότι η αναγκαία συνθήκη για να έχει η μέθοδος τουλάχιστον τάξη  $p=4$  για απλή ΔΕ είναι να έχει η  $\phi$  ανάπτυγμα της μορφής

$$\begin{aligned} \phi_*(t, y, h) = & f + hF/2 + h^2(Ff_y + G)/6 \\ & + h^3[(3f_{ty} + 3ff_{yy} + f_y^2)F + Gf_y + H]/24 + O(h^4) \end{aligned}$$

όπου τα  $F, G$  ορίστηκαν στην Άσκηση 5(α) και όπου

$$H = f_{ttt} + 3ff_{tty} + 3f^2_{tuy} + f^3_{yyy}$$

(β) Αναπτύσσοντας τα  $k^i$  σε σειρά Taylor περί το  $(t, y)$  δείξτε ότι

$$\begin{aligned} (iv) \quad k^i = & f + h[\tau_i f_t + (a_{i1}k^1 + a_{i2}k^2)f_y] \\ & + h^2[\tau_i^2 f_{tt} + 2\tau_i(a_{i1}k^1 + a_{i2}k^2)f_{ty} + (a_{i1}k^1 + a_{i2}k^2)^2 f_{yy}]/2 \\ & + h^3[\tau_i^3 f_{ttt} + 3\tau_i^2(a_{i1}k^1 + a_{i2}k^2)f_{tty} + 3\tau_i(a_{i1}k^1 + a_{i2}k^2)^2 f_{tuy} + \\ & + (a_{i1}k^1 + a_{i2}k^2)^3 f_{yyy}]/6 + O(h^4), \quad i=1, 2 \end{aligned}$$

Επειδή τα αναπτύγματα αυτά είναι "πεπλεγμένα" δεν μπορούμε να προχωρήσουμε με διαδοχικές αντικαταστάσεις όπως στην Άσκηση 5. Υποθέστε λοιπόν ότι υπάρχουν τα αναπτύγματα της μορφής

$$(v) \quad k^i = A_i + hB_i + h^2C_i + h^3D_i + O(h^4), \quad i=1, 2,$$



και αντικαθιστώντας τις (v) στις (iv) και εξισώνοντας δυνάμεις του h, βρείτε ένα σύστημα της μορφής

$$A_i = f$$

$$B_i = B_i(A_1, A_2)$$

$$C_i = C_i(A_1, A_2, B_1, B_2)$$

$$D_i = D_i(A_1, A_2, B_1, B_2, C_1, C_2),$$

που μπορεί να λυθεί με απλή αντικατάσταση. Λύστε το! Βρείτε τέλος το ανάπτυγμα της  $\Psi$  μέσω της (ii) (κρατήστε μέχρι και όρους  $O(h^3)$ ).

(γ) Εξισώνοντας όρους ίσων δυνάμεων του h στα ανάπτυγματα των  $\Psi$  και  $\Psi_*$  βρείτε ευνόηκες έπει ώστε, η μέθοδος (i)-(iii) να έχει, αντίστοιχα, τάξη  $p \geq 1, 2, 3$  και 4.

(δ) Βρείτε την γενική μορφή των ημιπεπλεγμένων μεθόδων με  $p \geq 2$  και με ίσα διαχώνια στοιχεία  $a_{11} = a_{22} = \lambda$ . (βλ. (18)!) Για ποιά  $\lambda$  έχουμε  $p \geq 3$ ; Υπάρχει τέτοια μέθοδος με  $p=4$ ;

(ε) Αποδείξτε ότι υπάρχει μόνο μία μέθοδος με  $q=2$ ,  $p \geq 4$ , δηλ. η (19). (Μπορεί να δείχθει ότι η τάξη της είναι ακριβώς 4, βλ. Πρόγραμμα 1(γ)).

8. Να αποδειχθεί ο ισχυρισμός (α) του Προτάματος 1. (Υπόδειξη: αποδείξτε, όπως στο πρώτο μέρος της απόδειξης του θεωρήματος 2, ότι  $E_n(0) = O(h^p)$ ).

9. Να αποδειχθεί ο ισχυρισμός (β) του Προτάματος 1. (Υπόδειξη: δείξτε ότι ισχύουν οι υποθέσεις του θεωρήματος 2 με  $s=q$  και  $r=p-q$ . Χρησιμοποιήστε την ταυτότητα

$$\int_0^1 \int_0^1 x^k \psi(t) dt dx = \left\{ \int_0^1 (1-x^{k+1}) \psi(x) dx \right\} / (k+1)$$

και δείξτε την (35) παίρνοντας ως  $\psi(t)$  τα πολυώνυμα βαθμού  $\leq q-1$  τέτοια ώστε για δεδομένο  $j$ ,  $\psi(\tau_j) = 0$  αν  $i \neq j$ ,  $\psi(\tau_j) = 1$  και χρησιμοποιώντας το αποτέλεσμα της άσκησης 10(α)).

10. (α) θεωρείστε τον κανόνα αριθμητικής ολοκλήρωσης

$$\int_0^1 \psi(\tau) d\tau \approx \sum_{j=1}^q b_j \psi(\tau_j)$$

με  $q$  διακριτά σημεία  $\tau_j$ ,  $1 \leq j \leq q$ . Δείξτε ότι η τάξη του δεν μπορεί να υπερβαίνει το  $2q-1$ .

(β) Δείξτε ότι η τάξη μιάς μεθόδου RK με  $q$  στάδια της μορφής (12) δεν μπορεί να υπερβαίνει το  $2q$  (Υπόδειξη: θεωρείστε μία Δ.Ε. της μορφής  $y'=f(t)$ ).

(γ) Δείξτε - ανακαλώντας αποτελέσματα της ολοκλήρωσης Gauss' βλ. [5.2] - ότι οι μέθοδοι που περιγράφονται στο Πρόλημα 1(γ) είναι όντως τάξεως  $2q$ . (Η μοναδικότητά τους για κάθε  $q$  είναι πιο δύσκολο να αποδειχθεί' βλ. Butcher (1964)).

11. ('Άσκηση μαρούθ). Για κάθε μία από τις μεθόδους: Euler, πεπλεγμένη Euler, (15)-(17), (18) (δύο περιπτώσεις: είτε  $\lambda$  ρίζα του  $\lambda^2 - \lambda + 1/6 = 0$  είτε όχι), (19), (20α), (21α), (51), συγκρίνετε την πραγματική τους τάξη με την τάξη  $p$  που δίνουν οι (ικανές) συνθήκες των (α), (β) του Προβλήματος 1 καθώς και οι συνθήκες του θεωρήματος 2, παίρνοντας έτσι μιά ιδέα για την ισχύ των συνθηκών αυτών για συννηθισμένες μεθόδους RK..

12. Το παρακάτω αποτέλεσμα δίνει κάποια βάση στον ισχυρισμό "μικρό τοπικό εφάλμα"  $\Rightarrow$  "μικρό ολικό εφάλμα" της Παρατήρησης 2. θεωρούμε την Δ.Ε.

$$w' = f(t, w), \quad t^n \leq t \leq t^{n+1},$$

και υποθέτουμε ότι η  $f$  είναι συνεχής για  $(t, w) \in [t^n, t^{n+1}] \times \mathbb{R}$ , και ότι  $|f(t, w_1) - f(t, w_2)| \leq L|w_1 - w_2|$  για κάθε  $t \in [t^n, t^{n+1}]$ ,  $w_1, w_2 \in \mathbb{R}$ . θεωρείστε δύο λύσεις  $u(t)$ ,  $y(t)$  της Δ.Ε. για  $t \in [t^n, t^{n+1}]$  που αντιστοιχούν στις αρχικές τιμές  $u(t^n)$ ,  $y(t^n)$ . Δείξτε ότι αν  $h_n = t^{n+1} - t^n$ ,  $h_n L < 1$ , τότε

$$\max_{t^n \leq t \leq t^{n+1}} |u(t) - y(t)| \leq |u(t^n) - y(t^n)| / (1 - Lh_n)$$

13. Για μία αυτόνομη Δ.Ε.  $y'=f(y)$  δείξτε ότι η μέθοδος του Calahan (50) έχει τάξη  $p=3$ . (Ορίστε την τάξη της σε αναλογία με ό,τι κάναμε για τις μεθόδους RK. Αναπτύξτε σε δυναμοσειρές του  $h$  τις ενδιάμεσες ποσότητες και αντικαταστήστε στην τελευταία γραμμή).

## 3.3 ΠΟΛΥΒΗΜΑΤΙΚΕΣ ΜΕΘΟΔΟΙ

Σ' αυτήν την παράγραφο θα εξετάσουμε μια δεύτερη μεγάλη κατηγορία μεθόδων για την αριθμητική λύση του προβλήματος αρχικών τιμών (3.1.1), τις λεγόμενες (γραμμικές) πολυβηματικές μεθόδους. Ένα παράδειγμα τέτοιας μεθόδου είναι το εξής: Έστω ότι η λύση  $y(t)$  του (3.1.1) είναι  $C^3$  στο  $[a,b]$ . Τότε, αν  $t, t \pm h \in [a,b]$ , το θεώρημα του Taylor μας δίνει:

$$\begin{aligned} y(t+h) &= y(t) + hy'(t) + \frac{h^2}{2} y''(t) + O(h^3) \\ y(t-h) &= y(t) - hy'(t) + \frac{h^2}{2} y''(t) + O(h^3). \end{aligned}$$

Αφαιρώντας κατά μέλη και χρησιμοποιώντας την Δ.Ε. έχουμε την σχέση

$$y(t+h) - y(t-h) = 2hy'(t) + O(h^3) = 2hf(t, y(t)) + O(h^3),$$

που δίνει (για τον συμβολισμό του ομοιόμορφου διαμερισμού βλ. παρ.3.1) την μέθοδο:

$$(1) \quad y^{n+1} - y^{n-1} = 2hf^n, \quad n=1, 2, \dots, N-1,$$

όπου, από δώ και εμπρός,  $f^k = f(t^k, y^k)$ ,  $0 \leq k \leq N$ .

Η μέθοδος (1) είναι ένα παράδειγμα πολυβηματικής (διβηματικής) μεθόδου: για τον προσδιορισμό της  $y^{n+1}$  απαιτεί γνώση των τιμών  $y^n, y^{n-1}$  των δύο προηγούμενων βημάτων. Επίσης η (1) απαιτεί δύο αρχικές τιμές: την  $y^0$ , που παίρνουμε από το πρόβλημα αρχικών τιμών (3.1.1), και την  $y^1$  που μπορούμε να υπολογίσουμε με μία "μονοβηματική" μέθοδο ή μια μέθοδο RK. Αλλάζοντας τον δείκτη στην (1) θα γράψουμε ευθέως την μέθοδο στην κανονική της μορφή

$$(1') \quad y^{n+2} - y^n = 2hf^{n+1}, \quad 0 \leq n \leq N-2.$$

### 3.3.2

Είναι προφανές ότι χρησιμοποιώντας αναπτύγματα Taylor τιμών  $y(t+mh)$ ,  $m=-1,-2,\dots$  γύρω από το  $t$  μπορούμε να βρούμε και άλλες πολυβηματικές μεθόδους. Μία άλλη τεχνική που χρησιμοποιείται για τον ίδιο σκοπό είναι η αριθμητική ολοκλήρωση. Π.χ. από την ταυτότητα

$$y(t^{n+2}) - y(t^n) = \int_{t^n}^{t^{n+2}} f(t, y(t)) dt,$$

προεχθίζοντας το ολοκλήρωμα του δευτέρου μέλους με τον κανόνα του Simpson, δηλ. χρησιμοποιώντας την σχέση

$$\int_{t^n}^{t^{n+2}} f(t, y(t)) dt \approx (t^{n+2} - t^n) [f(t^{n+2}, y(t^{n+2})) + 4f(t^{n+1}, y(t^{n+1})) + f(t^n, y(t^n))] / 6,$$

παίρνουμε την (επίσης διβηματική) μέθοδο του Simpson

$$(2) \quad y^{n+2} - y^n = (h/3)(f^{n+2} + 4f^{n+1} + f^n), \quad 0 \leq n \leq N-2,$$

η οποία είναι πεπλεγμένη, διότι απαιτεί την επίλυση του μη γραμμικού συστήματος  $y^{n+2} = (h/3)f(t^{n+2}, y^{n+2}) + g^n$  σε κάθε βήμα  $n$ ,  $0 \leq n \leq N-2$ .

Αντίθετα η μέθοδος (1) είναι άμεση.

Γενικά, μια (γραμμική) k-βηματική ( $k \geq 1$ ) μέθοδος για την λύση του (3.1.1) είναι της μορφής:

$$(3) \quad \begin{cases} y^0, y^1, \dots, y^{k-1} \text{ δεδομένα. (Συνήθως υπολογίζονται με μεθόδους RK).} \\ \alpha_k y^{n+k} + \alpha_{k-1} y^{n+k-1} + \dots + \alpha_0 y^n = h(\beta_k f^{n+k} + \beta_{k-1} f^{n+k-1} + \dots + \beta_0 f^n), \\ (\sum_{j=0}^k \alpha_j y^{n+j} = h \sum_{j=0}^k \beta_j f^{n+j}), \quad 0 \leq n \leq N-k, \end{cases}$$

## 3.3.3

όπου  $\{\alpha_j, \beta_j\}$ ,  $0 \leq j \leq k$ , δεδομένες πραγματικές σταθερές ανεξάρτητες του  $h$  ή του  $n$ . Θα υποθέτουμε συνήθως ότι  $\alpha_k = 1$  και ότι  $|\alpha_0| + |\beta_0| > 0$  ώστε να έχουμε πράγματι μια  $k$ -βηματική μέθοδο. Αν  $\beta_k = 0$  η μέθοδος θα λέγεται άμεση: ο υπολογισμός του  $y^{n+k}$  γίνεται με απλή αντικατάσταση των γνωστών τιμών  $y^{n+i}$ ,  $0 \leq i \leq k-1$ . Αν  $\beta_k \neq 0$  ο προσδιορισμός του  $y^{n+k}$  απαιτεί την λύση ενός  $m \times m$  μη γραμμικού συστήματος της μορφής

$$(4) \quad y^{n+k} = h\beta_k f(t^{n+k}, y^{n+k}) + g^n,$$

όπου  $g^n$  γνωστό διάνυσμα. Το (4) έχει προφανώς μοναδική λύση (από το θεώρημα ευστολής) αν  $h|\beta_k| L < 1$ , δηλ. αν πάρουμε αρκετά μικρό  $h$  το  $L$  είναι η σταθερά Lipschitz της  $f$  ως προς  $y$ . Είναι φανερό λοιπόν ότι οι πεπλεγμένες πολυβηματικές μέθοδοι έχουν πολύ μικρότερο κόστος ανά βήμα απ'ότι οι πεπλεγμένες μέθοδοι RK. Οι δε άμεσες πολυβηματικές μέθοδοι απαιτούν σε κάθε βήμα  $n$  ένα μόνο υπολογισμό της  $f$  ( $f^{n+k}$ ). Είναι λοιπόν οι πολυβηματικές μέθοδοι πολύ φθηνότερες από τις μεθόδους RK\*. Δες όμως την παράγραφο 3.4 για εαφή πλεονεκτήματα των (πεπλεγμένων) μεθόδων RK ως προς την "απόλυτη ευστάθεια" τους.

Η κλάση των μονοβηματικών μεθόδων ( $k=1$ ) δηλ. οι μέθοδοι

$$(5) \quad \alpha_1 y^{n+1} + \alpha_0 y^n = h(\beta_1 f^{n+1} + \beta_0 f^n), \quad 0 \leq n \leq N-1$$

με  $\alpha_1 = 1$ ,  $|\alpha_0| + |\beta_0| > 0$ , αποτελεί την "τομή" των πολυβηματικών μεθόδων με τις μεθόδους RK. Επει οι μέθοδοι Euler, πεπλεγμένη Euler και τραπεζίου, μπορούν να μελετηθούν και ως μονοβηματικές μέθοδοι.

\* Γι' αυτό και χρησιμοποιήθηκαν ευρύτατα, από τον 19ο αιώνα ακόμη, για την αριθμητική λύση προβλημάτων αστρονομίας.

### 3.3.4

Υπάρχουν πολλοί τρόποι κατασκευής πολυβηματικών μεθόδων. Είδαμε παραδείγματα χρήσης αναπτυγμάτων Taylor και αριθμητικής ολοκλήρωσης. Άλλες τεχνικές που χρησιμοποιούνται είναι χρήση του πολυωνύμου παρεμβολής και αριθμητικής παραγωγήσις. Π.χ. θεωρώντας το πολυώνυμο παρεμβολής  $P_{k,n}(t)$  βαθμού  $\leq k$  που παρεμβάλλεται στα σημεία  $t^{n+k}, t^{n+k-1}, \dots, t^n$  στις τιμές  $y(t^{n+k}), \dots, y(t^n)$ , αντίστοιχα, και υπολογίζοντας την παράγωγό του στο σημείο  $t^{n+k}$  έχουμε ότι

$$P'_{n,k}(t^{n+k}) \cong y'(t^{n+k}) = f(t^{n+k}, y(t^{n+k})).$$

Αντικαθιστώντας στην εκέση αυτή τα  $y(t^{n+k})$  με τα  $y^{n+k}$  παίρνουμε μια ενδιαφέρουσα κατηγορία μεθόδων, τις μεθόδους "οπισθοδρομικών διαφορών με k βήματα".

$$(6) \sum_{j=1}^k j^{-1} \nabla^j y^{n+k} = hf^{n+k}, \quad 0 \leq n \leq N-k,$$

όπου χρησιμοποιούμε τον συμβολισμό  $\nabla^1 y^n = y^n - y^{n-1}$ ,  $\nabla^2 y^n = \nabla^1(\nabla^1 y^n)$ , ... του λογισμού διαφορών. Οι μέθοδοι (6) σε κανονική μορφή (έτσι ώστε  $\alpha_k = 1$ ) γράφονται:

$$(6') \sum_{j=0}^k \alpha_j y^{n+j} = h\beta_k f^{n+k},$$

όπου

$$\text{για } k=1: \alpha_1=1, \alpha_0=-1, \beta_1=1 \text{ (παραλεχθέν Euler),}$$

$$\text{για } k=2: \alpha_2=1, \alpha_1=-4/3, \alpha_0=1/3, \beta_2=2/3,$$

$$\text{για } k=3: \alpha_3=1, \alpha_2=-18/11, \alpha_1=9/11, \alpha_0=-2/11, \beta_3=6/11$$

Γιά μία ευστηρατική μελέτη των τρόπων κατασκευής πολυβηματικών μεθόδων παραπέμπουμε στο βιβλίο του Henrici [3.4]. Συνήθως,  $k$ -βηματικές μέθοδοι της μορφής

$$(7) \quad y^{n+k} - y^{n+k-1} = h \sum_{j=0}^k \beta_j f^{n+j}$$

λέγονται μέθοδοι Adams. Ειδικότερα, αν  $\beta_k = 0$  (άμεσες), είναι γνωστές ως μέθοδοι Adams-Bashforth (1883). Αν  $\beta_k \neq 0$  παίρνουμε τις μεθόδους Adams-Moulton. Μέθοδοι της μορφής

$$(8) \quad y^{n+k} - y^{n+k-2} = h \sum_{j=0}^k \beta_j f^{n+j}$$

λέγονται συνήθως μέθοδοι του Nyström αν  $\beta_k = 0$  και μέθοδοι των Milne-Simpson αν  $\beta_k \neq 0$ .

Προχωρούμε τώρα στην ανάλυση των πολυβηματικών μεθόδων. Για να απλοποιήσουμε τα αποτελέσματα και τον συμβολισμό περιορίζομαστε στην περίπτωση μιας Δ.Ε., δηλ. θεωρούμε το πρόβλημα (3.1.1) στον  $\mathbb{R}^1$ , ακολουθώντας την ανάλυση του Henrici [3.4, Κεφ.5]. (Για ευετήματα ισχύουν εντελώς ανάλογα: βλ. π.χ., το βιβλίο του Gear [3.2, Κεφ.10].) Ανακαλούμε πρώτα ορισμένα αποτελέσματα από την θεωρία των εξισώσεων διαφορών με σταθερούς συντελεστές, εντελώς ανάλογα παρόμοιων αποτελεσμάτων για διαφορικές εξισώσεις ανωτέρας τάξης με σταθερούς συντελεστές. Θεωρούμε την ομογενή εξίσωση διαφορών

$$(9) \quad \alpha_k y^{n+k} + \alpha_{k-1} y^{n+k-1} + \dots + \alpha_0 y^n = 0, \quad n \geq 0,$$

όπου  $\alpha_j$ ,  $0 \leq j \leq k$  σταθερές με  $\alpha_k, \alpha_0 \neq 0$ . Ζητάμε να βρούμε λύσεις  $y^n$  της



(9) για  $n \geq 0$ , δηλ. γενικά μιγαδικές ακολουθίες  $\{y^n\}$ ,  $n \geq 0$  που κάνουν την (9) ταυτότητα. Υποθέτουμε ότι υπάρχουν λύσεις της μορφής  $y^n = z^n$  ( $n$ -ετή δύναμη του  $z$ ), όπου  $z \in \mathbb{C}$ ,  $z \neq 0$  βλέπουμε ότι πρέπει να ισχύει

$$\alpha_k z^{n+k} + \dots + \alpha_0 z^n = 0, \quad n \geq 0,$$

δηλ. ότι το  $z$  πρέπει να είναι ρίζα του πολυώνυμου

$$(10) \quad p(z) \equiv \alpha_k z^k + \alpha_{k-1} z^{k-1} + \dots + \alpha_0.$$

Αντίστροφα, κάθε ρίζα  $z$  του (10) ορίζει την λύση  $y^n = z^n$  της (9). Ισχύει δε το εξής βασικό αποτέλεσμα, για την (εύκολη) απόδειξη του οποίου βλ. π.χ. [3.4, σελ. 213-4]: Αν το πολυώνυμο  $p(z)$  έχει  $k$  διακριτές ρίζες  $z_1, \dots, z_k$ , τότε η γενική λύση της ομογενούς εξίσωσης διαφορών (9) δίδεται από τον γραμμικό συνδυασμό

$$(11) \quad y^n = c_1 z_1^n + c_2 z_2^n + \dots + c_k z_k^n$$

για οποιεσδήποτε σταθερές  $c_j \in \mathbb{C}$  (οι οποίες μπορούν να προσδιορισθούν μονοσήμαντα π.χ. αν δίνονται οι αρχικές τιμές  $y^j$ ,  $0 \leq j \leq k-1$  οδηγούμαστε έτσι σε ειδικές λύσεις της (9)). Αν το πολυώνυμο  $p(z)$  έχει  $m \leq k$  διακριτές ρίζες  $z_1, \dots, z_m$  με πολλαπλότητες  $p_1, \dots, p_m$  αντίστοιχα ( $p_1 + \dots + p_m = k$ ), τότε η γενική λύση της (9) δίνεται από του γραμμικό συνδυασμό

$$(12) \quad y^n = c_1 \zeta_1(n) + \dots + c_k \zeta_k(n), \quad c_i \in \mathbb{C}$$

όπου

$$\zeta_1(n) = z_1^n$$

$$\zeta_2(n) = nz_1^n$$

⋮

⋮

⋮

$$\zeta_{p_1}(n) = n(n-1)\dots(n-p_1+2)z_1^n$$

$$\zeta_{p_1+1}(n) = z_2^n$$

⋮

⋮

⋮

$$\zeta_{p_1+p_2}(n) = n(n-1)\dots(n-p_2+2)z_2^n$$

⋮

⋮

⋮

$$\zeta_k(n) = n(n-1)\dots(n-p_m+2)z_m^n.$$

Σημαντικό ρόλο στην ανάλυση της εύγκλισης των πολυημετικών μεθόδων παίζουν οι έννοιες της ευετάρθειας και της ευνέπειας. Το ξεκαθάρισμα των εννοιών αυτών και τα θεωρήματα που τις συνδέουν οφείλονται κυρίως στον Dahlquist (1956). Λέμε ότι η μέθοδος (3) είναι ευχκλίνουσα (ή ότι ευχκλίνει) αν για κάθε πραγματική συνάρτηση  $f(t,y)$  - που ικανοποιεί τις συνθήκες του θεωρήματος 3.1.1 στον  $\mathbb{R}^1$  - και κάθε  $y_0 \in \mathbb{R}^1$ , συμβολίζοντας με  $y(t)$  την λύση του προβλήματος (3.1.1) στον  $\mathbb{R}^1$ , έχουμε για κάθε  $t \in [a,b]$  ότι ισχύει

## 3.3.8

$$(13) \quad \lim_{h \rightarrow 0} y^n = y(t) \text{ όταν } h \rightarrow 0, n \rightarrow \infty, t^n = \alpha + nh \rightarrow t,$$

για κάθε λύση  $y^n$  της (3) που έχει αρχικές τιμές  $y^0, \dots, y^{k-1}$  τέτοιες ώστε

$$(14) \quad \lim_{h \rightarrow 0} y^j = y_0^j, \quad 0 \leq j \leq k-1.$$

Θα δούμε στην συνέχεια δύο σημαντικές συνέπειες του να είναι μια πολυβηματική μέθοδος ευκρίνους. Η πρώτη είναι η ευεταθεία της. Λέμε ότι η μέθοδος (3) είναι ευεταθείς αν οι ρίζες  $z_i$  του πολυωνύμου  $p(z)$  που δίνεται από την (10) ικανοποιούν την λεγόμενη συνθήκη των ριζών:

$$(15) \quad |z_i| \leq 1 \quad \forall i \quad \text{και} \quad |z_i| < 1 \text{ αν η } z_i \text{ είναι πολλαπλή ρίζα.}$$

**ΠΡΟΤΑΣΗ 1 (Dahlquist).** Αν η μέθοδος (3) ευκρίνεται, τότε είναι ευεταθείς.

Απόδειξη: Έστω ότι η μέθοδος (3) ευκρίνεται. Τότε θα ευκρίνεται και για το πρόβλημα αρχικών τιμών  $y' = 0, y(0) = 0$  του οποίου η λύση είναι  $y(t) \equiv 0$ . Για το πρόβλημα αυτό η (3) έχει την μορφή

$$(16) \quad \alpha_k y^{n+k} + \dots + \alpha_0 y^n = 0, \quad n \geq 0.$$

Για κάθε λύση της (16) λοιπόν τέτοια ώστε

$$(17) \quad \lim_{h \rightarrow 0} y^j = 0, \quad 0 \leq j \leq k-1,$$

θα ισχύει, για κάθε  $t \geq 0$ , ότι

$$(18) \quad \lim_{n \rightarrow \infty} y^n = 0, \quad h = t/n.$$

Εστω  $\zeta = re^{i\varphi}$ ,  $r \geq 0$ ,  $0 \leq \varphi < 2\pi$ , μία ρίζα του πολυωνύμου  $p(z)$ ,  
(10). θεωρείστε τους αριθμούς

$$(19) \quad y^n = h \operatorname{Re}(\zeta^n) = hr^n \cos n\varphi, \quad n=0,1,2,\dots,$$

οι οποίοι, επειδή  $\alpha_j \in \mathbb{R}^1$ , ικανοποιούν την (16) για  $n \geq 0$  με βάση τα όσα αναφέρθηκαν πιο πάνω για την εξίσωση διαφορών (9). Επίσης είναι προφανές ότι ικανοποιούν την (17). Συνεπώς για την ακολουθία (19) πρέπει να ισχύει η (18) για κάθε  $t > 0$ . Αν  $\varphi=0$  ή  $\varphi=\pi$  τότε από την (18) έχουμε αναγκαστικά ότι  $r \leq 1$ . Αν τώρα  $\varphi \neq 0$  και  $\varphi \neq \pi$ , παρατηρούμε ότι

$$\begin{aligned} (y^n)^2 - y^{n+1}y^{n-1} &= h^2 r^{2n} \cos^2 n\varphi - h^2 r^{2n} \cos(n+1)\varphi \cos(n-1)\varphi \\ &= h^2 r^{2n} \sin^2 \varphi, \end{aligned}$$

δηλ. ότι

$$[(y^n)^2 - y^{n+1}y^{n-1}] / \sin^2 \varphi = h^2 r^{2n}.$$

Το αριστερό μέλος αυτής της ισότητας τείνει στο μηδέν λόγω της (18) όταν  $n \rightarrow \infty$ . Άρα για κάθε  $t > 0$   $h^2 r^{2n} = t^2 (r^n/n)^2 \rightarrow 0$ ,  $n \rightarrow \infty$ . Αναγκαστικά λοιπού  $r \leq 1$ . Συμπεραίνουμε ότι "σύγκλιση της (3)"  $\Rightarrow$  "κάθε ρίζα  $\zeta$  του  $p(z)$  ικανοποιεί  $|\zeta| \leq 1$ ".

Εστω τώρα  $\zeta = re^{i\varphi}$  μια ρίζα του  $p(z)$  με πολλαπλότητα μεγαλύτερη της μονάδας. Από την θεωρία των λύσεων της ομογενούς εξίσωσης διαφορών (9) ε' αυτήν την περίπτωση, συμπεραίνουμε ότι η ακολουθία

$$(20) \quad y^n = h^{1/2} nr^n \cos n\varphi, \quad n=0,1,2,\dots$$

αποτελεί λύση της (16) που ικανοποιεί την (17). Συνεπώς, θα ικανοποιεί την (18) για κάθε  $t > 0$ . Αν  $\varphi = 0$  ή  $\varphi = \pi$  συμπεραίνουμε ότι  $|y^n| = h^{1/2} n r^n = t^{1/2} n^{1/2} r^n \rightarrow 0$ ,  $n \rightarrow \infty$  για  $t > 0$ . Συνεπώς  $r < 1$ . Αν  $\varphi \neq 0$  ή  $\varphi \neq \pi$ , εύκολα βλέπουμε ότι αν  $z^n = n^{-1} h^{-1/2} y^n$ , τότε

$$(21) \quad ((z^n)^2 - z^{n+1} z^{n-1}) / \sin^2 \varphi = r^{2n}.$$

Επειδή  $z^n = y^n / n h^{1/2} = y^n / (n t)^{1/2} \rightarrow 0 \quad \forall t > 0$  όταν  $n \rightarrow \infty$  (λόγω της (18)) η σύγκλιση στο 0 του αριστερού μέλους της (21) για  $n \rightarrow \infty$  δίνει  $r < 1$  πάλι. Συνεπώς κάθε πολλαπλή ρίζα  $\xi$  του  $p(z)$  πρέπει να βρίσκεται στον ανοικτό μοναδιαίο δίσκο, δηλ. να ικανοποιεί  $|\xi| < 1$ . @

Εξετάζοντας τις ρίζες του αντίστοιχου  $p(z)$  βλέπουμε ότι οι μέθοδοι Euler, πεπλεγμένη Euler και τραπεζίου (για τις οποίες  $p(z) = z-1$ ), η μέθοδος (1') και η μέθοδος Simpson (2) (για τις οποίες  $p(z) = z^2-1$ ) ικανοποιούν την συνθήκη των ριζών (15) και συνεπώς είναι ευσταθείς. (Ευσταθείς είναι γενικά οι μέθοδοι Adams (7) για τις οποίες  $p(z) = z^{k-1}(z-1)$  και οι μέθοδοι (8) για τις οποίες  $p(z) = z^{k-2}(z^2-1)$ ). Έχει αποδειχθεί (Cryer) ότι οι μέθοδοι οπισθοδρομικών διαφορών (6) είναι ευσταθείς αν και μόνο αν  $1 \leq k \leq 6$  (βλ. και Άσκηση 8).

Αν μία μέθοδος δεν είναι ευσταθής, τότε και για την απλούστατη Δ.Ε.  $y' = 0$ ,  $y(0) = 0$  με  $y^j = 0$ ,  $0 \leq j \leq k-1$ , λόγω εφαλμάτων ετραγχύλευσης, η λύση δεν θα είναι  $y^n = 0$ , πλθ. αλλά θα έχει γενικά μια συνιστώσα - βλ. (11), (12) - που δεν θα μένει φραχμένη καθώς αυξάνει το  $n$ . Στην πράξη αν υπολογίσουμε με μια ασταθή μέθοδο βλέπουμε πολύ γρήγορα μια "έκρηξη" της αριθμητικής λύσης  $y^n$  καθώς αυξάνει το  $n$ . Η χρήση μικρότερων  $h$  δεν βελτιώνει (ατίθεται χειροτερεύει) την κατάσταση - βλ. [5.3, σελ 188] για αριθμητικά παραδείγματα.

Θα δούμε αργότερα (Παρ. 3.4) ότι για πολλά ευδαφέρουτα προβλήματα η έννοια της "ευσταθείας" που εισαγάγαμε δεν είναι

επαρκής' δηλ. ότι οι μέθοδοι που ικανοποιούν την συνθήκη των ριζών (15) μπορεί να μην ευπεριφέρονται καθόλου ικανοποιητικά στην πράξη σε περιβάλλον εφαγμάτων ετρογχύλευσης. Προς το παρόν όμως στρεφόμαστε στην μελέτη της τάξης ακρίβειας και της ευνέπειας των πολυβηματικών μεθόδων.

Δεδομένης της  $k$ -βηματικής μεθόδου (3) (δηλ. των σταθερών  $\{\alpha_j, \beta_j\}$ ,  $0 \leq j \leq k$ ), θεωρούμε για  $a \leq t \leq b$  και για μια αρκετά ομαλή συνάρτηση  $y(t)$  την ποσότητα

$$(22) \quad L[y(t); h] = \sum_{j=0}^k \{\alpha_j y(t+jh) - h\beta_j y'(t+jh)\},$$

που παίζει τον ρόλο της ποσότητας  $y(t+h) - y(t) - h\Phi(t, y(t), h)$  των μεθόδων RK. Αναπτύσσοντας το δεύτερο μέλος της (22) σε σειρά Taylor γύρω απ' το σημείο  $t$  έχουμε

$$(23) \quad L[y(t); h] = C_0 y(t) + C_1 h y'(t) + C_2 h^2 y''(t) + \dots + C_j h^j y^{(j)}(t) + \dots,$$

όπου οι σταθερές  $C_j$  είναι ανεξάρτητες των  $h, t, y(t)$  και εξαρτώνται μόνο από την μέθοδο (3). Πάλιετα εύκολα αποδεικνύεται ότι

$$(24) \quad C_0 = \alpha_0 + \alpha_1 + \dots + \alpha_k = \sum_{j=0}^k \alpha_j,$$

$$C_1 = \alpha_1 + 2\alpha_2 + \dots + k\alpha_k - (\beta_0 + \dots + \beta_k) = \sum_{j=1}^k j\alpha_j - \sum_{j=0}^k \beta_j,$$

και

$$C_j = (\alpha_1 + 2^j \alpha_2 + 3^j \alpha_3 + \dots + k^j \alpha_k) / j! - (\beta_1 + 2^{j-1} \beta_2 + 3^{j-1} \beta_3 + \dots + k^{j-1} \beta_k) / (j-1)! \quad \text{για } j \geq 2.$$

Λέμε ότι η μέθοδος (3) έχει τάξη ακρίβειας  $p$  αν στο ανάπτυγμα (23),  $C_0 = C_1 = \dots = C_p = 0$  αλλά  $C_{p+1} \neq 0$ , δηλ. αν  $L[y(t); h] = C_{p+1} h^{p+1} y^{(p+1)}(t) + \dots$ . Εύκολα βλέπουμε ότι οι τάξεις των μεθόδων Euler, πεπλεγμένης Euler και τραπεζίου είναι αντίστοιχα  $p=1$ ,

1 και 2. Η τάξη της μεθόδου (1') είναι  $p=2$  ενώ της μεθόδου του Simpson (2) είναι  $p=4$ . Οι μέθοδοι (6) έχουν τάξη ακρίβειας  $p=k$ .

Λέμε ότι η μέθοδος (3) είναι ευνεπής αν έχει τάξη (ακρίβειας) τουλάχιστον 1. Από τις (24) βλέπουμε ότι η (3) είναι ευνεπής αν  $C_0=C_1=0$  δηλ. αν

$$(25) \sum_{j=0}^k \alpha_j = 0, \quad \sum_{j=1}^k j \alpha_j = \sum_{j=0}^k \beta_j.$$

Εισάγοντας, εκτός από το πολυώνυμο  $p(z)$  (10), το πολυώνυμο  $\epsilon(z)$ :

$$(26) \epsilon(z) = \beta_k z^k + \beta_{k-1} z^{k-1} + \dots + \beta_0,$$

βλέπουμε ότι οι συνθήκες (25) της ευνεπείας μίας μεθόδου γράφονται ισοδύναμα ως

$$(25') p(1) = 0, \quad p'(1) = \epsilon(1).$$

**ΠΡΟΤΑΣΗ 2.** Αν η μέθοδος (3) ευκλίνει τότε είναι ευνεπής.

Απόδειξη: Έστω ότι η (3) ευκλίνει. Τότε θα ευκλίνει και για το πρόβλημα  $y'=0$ ,  $y(0)=1$  του οποίου η λύση είναι  $y(t)=1$ . Για το πρόβλημα αυτό η (3) πάλι έχει την μορφή (16). Επειδή η μέθοδος ευκλίνει, η λύση  $\{y^n\}$  της (16) που αντιστοιχεί στις αρχικές τιμές  $y^j=1$ ,  $0 \leq j \leq k-1$  (για την οποία δηλ. ισχύουν τετριμμένα οι (14)), πρέπει να ικανοποιεί την (13) δηλ. να ισχύει ότι για κάθε  $t > 0$

$$(26) \lim_{n \rightarrow \infty} y^n = 1, \quad n \rightarrow \infty, \quad h \rightarrow 0, \quad nh=t.$$

Λόγω όμως της (16) - που δεν εξαρτάται από το  $h$  - και των αρχικών τιμών  $y^j=1$ ,  $0 \leq j \leq k-1$  - που δεν εξαρτώνται από το  $h$  - η ακολουθία  $\{y^n\}$  δεν εξαρτάται από το  $h$  ή το  $t$ . Συνεπώς η (26) λέει απλώς ότι

$$(26') \lim_{n \rightarrow \infty} y^n = 1, \quad n \rightarrow \infty.$$

Άρα, αν το  $n$  τείνει στο  $\infty$  στην (16) παίρνουμε ότι  $\alpha_k + \alpha_{k-1} + \dots + \alpha_0 = 0$ , δηλ. ότι  $C_0 = 0$ .

Γιά να αποδείξουμε ότι  $C_1=0$  θεωρούμε το πρόβλημα αρχικών τιμών  $y'=1$ ,  $y(0)=0$  με λύση  $y(t)=t$ . Η (3) παίρνει τώρα την μορφή

$$(27) \quad \alpha_k y^{n+k} + \dots + \alpha_0 y^n = h(\beta_k + \dots + \beta_0).$$

Επειδή η (3) συγκλίνει έχουμε ότι κάθε λύση της (27) για την οποία

$$(28) \quad \lim_{h \rightarrow 0} y^j = 0, \quad 0 \leq j \leq k-1,$$

θα πρέπει να ικανοποιεί, για κάθε  $t > 0$ , την

$$(29) \quad \lim_{h \rightarrow 0} y^n = t, \quad nh=t.$$

θεωρούμε την ακολουθία  $y^n$ ,  $n \geq 0$  που ορίζεται ως

$$(30) \quad y^n = nhK, \quad K \equiv \varepsilon(1)/p'(1), \quad n \geq 0.$$

(Από την Πρόταση 1 λόγω της ευεστάθειας της μεθόδου έχουμε  $p'(1) \neq 0$ ). Προφανώς για την (30) ισχύει η (28). Έχουμε επίσης ότι

$$\sum_{j=0}^k \alpha_j y^{n+j} = hK \sum_{j=0}^k (n+j)\alpha_j = (\text{λόγω της } C_0=0) =$$

$$hK \sum_{j=0}^k j\alpha_j = hKp'(1) = h \sum_{j=0}^k \beta_j,$$

δηλ. ότι η (30) είναι λύση της (27). Συνεπώς πρέπει να ικανοποιεί την (29) η οποία μας δίνει ότι για  $t > 0$   $t = \lim_{h \rightarrow 0} nhK = \lim_{h \rightarrow 0} tK = tK$ ,

δηλ. ότι  $K=1 \Leftrightarrow p'(1)=\varepsilon(1) \Leftrightarrow C_1=0$ . @

Συνεπώς η ευεστία μιάς μεθόδου είναι αναγκαία για την σύγκλιση της. Στο [5.3, εελ. 189] βλέπουμε ένα τυπικό παράδειγμα υπολογισμού με μιά ευεσταθή αλλά μη ευεπή μέθοδο: Το εφάλμα μεγαλώνει καθώς  $h \rightarrow 0$  αλλά η ακολουθία  $y^n$  μένει φραχμένη· πλησιάζει δε την λύση κάποιου άλλου προβλήματος αρχικών τιμών!



Όπως θα δούμε παρακάτω, (πόρισμα του θεωρήματος 1), ευγένεια+ευστάθεια  $\Rightarrow$  εύγκλιση. Δηλ. η {ευγένεια και ευστάθεια} είναι ικανή και αναγκαία συνθήκη για εύγκλιση. Ας σημειώσουμε όμως το εξής σημαντικό αποτέλεσμα του Dahlquist (1956) το οποίο περιορίζει την δυνατή τάξη ακρίβειας μίας ευσταθούς μεθόδου. Η απόδειξη μπορεί να βρεθεί στο βιβλίο του Henrici [3.4, σελ. 229]:

**ΠΡΟΤΑΣΗ 3.** Η μέγιστη τάξη μίας ευσταθούς  $k$ -βηματικής μεθόδου της μορφής (3) είναι  $p=k+1$  αν  $k$  περιττός και  $p=k+2$  αν  $k$  άρτιος. @

Μέθοδοι βέλτιστης τάξης για δεδομένο αριθμό βημάτων  $k$  λέγονται λοιπόν οι μέθοδοι με τάξη  $p=k+1$  αν  $k$  περιττός και  $p=k+2$  αν  $k$  άρτιος. Ο προσδιορισμός τους παρουσιάζει θεωρητικό και πρακτικό προφανώς ενδιαφέρον (βλ. Ασκήσεις 3,6). Προφανώς η μέθοδος του τραπεζίου ( $k=1, p=2$ ) και η μέθοδος του Simpson ( $k=2, p=4$ ) είναι μέθοδοι βέλτιστης τάξης για  $k=1,2$  αντίστοιχα.

Προχωρούμε τώρα στην απόδειξη του βασικού αποτελέσματος αυτής της παραγράφου, δηλ. στην απόδειξη ενός φράγματος άριστης τάξης εύγκλισης για το σφάλμα της μεθόδου (3). Πρώτα δύο προκαταρκτικά αποτελέσματα:

**Λήμμα 1.** Έστω ότι το πολυώνυμο  $p(z) = \sum_{j=0}^k a_j z^j$  ικανοποιεί την συνθήκη των ριζών (15). Ορίζουμε τους συντελεστές  $\gamma_j, j \geq 0$ , από το ανάπτυγμα

$$(31) \quad 1/(a_k + a_{k-1}z + \dots + a_0 z^k) = \gamma_0 + \gamma_1 z + \gamma_2 z^2 + \dots$$

Τότε τα  $\gamma_j$  ικανοποιούν τις συνθήκες

$$(32) \quad \begin{aligned} a_k \gamma_0 &= 1, \\ a_k \gamma_j + a_{k-1} \gamma_{j-1} + \dots + a_{k-j} \gamma_0 &= 0 \text{ αν } 1 \leq j \leq k, \\ &\vdots \\ a_k \gamma_j + a_{k-1} \gamma_{j-1} + \dots + a_0 \gamma_{j-k} &= 0 \text{ αν } j > k. \end{aligned}$$

Επιπλέον,

$$(33) \quad \Gamma \equiv \sup_{j \geq 0} |\gamma_j| < \infty.$$

Απόδειξη: Έστω  $\hat{p}(z) = \alpha_k + \alpha_{k-1}z + \dots + \alpha_0 z^k = z^k p(z^{-1})$ . Επειδή, λόγω ευστάθειας, το πολυώνυμο  $p(z)$  δεν έχει ρίζες με  $|z| > 1$  και επειδή  $\alpha_k \neq 0$ , συμπεραίνουμε ότι το πολυώνυμο  $\hat{p}(z)$  δεν έχει ρίζες  $z$  τέτοιες ώστε  $|z| < 1$ . Συνεπώς η συνάρτηση  $1/\hat{p}(z)$  είναι αναλυτική στον ανοιχτό δίσκο  $D = \{z \in \mathbb{C} : |z| < 1\}$ . Θεωρούμε το ανάπτυγμα της (31) κατά Taylor γύρω απ' το μηδέν. Οι ταυτότητες (32) προκύπτουν εύκολα από την (31) στην μορφή

$$1 = (\alpha_k + \alpha_{k-1}z + \dots + \alpha_0 z^k)(\gamma_0 + \gamma_1 z + \gamma_2 z^2 + \dots),$$

μετά από πολλαπλασιασμό των σειρών του δευτέρου μέλους και εξίσωση δυνάμεων του  $z$  των δύο μελών. (Επειδή  $\alpha_k \neq 0$  οι συνθήκες (32) προσδιορίζουν μονοσήμαντα τα  $\gamma_j$ ,  $j \geq 0$  συναρτήσει των  $\alpha_j$ ,  $0 \leq j \leq k$ ).

Γιά να αποδείξουμε την εκτίμηση (33) για τους συντελεστές της σειράς Taylor του  $(\hat{p}(z))^{-1}$ , παρατηρούμε πρώτα ότι λόγω της (15) οι ρίζες  $z_1, \dots, z_m$  του  $p(z)$  που ικανοποιούν  $|z_j| = 1$  είναι απλές (αν καν υπάρχουν). Συμπεραίνουμε ότι οι μόνοι πόλοι της συνάρτησης  $(\hat{p}(z))^{-1}$  στην περιφέρεια  $|z| = 1$  είναι οι απλοί πόλοι στα σημεία  $z_j^{-1}$ ,  $1 \leq j \leq m$ . Άρα, από γνωστό μας αποτέλεσμα της θεωρίας μιγαδικών συναρτήσεων, υπάρχουν σταθερές  $A_j$ ,  $1 \leq j \leq m$  τέτοιες ώστε η συνάρτηση

$$(34) \quad f(z) = (\hat{p}(z))^{-1} - A_1/(z - z_1^{-1}) - \dots - A_m/(z - z_m^{-1}),$$

να είναι αναλυτική για  $|z| \leq 1$ . Θεωρούμε την σειρά Taylor της  $f(z)$  περί το μηδέν:

$$f(z) = \sum_{n=0}^{\infty} z^n f^{(n)}(0)/n!$$

Από το θεώρημα του Cauchy αν  $C = \{z : |z| = 1\}$  έχουμε

$$f^{(n)}(0) = n! (2\pi i)^{-1} \int_C (f(z)/z^{n+1}) dz,$$

από την οποία συμπεραίνουμε ότι  $|f^{(n)}(0)|/n! \leq \sup_{z \in C} |f(z)| = M < \infty$ , δηλ. ότι οι συντελεστές  $f^{(n)}(0)/n!$  της σειράς Taylor της  $f$  περί το μηδέν είναι ομοιόμορφα φραγμένοι από κάποιο  $M < \infty$ . Επειδή εξ άλλου για  $|z| < 1$  επειδή  $|z_i| = 1$

$$-A_i/(z-z_i^{-1}) = A_i z_i (1+z_i z+(z_i z)^2+\dots+(z_i z)^n+\dots),$$

βλέπουμε ότι οι συντελεστές της σειράς Taylor της συνάρτησης  $-A_i/(z-z_i^{-1})$  περί το  $z=0$  φράσσονται από  $|A_i|$ . Συμπεραίνουμε από την (34) ότι οι συντελεστές  $\gamma_j$  της σειράς Taylor της  $1/\hat{p}(z)$  περί το 0 φράσσονται ομοιόμορφα:

$$\Gamma = \sup_{j \geq 0} |\gamma_j| \leq M + \sum_{i=1}^m |A_i| < \infty @$$

Θα χρησιμοποιήσουμε το παραπάνω λήμμα για να εκτιμήσουμε α priori τις λύσεις μιάς μη ομογενούς εξίσωσης διαφορών που εχειζεται με την πολυβηματική μέθοδο (3):

**Λήμμα 2.** Έστω ότι η  $k$ -βηματική μέθοδος (3) είναι ευσταθής. Έστω  $\lambda^n$ ,  $0 \leq n \leq N-k$  δεδομένες εταθερές και έστω  $\beta_i^n$ ,  $0 \leq i \leq k$ ,  $0 \leq n \leq N-k$  δεδομένοι αριθμοί, τέτοιοι ώστε  $|\beta_i^n| \leq B < \infty \quad \forall i, n$ . Θεωρείστε την εξίσωση διαφορών

$$(35) \quad \alpha_k \psi^{n+k} + \alpha_{k-1} \psi^{n+k-1} + \dots + \alpha_0 \psi^n = \\ = h(\beta_k^n \psi^{n+k} + \beta_{k-1}^n \psi^{n+k-1} + \dots + \beta_0^n \psi^n) + \lambda^n, \quad 0 \leq n \leq N-k.$$

Τότε υπάρχει  $h_0 > 0$  τέτοιο ώστε για  $0 \leq h \leq h_0$ :

$$(36) \quad \max_{0 \leq n \leq N} |\psi^n| \leq C(N \max_n |\lambda^n| + \max_{0 \leq j \leq k-1} |\psi^j|),$$

όπου η σταθερά  $C$  εξαρτάται από τα  $b-a$ ,  $h_0$ ,  $\theta$  και την μέθοδο (3) αλλά όχι από τα  $h, \lambda^n, \psi^n, N$ .

Απόδειξη: Για  $k \leq m$  θεωρούμε την (35) για  $n=0, 1, 2, \dots, m-k$ .

[Για  $n=m-k-j$ ,  $0 \leq j \leq m-k$  πολλαπλασιάζουμε την (35) επί  $\delta_j$  (που ορίστηκε στο λήμμα 1) και προσθέτουμε κατά μέλη τις εξισώσεις που προκύπτουν. Από το αριστερό μέλος παίρνουμε το άθροισμα

$$(37) \quad \begin{aligned} S_m &\equiv \delta_0 (\alpha_k \psi^m + \alpha_{k-1} \psi^{m-1} + \dots + \alpha_0 \psi^{m-k}) + \delta_1 (\alpha_k \psi^{m-1} + \alpha_{k-1} \psi^{m-2} \\ &\quad + \dots + \alpha_0 \psi^{m-k-1}) + \dots + \delta_{m-k} (\alpha_k \psi^k + \alpha_{k-1} \psi^{k-1} + \dots + \alpha_0 \psi^0) = \\ &\quad \alpha_k \delta_0 \psi^m + (\alpha_k \delta_1 + \alpha_{k-1} \delta_0) \psi^{m-1} + \dots + (\alpha_k \delta_{m-k} + \\ &\quad \alpha_{k-1} \delta_{m-k-1} + \dots + \alpha_0 \delta_{m-2k}) \psi^k + (\alpha_{k-1} \delta_{m-k} + \dots + \alpha_0 \delta_{m-2k+1}) \psi^{k-1} \\ &\quad + \dots + \alpha_0 \delta_{m-k} \psi^0 = (\text{χρησιμοποιώντας τις ταυτότητες (32)}) \\ &\quad = \psi^m + (\alpha_{k-1} \delta_{m-k} + \dots + \alpha_0 \delta_{m-2k+1}) \psi^{k-1} + \dots + \alpha_0 \delta_{m-k} \psi^0. \end{aligned}$$

Από το δεξιό μέλος παίρνουμε μετά από λίγες ανακατατάξεις όρων

$$(38) \quad \begin{aligned} S_m &= h \{ (\beta_k^{m-k} \delta_0 \psi^m + (\beta_{k-1}^{m-k} \delta_0 + \beta_k^{m-k-1} \delta_1) \psi^{m-1} + \dots + \\ &\quad (\beta_0^{m-k} \delta_0 + \dots + \beta_k^{m-2k} \delta_k) \psi^{m-k} + \dots + \beta_0^0 \delta_{m-k} \psi^0) \\ &\quad + \lambda^{m-k} \delta_0 + \lambda^{m-k-1} \delta_1 + \dots + \lambda^0 \delta_{m-k} \}. \end{aligned}$$

Εξισώνοντας τις (37) και (38) και λύοντας ως προς  $\psi^m$  (χρησιμοποιώντας ότι  $\delta_0 = 1/\alpha_k = 1$  από την (32)) έχουμε:

$$\begin{aligned} (1 - h\beta_k^{m-k}) \psi^m &= h \{ (\beta_{k-1}^{m-k} \delta_0 + \beta_k^{m-k-1} \delta_1) \psi^{m-1} + \dots + \beta_0^0 \delta_{m-k} \psi^0 \} \\ &\quad - \{ (\alpha_{k-1} \delta_{m-k} + \dots + \alpha_0 \delta_{m-2k+1}) \psi^{k-1} + \dots + \alpha_0 \delta_{m-k} \psi^0 \} \\ &\quad + \lambda^{m-k} \delta_0 + \lambda^{m-k-1} \delta_1 + \dots + \lambda^0 \delta_{m-k}, \end{aligned}$$

από την οποία, χρησιμοποιώντας την (33) παίρνουμε

$$|1 - h\alpha_k^{m-k}| |\psi^m| \leq C_1 h \sum_{j=0}^{m-1} |\psi^j| + C_2 \sum_{j=0}^{k-1} |\psi^j| + \Gamma H \max_{0 \leq j \leq m-k} |\lambda^j|,$$

όπου  $C_1 = (k+1)B\Gamma$ ,  $C_2 = \Gamma \sum_{j=0}^k |\alpha_j|$ . Συνεπώς για  $h \leq h_0$  όπου  $h_0 < B^{-1}$ , συμπεραίνουμε ότι υπάρχει σταθερά  $\theta$  τ.ω.

$$(39) \quad |\psi^m| \leq C' (h \sum_{j=0}^{m-1} |\psi^j| + H \max_j |\lambda^j| + \sum_{j=0}^{k-1} |\psi^j|), \quad k \leq m \leq N.$$

θέτουμε  $A = C' (H \max_j |\lambda^j| + \sum_{j=0}^{k-1} |\psi^j|)$ . Τότε ισχύει (τετριμμένα) ότι

$$(40) \quad |\psi^j| \leq A(1+hC')^j, \quad 0 \leq j \leq k-1.$$

θα δείξουμε επαγωγικά ότι η (40) ισχύει για  $0 \leq j \leq N$ . Πράγματι, έστω ότι έχουμε

$$(41) \quad |\psi^j| \leq A(1+hC')^j, \quad 0 \leq j \leq m-1.$$

Η (39) και η (41) δίνουν τότε

$$\begin{aligned} |\psi^m| &\leq C' h A \sum_{j=0}^{m-1} (1+hC')^j + A = A(C'h \sum_{j=0}^{m-1} (1+hC')^j + 1) \\ &= A(C'h \{(1+hC')^m - 1\} / C'h + 1) = A(1+hC')^m. \end{aligned}$$

Άρα ισχύει η (41) για  $j=m$  το επαγωγικό βήμα τελείωσε. Συνεπώς για κάθε  $m$ ,  $0 \leq m \leq N$ , έχουμε

$$|\psi^m| \leq A(1+hC')^m \leq A e^{hmC'} \leq A e^{(b-a)C'}.$$

Άρα ισχύει η (36) με  $C = e^{(b-a)C'}$ . @

Ερχόμαστε τώρα στο κεντρικό αποτέλεσμα αυτής της παραγράφου:

**ΘΕΩΡΗΜΑ 1.** Υποθέτουμε ότι η  $k$ -βηματική μέθοδος (3) είναι ευстаθής και έχει τάξη ακρίβειας  $p \geq 1$ . Έστω ότι η λύση  $y$  του προβλήματος αρχικών τιμών (για μία απλή Δ.Ε.) (3.1.1) ανήκει στον χώρο  $C^{p+1}[a, b]$ . Υπάρχει τότε  $h_0 > 0$  τέτοιο ώστε για  $0 \leq h \leq h_0$  να ισχύει η εκτίμηση

$$(42) \max_{0 \leq n \leq N} |y^n - y(t^n)| \leq C (\max_{0 \leq j \leq k-1} |y^j - y(t^j)| + h^p \max_{a \leq t \leq b} |y^{(p+1)}(t)|),$$

όπου  $C$  σταθερά ανεξάρτητη των  $y^n$ ,  $y(t)$ ,  $h, N$ .

Απόδειξη: Από την υπόθεσή μας ότι η (3) έχει τάξη  $p$ , συμπεραίνουμε ότι για τις ποσότητες

$$(43) p^n = \sum_{j=0}^k [\alpha_j y(t^{n+jh}) - h \beta_j y'(t^{n+jh})], \quad 0 \leq n \leq N-k,$$

ισχύει, από το θεώρημα του Taylor με υπόλοιπο, ότι για κάποια σταθερά  $C'$ :

$$(44) \max_{0 \leq n \leq N-k} |p^n| \leq C' h^{p+1} \max_{a \leq t \leq b} |y^{(p+1)}(t)|.$$

Θεωρούμε τώρα το σφάλμα  $e^n = y^n - y(t^n)$ ,  $0 \leq n \leq N$ , το οποίο ικανοποιεί την εξίσωση, για  $0 \leq n \leq N-k$ ,

$$(45) \alpha_k e^{n+k} + \alpha_{k-1} e^{n+k-1} + \dots + \alpha_0 e^n = \\ = (\alpha_k y^{n+k} + \dots + \alpha_0 y^n) - (\alpha_k y(t^{n+k}) + \dots + \alpha_0 y(t^n)) =$$

$$= (από τις (3), (43)), = h(\beta_k f^{n+k} + \dots + \beta_0 f^n) - h[\beta_k f(t^{n+k}, y(t^{n+k})) + \dots$$

$$\dots + \beta_0 f(t^n, y(t^n))] - p^n =$$

$$= h\{\beta_k [f^{n+k} - f(t^{n+k}, y(t^{n+k}))] + \dots + \beta_0 [f^n - f(t^n, y(t^n))]\} - p^n.$$

Ορίζουμε τώρα  $g^n$ ,  $0 \leq n \leq N$  ως

$$(46) \quad g^n = \begin{cases} [f^n - f(t^n, y(t^n))]/\varepsilon^n & \text{αν } \varepsilon^n \neq 0 \\ 0 & \text{αν } \varepsilon^n = 0, \end{cases}$$

οπότε η (45) γράφεται

$$(47) \quad \alpha_k \varepsilon^{n+k} + \dots + \alpha_0 \varepsilon^n = h \{ \beta_k g^{n+k} \varepsilon^{n+k} + \dots + \beta_0 g^n \varepsilon^n \} - p^n, \quad 0 \leq n \leq N-k.$$

Η (46) δίνει τώρα για  $\varepsilon^n \neq 0$ , με χρήση της συνθήκης Lipschitz της  $f$ ,

$$(48) \quad |g^n| \leq L |y^n - y(t^n)| / |\varepsilon^n| = L, \quad 0 \leq n \leq N.$$

που ισχύει βέβαια και για  $\varepsilon^n = 0$ . Ορίζουμε τους αριθμούς  $\beta_i^n$ ,  $0 \leq i \leq k$ ,  $0 \leq n \leq N-k$ , από τις σχέσεις  $\beta_i^n = \beta_i g^{n+i}$ ,  $0 \leq i \leq k$ ,  $0 \leq n \leq N-k$ . Λόγω της (48) έχουμε

$$(49) \quad \max_{i,n} |\beta_i^n| = \max_i |\beta_i| \max_n |g^n| \leq \max_i |\beta_i| \cdot L \equiv B < \infty$$

Η (47) γράφεται λοιπόν ως

$$\alpha_k \varepsilon^{n+k} + \dots + \alpha_0 \varepsilon^n = h \{ \beta_k \varepsilon^{n+k} + \dots + \beta_0 \varepsilon^n \} - p^n, \quad 0 \leq n \leq N-k,$$

δηλ. ως μία εξίσωση διαφορών για τα  $\varepsilon^n$  με μεταβλητούς συντελεστές της μορφής (35) του Λήμματος 2. Εφαρμόζοντας τώρα την (36) παίρνοντας υπ' όψιν την (44) και ότι  $h = (b-a)$ , καταλήγουμε στην (42). @

Το θεώρημα 1 αποτελεί εχεβόλ το αντίστροφο των προτάσεων 1 και 2' μας διαβεβαιώνει ότι για τουλάχιστον  $y \in C^2[a, b]$ , "ευεστία και συνέπεια"  $\Rightarrow$  "εύγκλιση". (Μας δίνει βέβαια και εκτίμηση του εφάλματος

και γιαυτό είναι πολύ χρήσιμο). Μπορεί να αποδειχθεί, βλ. Henrici [3.4, Παρ. 5.3-3] ότι "ευεταθία και ευγένεια" => "εύγκλιση", χωρίς να υποθέσουμε ότι  $y \in C^2[a, b]$ , δηλ. μόνο με τις συνθήκες στην  $f$  του θεωρήματος 3.1.1 που εγγυώνται  $y \in C^1[a, b]$ . Το αποτέλεσμα αυτό και οι προτάσεις 1 και 2 ενοψίζονται λοιπόν με την διατύπωση "ευεταθία και ευγένεια"  $\Leftrightarrow$  "εύγκλιση".

Είναι φανερό από την (42) ότι για να πάρουμε συνολικό εφάλμα  $O(h^p)$  αρκεί  $|y^j - y(t^j)| = O(h^p)$ ,  $0 \leq j \leq k-1$ . Αρκεί δηλ. οι τιμές  $y^j$ ,  $0 \leq j \leq k-1$ , ( $y^0 = y_0$ ) να υπολογισθούν με μία μέθοδο RK τάξης  $p-1$  (της οποίας το τοπικό εφάλμα, όσας  $O(h^p)$ , εξασφαλίζει ότι  $|y^j - y(t^j)| = O(h^p)$  για μικρό αριθμό βημάτων  $0 \leq j \leq k-1$ .)

#### Παρατηρήσεις

1. Η συνθήκη των ριζών (15) είναι βέβαια εύκολο να ελεγχθεί π.χ. για δευτεροβάθμια πολυώνυμα  $p(z)$ . Για πολυώνυμα βαθμού  $k > 2$  θα ήταν χρήσιμο να βρούμε αναλυτικές συνθήκες πάνω στους συντελεστές του  $p(z)$  για την ιεχύ της. Μία χρήσιμη εχετική θεωρία είναι η λεγόμενη Θεωρία Schur. Λέμε (βλ. π.χ. Miller, J. Inst. Math. Applies, 8 (1971), 397-406) ότι ένα πολυώνυμο με μιγαδικούς συντελεστές  $a_j$ , βαθμού  $k$ ,

$$(50) \quad p(z) = a_0 + a_1 z + \dots + a_k z^k,$$

όπου  $a_0 \neq 0$ ,  $a_k \neq 0$  είναι (πολυώνυμο) Schur αν όλες οι ρίζες του περιέχονται στον ανοικτό μοναδιαίο δίσκο  $D = \{z \in \mathbb{C} : |z| < 1\}$ . Λέμε ότι το  $p$  είναι απλό (πολυώνυμο) von Neumann αν όλες οι ρίζες του βρίσκονται στον κλειστό δίσκο  $\bar{D}$  και αν μόνο απλές ρίζες του βρίσκονται πάνω στην περιφέρεια  $|z|=1$ . Π' άλλα λόγια η συνθήκη των ριζών (15) είναι ισοδύναμη με την συνθήκη να είναι το  $p(z)$  απλό von Neumann.

Δεδομένου του  $p$  θεωρούμε το πολυώνυμο

$$(51) \quad p^*(z) \equiv \sum_{j=0}^k \bar{a}_{k-j} z^j = \bar{a}_k + \bar{a}_{k-1} z + \dots + \bar{a}_0 z^k,$$



όπου  $\bar{z}$  ο συζυγής του  $z$ . Προφανώς έχουμε

$$p^*(z) = z^k \bar{p}(z^{-1}).$$

Επίσης θεωρούμε το λεγόμενο "ανηγμένο πολυώνυμο"

$$(52) \quad p_1(z) = (p^*(0)p(z) - p(0)p^*(z))/z,$$

βαθμού  $\leq k-1$ , και συμβολίζουμε με  $p'(z)$  την παράγωγο του  $p$ . Τότε ισχύουν τα εξής:

- (i) Το πολυώνυμο  $p$  είναι Schur αν και μόνο αν  $|p^*(0)| > |p(0)|$  και το  $p_1$  είναι Schur.
- (ii) Το πολυώνυμο  $p$  είναι απλό von Neumann αν και μόνο αν είτε
  - (α)  $|p^*(0)| > |p(0)|$  και το  $p_1$  είναι απλό von Neumann
  - είτε
  - (β)  $p_1 \equiv 0$  και το  $p'$  είναι Schur.

Με την θεωρία αυτή ανάγουμε λοιπόν το ερώτημα σε ανάλογο ερώτημα για ένα πολυώνυμο βαθμού κατά ένα μικρότερο και προχωρούμε κατά τον ίδιο τρόπο. Τα κριτήρια αυτά είναι αρκετά εύχρηστα για π.χ.  $k=3$  ή  $4$ .

Αν και στην παράγραφο αυτή δεν μελετούμε πολυώνυμα Schur (εκτός π.χ αν εφαρμόσουμε το κριτήριο (ii.β)) θα δούμε στην Παρ. 3.4 ότι για την λεγόμενη απόλυτη ευστάθεια μίας πολυβηματικής μεθόδου, ανάλογα πολυώνυμα πρέπει να είναι Schur (και όχι απλά von Neumann). Το εξής κριτήριο (Routh-Hurwitz, βλ. Lambert, [3.6, σελ. 80]) είναι επίσης χρήσιμο (για  $k=2,3,4$ , κυρίως, στην πράξη). Θεωρούμε το πολυώνυμο  $p(z)$ , (50), και κάνουμε την αλλαγή των μεταβλητών

$$(53) \quad w = (1+z)/(1-z), \quad z = (w-1)/(w+1)$$

που απεικονίζει (επί) την περιφέρεια  $|w|=1$  στον φανταστικό άξονα  $\operatorname{Re} z = 0$ , τον δίσκο  $D = \{w: |w| < 1\}$  στο ημιεπίπεδο  $\operatorname{Re} z < 0$ , το σημείο  $w=1$  στο

$z=0$ , και το  $w=-1$  στο επ' άπειρου σημείο  $z=\infty$ . Βλέπουμε εύκολα τότε ότι το  $n(z)$  είναι Schur αν και μόνο αν το

$$\tilde{n}(z) = (1-z)^k n((1+z)/(1-z)) = b_0 z^k + \dots + b_k$$

έχει όλες τις ρίζες του με αρνητικά πραγματικά μέρη. Για να ισχύει αυτό ικανές και αναγκαίες ευθήκες είναι οι ευθήκες των Routh-Hurwitz στους συντελεστές  $b_i$ , που για  $k=2,3,4$  είναι οι

$$k=2: b_i > 0, 0 \leq i \leq 2,$$

$$k=3: b_i > 0, 0 \leq i \leq 3, b_1 b_2 - b_3 b_0 > 0,$$

$$k=4: b_i > 0, 0 \leq i \leq 4, b_1 b_2 b_3 - b_0 b_3^2 - b_4 b_1^2 > 0,$$

2. Ένας σχετικά απλός τρόπος για τον προσδιορισμό της τάξης  $p$  και της λεγόμενης "εταθερά εφάλματος"  $c^*$  (που θα οριστεί παρακάτω) μιάς πολυημετικής μεθόδου είναι ο εξής: Υποθέτουμε κατ' αρχήν ότι τα πολύνομα  $p$  και  $\epsilon$  δεν έχουν κοινό παράγοντα και ότι η μέθοδος είναι ευνεής, οπότε  $p(1)=0$ ,  $\epsilon(1)=p'(1)$ . Συνεπώς  $\epsilon(1) \neq 0$  (αλλιώς τα  $p, \epsilon$  θα είχαν τον κοινό παράγοντα  $(z-1)$ ). Έστω ότι η μέθοδος (3) έχει τάξη  $p$ . Τότε για κάθε ομαλή συνάρτηση  $\psi(t)$  έχουμε ((22), (23)), ότι όταν  $h \rightarrow 0$

$$L[\psi(t); h] \equiv \sum_{j=0}^k \{a_j \psi(t+jh) - h \beta_j \psi'(t+jh)\} \sim C_{p+1} h^{p+1} \psi^{(p+1)}(t),$$

όπου με το σύμβολο  $\sim$  εννοούμε ότι ο πρώτος μη μηδενικός όρος του αναπτύγματος του  $L[\psi(t); h]$  σε δυνάμεις του  $h$  για  $h \rightarrow 0$  είναι ο  $C_{p+1} h^{p+1} \psi^{(p+1)}(t)$ . Επειδή τα  $C_j, p$  είναι ανεξάρτητα της  $\psi(t)$  διαλέγουμε  $\psi(t) = e^t$ , οπότε η παραπάνω σχέση γράφεται

$$\sum_{j=0}^k (a_j e^{t+jh} - h \beta_j e^{t+jh}) \sim C_{p+1} h^{p+1} e^t,$$

από την οποία, θέτοντας  $e^h = \zeta$ , δηλ.  $h = \ln \zeta$ ,  $\zeta > 1$  παίρνουμε

$$\sum_{j=0}^k (\alpha_j \zeta^j - \ln \zeta - \beta_j \zeta^j) \sim c_{p+1} (\ln \zeta)^{p+1}, \quad \zeta \downarrow 1.$$

Αναπτύσσοντας τον  $\ln \zeta$  σε δυνάμεις του  $\zeta - 1$  παίρνουμε

$$\ln \zeta - [p(\zeta)/\epsilon(\zeta)] \sim c_{p+1} (\zeta - 1)^{p+1} / \epsilon(\zeta), \quad \zeta \downarrow 1,$$

και επειδή, καθώς  $\zeta \rightarrow 1$ ,  $\epsilon(\zeta) = \epsilon(1) + O((\zeta - 1))$  και  $\epsilon(1) \neq 0$ , έχουμε ότι  $(\epsilon(\zeta))^{-1} = (\epsilon(1))^{-1} (1 + O((\zeta - 1)))$ . Συνεπώς, ορίζοντας την "εταθερά εφάλματος"  $c^*$  της μεθόδου ως

$$(54) \quad c^* = c_{p+1} / \epsilon(1),$$

παίρνουμε από την παραπάνω σχέση ότι

$$(55) \quad \ln \zeta - [p(\zeta)/\epsilon(\zeta)] \sim c^* (\zeta - 1)^{p+1}, \quad \zeta \downarrow 1.$$

Άρα, για να βρούμε τα  $c^*, p$ , αναπτύσσουμε τον  $\ln \zeta$  σε δυνάμεις του  $\zeta - 1$ :

$$\ln \zeta = (\zeta - 1) - (\zeta - 1)^2 / 2 + (\zeta - 1)^3 / 3 - \dots + (-1)^{n-1} (\zeta - 1)^n / n + \dots$$

καθώς και το πηλίκο  $p(\zeta)/\epsilon(\zeta)$  (παρατηρούμε ότι  $p(\zeta)/\epsilon(\zeta) = (\zeta - 1) + O((\zeta - 1)^2)$ ), και βρίσκουμε τον πρώτο μη μηδενικό όρο της διαφοράς των δύο αναπτυγμάτων.

3. Στην περίπτωση μίας πεπλεγμένης πολυβηματικής μεθόδου πρέπει, όπως είδαμε, να λύσουμε σε κάθε βήμα ένα μη γραμμικό σύστημα για τον υπολογισμό του αγνώστου  $y^{n+k}$ . Στην περίπτωση των πολυβηματικών μεθόδων είναι πολύ συνηθισμένο στην πράξη να γίνεται μία αρκετά ακριβής πρόβλεψη  $y_0^{n+k}$  του  $y^{n+k}$  με μία βοηθητική άμεση πολυβηματική μέθοδο και να χρησιμοποιείται η κύρια, (πεπλεγμένη) μέθοδος για τον υπολογισμό διορθώσεων  $y_j^{n+k}$ ,  $j \geq 1$  με απλή επανάληψη. Προκύπτουν

έτσι ζευγάρια μεθόδων, οι λεγόμενες μέθοδοι πρόβλεψης-διόρθωσης (Π-Δ) που μπορούν να γραφούν γενικά στη μορφή

$$\text{Πρόβλεψη (Π): } y_0^{n+k} + \sum_{j=0}^{k-1} \tilde{\alpha}_j y^{n+j} = h \sum_{j=0}^{k-1} \tilde{\beta}_j f^{n+j}, \quad (\tilde{\alpha}_k = 1)$$

$$\begin{aligned} \text{Διόρθωση (Δ): } y_{\omega_i}^{n+k} + \sum_{j=0}^{k-1} \alpha_j y^{n+j} &= h \beta_k f(t^{n+k}, y_i^{n+k}) + \\ &+ h \sum_{j=0}^{k-1} \beta_j f^{n+j}, \quad i=0, 1, 2, \dots \quad (\alpha_k = 1). \end{aligned}$$

Η μέθοδος πρόβλεψης (Π) είναι μία άμεση k-βηματική μέθοδος  $\{\tilde{\alpha}_j, \tilde{\beta}_j\}_{j=0}^k$ ,  $\tilde{\beta}_k = 0$ , που παράγει μία αρχική τιμή  $y_0^{n+k}$  για την μέθοδο διόρθωσης (Δ). Η μέθοδος διόρθωσης χρησιμοποιείται είτε ως απλή επαναληπτική μέθοδος για την λύση του μη γραμμικού ευετήματος  $y^{n+k} = h \beta_k f^{n+k} + g^n$ , οπότε η διόρθωση επαναλαμβάνεται μέχρις ότου π.χ.

$\|y_{\omega_i}^{n+k} - y_i^{n+k}\| \leq \epsilon$  για κάποιο δεδομένο μικρό  $\epsilon$  της τάξεως του εφάλματος στρογγύλευσης, είτε για να παράγει ένα  $y_i^{n+k}$  για μικρό  $i$  (συνήθως  $i=1$  ή  $2$ ) το οποίο ορίζουμε ως  $y^{n+k}$  και προχωρούμε στο επόμενο βήμα. Στην δεύτερη περίπτωση το  $y^{n+k}$  δεν είναι βέβαια η ακριβής λύση του μη γραμμικού ευετήματος που παριστάνει η πεπλεγμένη μέθοδος (Δ) αλλά μία προσέγγιση, η οποία πολλές φορές (αν το ζευγάρι έχει κατάλληλες ιδιότητες ευετάθειας και ακρίβειας) είναι σχεδόν εξ ίσου ακριβής ως προσέγγιση της  $y(t^{n+k})$  όσο και η λύση του μη γραμμικού ευετήματος.

Δύο γνωστά ζευγάρια (Π)-(Δ) είναι:

(i) Euler-Τραπεζίου

$$(Π): y^{n+1} - y^n = h f^n$$

$$(Δ): y^{n+1} - y^n = h(f^{n+1} + f^n)/2$$

(ii) Μέθοδος του Milne

$$(Π): y^{n+4} - y^n = 4h(2f^{n+3} - f^{n+2} + 2f^{n+1})/3$$

$$(Δ): y^{n+4} - y^{n+2} = h(f^{n+4} + 4f^{n+3} + f^{n+2})/3$$

Αν η τάξη ακρίβειας της  $(\Pi)$  είναι  $\tilde{p}$  και της  $(\Delta)$  είναι  $p$ , τότε, ανάλογα με τον αριθμό  $m \geq 1$  των διορθώσεων (δηλ. της εφαρμογής της  $(\Delta)$  για  $i=0,1,2,\dots,m-1$ ), μπορούμε να υπολογίσουμε την τάξη ακρίβειας του ζεύγους (που ορίζεται όπως π.χ. ορίζουμε την τάξη ακρίβειας των μεθόδων RK). Π.χ. αν  $m \geq 1$  και  $\tilde{p} \geq p-1$  τότε η τελική τάξη ακρίβειας είναι  $p$ . Αν  $\tilde{p} = p-2$  και  $m=1$  η τελική τάξη είναι  $p-1$ . Αν  $\tilde{p} = p-2$  και  $m > 1$  τότε η τελική τάξη είναι  $p$ . Π.χ. η μέθοδος Euler-τραπεζίου με  $\tilde{p}=1$ ,  $p=2$  έχει τελική τάξη ακρίβειας 2 για οποιοδήποτε αριθμό διορθώσεων  $m \geq 1$ . Για την μέθοδο του Milne  $\tilde{p}=p=4$ , οπότε και η τελική τάξη είναι 4. Η μέθοδος  $(\Delta)$  εξ άλλου καθορίζει την ευστάθεια του ζεύγους το οποίο είναι ευσταθές αν και μόνο αν η  $(\Delta)$  είναι ευσταθής. Π' άλλα λόγια δεν χρειάζεται και η  $(\Pi)$  να είναι ευσταθής.

4. Όπως γίνεται για τις μεθόδους RK, είναι δυνατόν να μεταβάλουμε το βήμα  $h_j = t^{j+1} - t^j$  και για τις πολυβηματικές μεθόδους. Αυτό αυξάνει την λοχιστική των μεθόδων γιατί θα πρέπει (με παρεμβολή π.χ.) κάθε φορά που αλλάζει το βήμα να βρίσκονται και προεχθίσεις σε διαφορετικές προηγούμενες τιμές του  $t$ . Στρατηγικές για τον έλεγχο του μεγέθους του  $h_j$  ετηρίζονται και πάλι σε κάποιου είδους εκτίμηση του τοπικού εφάλματος, η οποία είναι αρκετά απλή αν χρησιμοποιούμε ζεύγη μεθόδων πρόβλεψης-διόρθωσης. Καλά προγράμματα που ετηρίζονται σε μεθόδους πρόβλεψης-διόρθωσης εκτός από μεταβλητό βήμα και εκτίμηση του τοπικού εφάλματος έχουν την δυνατότητα να μεταβάλλουν και την τάξη ακρίβειάς τους, δηλ. να καταφεύχουν σε διαφορετικά ζεύγη  $\Pi$ - $\Delta$ , μεγαλύτερης (μικρότερης) ευνολικής τάξης ακρίβειας, αν η λύση μεταβάλλεται - με δείκτη το μέγεθος του τοπικού εφάλματος - πολύ (λίγο) από βήμα σε βήμα. Χαρακτηριστικά τέτοια προγράμματα είναι το πρόγραμμα DIFSUB του Gear, βλ. [3.2], το πρόγραμμα DVDF του Krogh κ.ά, που υπάρχουν σε πολλές βιβλιοθήκες αλγορίθμων.

### Ασκήσεις 3.3

1. Αν  $e(z) = z^2$ , βρείτε  $p(z)$  (με  $a_k$  πιθανώς όχι 1) τέτοιο ώστε:

(α). Το  $p(z)$  να είναι δευτέρου βαθμού και η τάξη της μεθόδου  $(p, \epsilon)$  να είναι 2.

(β). Το  $p(z)$  να είναι τρίτου βαθμού και η τάξη της μεθόδου  $(p, \epsilon)$  να είναι 3.

(γ). Είναι ευσταθείς οι δύο μέθοδοι;

2. Αν  $p(z) = z^4 - 1$ , βρείτε  $\epsilon(z)$  βαθμού 4 τέτοιο ώστε η μέθοδος  $(p, \epsilon)$  να έχει μέγιστη τάξη. Ποιά είναι η τάξη της και ποιά η σταθερά του εφάλματος;

3. Βρείτε τους συντελεστές  $\{\alpha_j, \beta_j\}$   $0 \leq j \leq 2$ ,  $\alpha_2 = 1$  της γενικής διβηματικής μεθόδου έτσι ώστε η τάξη ακρίβειάς της να είναι  $p \geq 2, 3$  ή

4. Υπάρχει μέθοδος με  $p=4$ ; Με  $p \geq 5$ ;

4. Βρείτε τους συντελεστές  $\alpha_j, \beta_j$  της τριβηματικής μεθόδου μέγιστης τάξης ακρίβειας. Ποιά είναι η τάξη της; Είναι ευσταθής;

5. Εκφράστε τους συντελεστές της οικογένειας των διβηματικών μεθόδων τάξης  $\geq 3$  συναρτήσει της παραμέτρου  $\beta_0$ . Για ποιές τιμές του  $\beta_0$  είναι ευσταθείς οι μέθοδοι; Εκφράστε την σταθερά εφάλματος ως συνάρτηση του  $\beta_0$  (για τιμές του  $\beta_0$  που δίνουν ευσταθείς μεθόδους).

Τι παρατηρείτε για την μέθοδο του Simpson. ;

6. Βρείτε τις ευσταθείς 4-βηματικές μεθόδους με τάξη ακρίβειας  $p=6$ .

7. Αποδείξτε ότι οι μέθοδοι οπισθοδρομικών διαφορών με  $k$  βήματα (6) έχουν τάξη ακρίβειας  $k$ .

8. Αποδείξτε, χρησιμοποιώντας την θεωρία Schur (Παρατήρηση 1) ότι οι  $k$ -βηματικές μέθοδοι οπισθοδρομικών διαφορών (6) είναι ευσταθείς για  $k=1, 2, 3, 4$ . (Είναι γνωστό ότι οι μέθοδοι αυτές είναι ευσταθείς αν και μόνο αν  $1 \leq k \leq 6$ ).

9. Σε αναλογία με ό,τι κάναμε στις μεθόδους RK ορίστε την τάξη ακρίβειας και την ευστάθεια της μεθόδου (Π)-(Δ) Euler - τραπεζίου (βλ. Παρατήρηση 3(1)) και αποδείξτε ότι η τάξη ακρίβειας είναι  $p=2$  και ότι η μέθοδος είναι ευσταθής.

10. Για δευτεροβάθμιες Δ.Ε. της μορφής  $y''=f(t,y)$  - για τις οποίες δηλ. η  $f$  δεν είναι συνάρτηση του  $y'$  - είναι δυνατόν να αναπτύξουμε ανάλογη θεωρία  $k$ -βηματικών μεθόδων της μορφής

$$\sum_{j=0}^k \alpha_j y^{n+j} = h^2 \sum_{j=0}^k \beta_j f^{n+j}. \text{ Η συνθήκη ευστάθειας τώρα στο πολώνυμο}$$

$p(z) = \sum_{j=0}^k \alpha_j z^j$  είναι ότι όλες οι ρίζες του  $z_i$  πρέπει να ικανοποιούν  $|z_i| \leq 1$  και ότι η πολλαπλότητα των ριζών με  $|z_i|=1$  πρέπει να είναι το πολύ 2. Αν ισχύουν αυτές οι συνθήκες, δείξτε ότι στο ανάπτυγμα

$$[1/(\alpha_k + \alpha_{k-1}z + \dots + \alpha_0 z^k)] = \gamma_0 + \gamma_1 z + \gamma_2 z^2 + \dots$$

οι συντελεστές  $\gamma_i$  αυξάνουν "γραμμικά", δηλ. ότι υπάρχουν σταθερές  $c_1$  και  $c_2$  τέτοιες ώστε

$$|\gamma_i| \leq c_1 i + c_2, \quad i=0,1,2,\dots$$

### 3.4 ΑΚΑΜΠΤΑ ΠΡΟΒΛΗΜΑΤΑ. ΑΠΟΛΥΤΗ ΕΥΣΤΑΘΕΙΑ ΚΑΙ ΓΕΝΙΚΕΥΣΕΙΣ ΤΗΣ

Θεωρούμε το πρόβλημα αρχικών τιμών για την απλή Δ.Ε.

$$(1) \begin{cases} y' = \lambda y, & t \geq 0, \\ y(0) = 1, \end{cases}$$

και υποθέτουμε ότι η σταθερά  $\lambda$  είναι τέτοια ώστε  $\lambda < 0$ , δηλ.  $\lambda < 0$  με  $|\lambda| \gg 1$ . Προφανώς η λύση  $y(t) = e^{\lambda t}$  του (1) τείνει πολύ γρήγορα στο μηδέν καθώς αυξάνεται το  $t$ . Ας προεχθίσουμε το πρόβλημα (1) με την μέθοδο του Euler με κάποιο σταθερό βήμα  $h$ . Προφανώς προκύπτει η ακολουθία  $\{y^n\}$ ,  $n \geq 0$  όπου  $y^{n+1} = (1+h\lambda)y^n$ ,  $n \geq 0$ ,  $y^0 = 1$ , δηλ. η ακολουθία

$$(2) y^n = (1+h\lambda)^n, \quad n \geq 0.$$

Έστω ότι σε κάποιο βήμα  $k$  του αλγορίθμου υπολογίζουμε, π.χ. λόγω εφαλμάτων ετρογχύλευσης, τον αριθμό  $z^k$  αντί του  $y^k$  και ότι για  $n > k$  υπολογίζουμε τα  $z^n$  ακριβώς. Θεωρούμε δηλ. την ακολουθία  $\{z^n\}$ ,  $n \geq 0$ :

$$z^n = y^n, \quad 0 \leq n \leq k-1, \quad z^{n+1} = (1+h\lambda)z^n, \quad n \geq k, \quad z^k \neq y^k.$$

Συνεπώς έχουμε

$$y^n - z^n = (1+h\lambda)^{n-k} (y^k - z^k), \quad n \geq k,$$

δηλ. ότι

$$(3) |y^n - z^n| = |1+h\lambda|^{n-k} |y^k - z^k|, \quad n \geq k,$$

από την οποία βλέπουμε ότι αν  $|1+h\lambda| < 1$  (δηλ. αν  $0 < h < -2/\lambda$ ) η διαταραχή  $|y^k - z^k|$  στο βήμα  $k$  δεν θα έχει εσβαρές επιπτώσεις. Πράγματι, τότε,  $|1+h\lambda|^{n-k} \rightarrow 0$ ,  $n \rightarrow \infty$ , δηλ. η προκαλούμενη μεταβολή  $|y^n - z^n|$  (θεωρούμενη ως διαταραχή της "ακριβούς" λύσης  $y^n$  του



διακριτού προβλήματος) θα γίνεται όλο και πιο μικρή καθώς αυξάνεται το  $n$ . Αν  $|1+h\lambda|=1$  το εφάλμα  $|y^k-z^k|$  διατηρείται ( $|y^n-z^n| = |y^k-z^k|$ ,  $n \geq k$ ), ενώ αν  $|1+h\lambda| > 1$ , θα έχει καταστροφικές επιπτώσεις γιατί τότε  $|y^n-z^n| \rightarrow \infty$  εκθετικά με το  $n$ . Ας σημειωθεί ότι αν π.χ. υπολογίζουμε στο διάστημα  $0 \leq t \leq T$ , έχουμε, για  $h=T/N$ , ότι όντως ισχύει η εκτίμηση

$$\begin{aligned} \max_{0 \leq n \leq N} |y^n - z^n| &\leq |1+h\lambda|^{N-k} |y^k - z^k| \leq (1+h|\lambda|)^N |y^k - z^k| \\ &\leq e^{|\lambda|T} |y^k - z^k|, \end{aligned}$$

(όπως προβλέπεται από την Πρόταση 3.2) που εκφράζει στην απλή μας περίπτωση την "ευστάθεια" της μεθόδου του Euler ως μεθόδου RK, με την έννοια της "ευστάθειας" που εισαγάγαμε στην παράγραφο 3.2. Όμως η σταθερά Lipschitz  $L=|\lambda|$  του προβλήματός μας είναι πολύ μεγάλη έτσι ώστε το φράγμα  $e^{|\lambda|T}$  να μην έχει πρακτική αξία στην περίπτωση μας. Το παράδειγμα αυτό μας δείχνει όμως ότι μία νέα ευνοϊκή, δηλ. η  $|1+h\lambda| < 1$ , εξασφαλίζει μία πραγματική ευστάθεια της μεθόδου με την έννοια ότι κατατείνει τυχόν "εφάλματα"  $|y^k-z^k|$  και δεν τους επιτρέπει να επηρεάσουν σημαντικά τους υπολογισμούς για  $n > k$ . Η ευνοϊκή αυτή αποτελεί ένα σοβαρό περιορισμό στο μέγεθος του βήματος  $h$ .

Σημειώστε ότι η ευνοϊκή  $|1+h\lambda| < 1$  εξασφαλίζει επίσης ότι η ακολουθία  $\{y^n\}$  που παράγει η μέθοδος του Euler έχει την ιδιότητα (βλ. (2)) ότι  $y^n \rightarrow 0$ ,  $n \rightarrow \infty$  (Τονίζουμε ότι παίρνουμε το όριο για σταθερό  $h$  όταν  $n \rightarrow \infty$ ). Δηλ. είναι ικανή και αναγκαία ευνοϊκή έτσι ώστε η λύση  $\{y^n\}$  του διακριτού προβλήματος (για κάθε σταθερό  $h$ ) να μιμείται την συμπεριφορά της λύσης του συνεχούς προβλήματος (1),  $y(t) = e^{\lambda t}$ , για την οποία προφανώς  $y(t) \rightarrow 0$ ,  $t \rightarrow \infty$ .

Θεωρείστε τώρα την ακολουθία  $\{y^n\}$  που παράγει η πεπλεγμένη μέθοδος του Euler για το ίδιο πρόβλημα, δηλ. η μέθοδος  $y^{n+1} - y^n = h\lambda y^{n+1}$ ,  $n \geq 0$ ,  $y^0 = 1$ , που γράφεται ως

$$(4) \quad y^n = (1-h\lambda)^{-n}, \quad n \geq 0.$$

Για την (4) το ανάλογο της σχέσης (3) είναι

$$(5) \quad |y^n - z^n| = |1 - h\lambda|^{-(n-k)} |y^k - z^k|, \quad n \geq k.$$

Επειδή  $\lambda < 0$  συμπεραίνουμε ότι  $|1 - h\lambda|^{-1} < 1 \quad \forall h > 0$ , δηλ. ότι η (5) για κάθε  $h > 0$  δίνει  $|y^n - z^n| \ll |y^k - z^k|$ , αν  $n \gg k$ : η πεπλεγμένη μέθοδος του Euler δεν επηρεάζεται εχεθόν καθόλου από εφάλματα  $|y^k - z^k| \neq 0$ . Ισοδύναμα,  $\forall h > 0$ , η ακολουθία (4) ικανοποιεί  $y^n \rightarrow 0, n \rightarrow \infty$  δηλ. το διακριτό ανάλογο της ιδιότητας  $y(t) \rightarrow 0, t \rightarrow \infty$  της λύσης του συνεχούς προβλήματος (1).

Ας θεωρήσουμε τώρα το πρόβλημα

$$(6) \quad \begin{cases} y' = \lambda y + f'(t) - \lambda f(t), & t \geq 0, \\ y(0) = y_0, \end{cases}$$

του οποίου η λύση είναι προφανώς

$$(7) \quad y(t) = f(t) + e^{\lambda t}(y_0 - f(0)), \quad t \geq 0.$$

Υποθέτουμε ξανά ότι  $\lambda < 0$  και ότι η  $f(t)$  είναι μία δεδομένη ομαλή συνάρτηση που δεν μεταβάλλεται γρήγορα με το  $t$ . Από τον τύπο (7) βλέπουμε ότι η λύση  $y(t)$  εκφράζεται ως άθροισμα του όρου  $f(t)$  και του εφήμερου όρου  $e^{\lambda t}(y_0 - f(0))$  που τείνει στο 0 πολύ γρήγορα (δηλ. που εξαφανίζεται εχεθόν αμέσως και δεν συμβάλλει στην λύση) καθώς αυξάνει το  $t$ . Προφανώς ο σημαντικός όρος για την λύση για  $t > 0$  είναι ο  $f(t)$ .

Γράφοντας την (7) για  $t+h$  αντί  $t$  και απαλείφοντας του αρχικό όρο  $y_0 - f(0)$  παίρνουμε την ταυτότητα

$$(8) \quad y(t+h) = f(t+h) + e^{\lambda h}(y(t) - f(t)), \quad t, h \geq 0.$$

Η μέθοδος του Euler για το βήμα  $t^n \rightarrow t^{n+1}$  γράφεται ως  $y^{n+1} = y^n + h(\lambda y^n + f'(t^n) - \lambda f(t^n))$ , δηλ. ως

$$(9) \quad y^{n+1} = f(t^n) + hf'(t^n) + (1+h\lambda)(y^n - f(t^n)).$$

Παρατηρούμε ότι  $e^{\lambda h} = 1 + h\lambda + O(h^2\lambda^2)$  και ότι  $f(t^n+h) = f(t^n) + hf'(t^n) + O(h^2)$ , δηλ. ότι άυτως η (9) είναι προσέγγιση της (8) για  $t=t^n$ . Ενώ όμως στην (8) η διαφορά  $(y(t)-f(t))$  - που μπορεί να θεωρηθεί ως "διαταραχή" στην τιμή της "θεμελιώδους" λύσης στο  $t+h$ ,  $y(t+h) \approx f(t+h)$  - πολλαπλασιάζεται επί  $e^{\lambda h} \ll 1$  (για  $\lambda h \ll 0$ ) και δεν διαταράσσει ευσεπώς εχεδόν καθόλου την εχέση  $y(t+h) \approx f(t+h)$ , στην (9), αν  $|1+h\lambda| \geq 1$ , η διαφορά  $y^n - f(t^n)$  επηρεάζει σημαυτικά την εχέση  $y^{n+1} \approx f(t^{n+1}) + hf'(t^n)$  δηλ. την  $y^{n+1} \approx f(t^{n+1})$ . Βλέπουμε δηλ. πάλι του σημαυτικό ρόλο της ευσηήκης  $|1+h\lambda| < 1$  για την μέθοδο του Euler για αυτό το πρόβλημα. Η πεπλεγμένη μέθοδος του Euler, όπως μπορούμε εύκολα να δούμε, δίνει την εχέση

$$y^{n+1} = (1-h\lambda)^{-1} [f(t^n) + hf'(t^{n+1}) - \lambda hf(t^{n+1})] + (1-h\lambda)^{-1} (y^n - f(t^n)),$$

που μας δείχνει - παρατηρείτε ότι ο πρώτος όρος είναι μία  $O(h^2)$  προσέγγιση της  $f(t^{n+1})$  - ότι για κάθε  $h > 0$  η "διαταραχή"  $y^n - f(t^n)$  εμποδίζεται να διαταράξει την εχέση  $y^{n+1} \approx f(t^{n+1})$  επειδή πολλαπλασιάζεται επί του παράγοντα  $(1-h\lambda)^{-1} \ll 1$ .

Το πρόβλημα (6) αποτελεί τυπικό δείγμα προβλήματος αρχικών τιμών για μία άκαμπτη Δ.Ε.. Τα χαρακτηριστικά των ακάμπτων προβλημάτων διαφαίνονται ήδη από τύπους όπως οι (7) και (8): Η λύση τους είναι βασικά ομαλή ( $y(t) \approx f(t)$ ) και μεταβάλλεται αρχά με του χρόνο για  $t > \epsilon > 0$ , ε "μικρό", αλλά περιέχει μία ευαιστώσα η οποία αποεβένεται εκθετικά καθώς αυξάνει το  $t$ , δηλ. που πρακτικά παύει να ευαισφέρει στη λύση για  $t > \epsilon$ , βλ. (7). Απ' την άλλη μεριά η (8) μας λέει ότι για κάθε  $t$  (π.χ. μεγάλο) η λύση  $y(t+h)$  είναι περίπου ίση με την τιμή  $f(t+h)$  του θεμελιώδους "φορέα" της· η εχέση όμως  $y(t+h) \approx f(t+h)$  διαταράσσεται από μία μικρή μεταβολή (στο πρόβλημά μας την  $e^{h\lambda}(y(t)-f(t))$ ) που προκύπτει από το ότι  $y(t) \neq f(t)$ , αλλά που αποεβένεται από του παράγοντα  $e^{\lambda h}$ . (Θυσιαστικά ο (8) είναι ο ίδιος τύπος με του (7) αλλά με αρχικές ευσηήκες στο σημείο  $t$ ). Αν λοιπόν η αριθμητική μέθοδος προσεγγίζει το εκθετικό  $e^{\lambda h}$  με όχι καλό τρόπο για  $h\lambda \ll 0$ , όπως π.χ. συμβαίνει με την μέθοδο του Euler (για την οποία η

προέχγιση  $1+h\lambda$  είναι μεν αποδεκτή για  $|h\lambda| \ll 1$ , αλλά ανακριβής και ποιοτικά τελείως λανθασμένη για  $h\lambda \ll 0$  οπότε  $|1+h\lambda| \gg 1$ ), τότε περιμένουμε σοβαρή διαταραχή της λύσης  $y^n$  (αστάθεια) που θα προέλθει στο παράδειγμά μας από του όρο  $(1+h\lambda)(y^n - f(t^n))$ . Πρέπει δηλ. να επιβάλουμε στο  $h$  την περιοριστική συνθήκη  $|1+h\lambda| < 1 \Leftrightarrow 0 < h < -2/\lambda$ , η οποία πρέπει να ισχύει ε' όλη την διάρκεια των υπολογισμών παρά το γεγονός ότι ο παράγοντας  $e^{\lambda t}(y_0 - f(0))$  της λύσης που προκαλεί αυτόν τον περιορισμό στο  $h$  δεν ευνοιάζει πρακτικά τίποτε στη λύση για  $t > 0$ ! Στις σημειώσεις [5.3], εσλ. 163 et seq. γίνεται η εφαρμογή της μεθόδου του Euler στο παράδειγμα (6) με  $\lambda = -1000$ ,  $f(t) = t^2$ ,  $y(0) = 0$ . Παρατηρούμε ότι  $y_0 - f(0) = 0$ , δηλ. ότι ο παράγοντας  $e^{\lambda t}(y_0 - f(0))$  δεν υπάρχει καθόλου στην λύση  $y(t) = t^2$ . ο τύπος (9) εξακολουθεί βέβαια να ισχύει. Αναπόφευκτα το αριθμητικό πείραμα δείχνει καταστροφικά γρήγορα αυξανόμενη αστάθεια αν  $h > .002$  και καλά αποτελέσματα - με  $O(h^2)$  βέβαια σφάλμα - για  $h < .002$  όπως προβλέπει η συνθήκη  $|1 - 1000h| < 1$ . Αντίθετα η πεπλεγμένη μέθοδος του Euler προεχγίζει το εκθετικό  $e^{\lambda h}$  με την ποιοτικά εωστή για  $\lambda h \ll 0$  ρητή προέχγιση  $(1 - \lambda h)^{-1}$  (Πάλι τάξης  $O(h^2 \lambda^2)$  κοντά στο 0 αλλά που μιμείται την ιδιότητα  $e^{\lambda h} \ll 1$  για  $\lambda h \ll 0$ ). Συνεπώς για κάθε  $h > 0$  η πεπλεγμένη μέθοδος δεν θα έχει φαινόμενα αστάθειας προερχόμενα από του όρο  $(1 - \lambda h)^{-1}(y^n - f(t^n))$ .

Μετά από τις εισαγωγικές αυτές παρατηρήσεις προχωρούμε ε' ένα ορισμό, που σφείζεται (όπως και μεγάλο μέρος της αρχικής έρευνας για αυτά τα θέματα) στον Dahlquist (1963). Λέμε ότι μιά μέθοδος (για κάποιο ορισμένο βήμα  $h$ ) είναι απόλυτα ευσταθής (για απλές Δ.Ε.) αν, όταν εφαρμοσθεί στο πρόβλημα (1) με  $\lambda < 0$ , έχει την ιδιότητα να δίνει προεχγίσεις  $y^n$ ,  $n \geq 0$  τέτοιες ώστε  $y^n \rightarrow 0$ ,  $n \rightarrow \infty$  (για σταθερό  $h$ ). Οι τιμές του γινομένου  $\lambda h$  για τις οποίες η μέθοδος έχει αυτή την ιδιότητα αποτελούν την περιοχή απόλυτης ευστάθειας της μεθόδου που είναι, ευσεπώς, υποσύνολο (για απλές Δ.Ε.) του αρνητικού πραγματικού ημίξονα. Συνήθως ενδιαφερόμαστε να εντοπίσουμε (αν υπάρχει) κάποιο διάστημα απόλυτης ευστάθειας με μέγιστο μήκος αλλά της μορφής  $(r, 0)$ ,  $-\infty < r < 0$ . Τότε, αν  $h\lambda \in (r, 0)$ , δηλ. αν επιλέξουμε  $h$  τέτοιο ώστε  $0 < h < r/\lambda$ , η μέθοδος θα είναι απόλυτα ευσταθής, εξ ορισμού. Παραδείγματος χάριν είδαμε ότι η μέθοδος του Euler έχει

διάστημα απόλυτης ευστάθειας το  $(-2,0)$  και η πεπλεγμένη μέθοδος του Euler το  $(-\infty,0)$ , δηλ. είναι απόλυτα ευσταθής για κάθε  $h$ . Δεν είναι δύσκολο να δούμε ότι η μέθοδος του τραπεζίου (3.2.5) έχει διάστημα απόλυτης ευστάθειας επίσης ίσο με  $(-\infty,0)$  (δηλ. είναι και αυτή ιδιαίτερα κατάλληλη για προβλήματα όπως το (1) με  $\lambda < 0$ ). Αντίθετα η μέθοδος του μέσου (άμεση) (3.2.6) έχει  $(-2,0)$ , οι άμεσες μέθοδοι RK τρίτης τάξης (3.20 α,β) έχουν  $(-2.51,0)$ , οι κλασικές μέθοδοι RK (3.21α,β) έχουν  $(-2.76;0)$  κ.ο.κ. Θα μελετήσουμε πιο συστηματικά την απόλυτη ευστάθεια των μεθόδων RK και των πολυβηματικών μεθόδων παρακάτω.

Μία παρατήρηση: εξετάζουμε το πρόβλημα  $y' = \lambda y$ ,  $\lambda < 0$ , γιατί το αντίστοιχο με  $\lambda > 0$  έχει εκθετικά αυξανόμενες λύσεις  $e^{\lambda t}$  για τις οποίες για  $\lambda$  ή  $t \gg 0$  είναι δύσκολο να διακρίνει κανείς την λύση από την αεταθή διαταραχή της. Γεννάται επίσης το εύλογο ερώτημα κατά πόσο η συμπεριφορά μιάς αριθμητικής μεθόδου στην περίπτωση του απλού παραδείγματος (1) μπορεί να μας οδηγήσει σε συμπεράσματα για την συμπεριφορά της σε πιο πολύπλοκα, μη γραμμικά προβλήματα. Θα μελετήσουμε λοιπόν αρχότερα, επεκτάσεις της έννοιας της απόλυτης ευστάθειας για (κατάλληλα) μη γραμμικά προβλήματα. Προς το παρόν ας θεωρήσουμε την ευσταθή συμπεριφορά της ε' ένα απλό πρόβλημα όπως το (1) ως αναγκαία συνθήκη για ευσταθή συμπεριφορά σε πιο πολύπλοκα προβλήματα. Δεν είναι όμως εαφές ότι είναι και ικανή συνθήκη αν και κανείς θα μπορούσε εύκολα να διατυπώσει την εικασία ότι ευσταθής συμπεριφορά για το (1) συνεπάγεται ευσταθή συμπεριφορά για προβλήματα με μεγάλες (απόλυτως) σταθερές Lipschitz, που όμως έχουν φθίνουσες λύσεις, όπως π.χ.  $y' = \lambda(t)y$  με  $\lambda(t) < 0$ ,  $\max_t |\lambda(t)| \gg 1$  ή  $y' = f(y)$  με  $f'(y) < 0$  αλλά με  $\max_y |f'(y)| \gg 1$  κ.ο.κ. Αυτό είναι για ορισμένες μεθόδους σωστό αλλά, όπως θα δούμε, υπάρχουν απόλυτα ευσταθείς μέθοδοι που δεν είναι ευσταθείς για ανάλογα μη γραμμικά προβλήματα ή ακόμα και για προβλήματα με μεταβλητούς συντελεστές  $\lambda(t)$ .

Μία σχετική απλή επέκταση όμως του ορισμού της απόλυτης ευστάθειας αρκεί για την ρεαλιστική μελέτη της ευστάθειας μιάς μεθόδου όταν εφαρμοσθεί σε γραμμικά ευστήματα δ.Ε. με σταθερούς συντελεστές, δηλ. σε ευστήματα της μορφής  $y' = Ay + F(t)$ , όπου  $A \in \mathbb{R}^{m \times m}$  ( $a_{ij}$  σταθερές),  $y \in \mathbb{R}^m$ . θεωρούμε το εξής πρόβλημα αρχικών τιμών

γιά το ομογενές γραμμικό σύστημα

$$(10) \quad \begin{cases} y' = Ay, & t \geq 0, \\ y(0) = y_0 \end{cases}$$

και υποθέτουμε για απλούστευση ότι ο  $A$  διαγωνοποιείται, δηλ. ότι υπάρχει αντιστρέψιμος πίνακας  $S$  τέτοιος ώστε  $S^{-1}AS = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ , όπου  $\lambda_i$ ,  $1 \leq i \leq m$  οι (εν γένει μιγαδικές) ιδιοτιμές του  $A$ . Η αλλαγή των μεταβλητών  $y = S\psi$  δίνει το ισοδύναμο σύστημα

$$(10') \quad \begin{cases} \psi' = \Lambda\psi, & t \geq 0 \\ \psi(0) = \psi_0 = S^{-1}y_0, \end{cases}$$

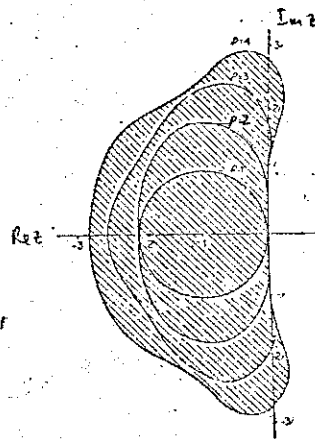
του οποίου η λύση είναι προφανώς  $\psi_i = (\psi_0)_i e^{\lambda_i t}$ . Συμπεραίνουμε π.χ. ότι η λύση  $y(t)$  του (10) έχει την ιδιότητα  $y(t) \rightarrow 0$ ,  $t \rightarrow \infty$  για κάθε  $y_0 \in \mathbb{R}^m$  αν και μόνο αν  $\text{Re} \lambda_i < 0$ ,  $1 \leq i \leq m$ . (Το ίδιο ισχύει, όπως εύκολα μπορούμε να δούμε χρησιμοποιώντας την μορφή Jordan του  $A$ , και στην γενική περίπτωση ενός οποιοδήποτε, όχι αναγκαστικά διαγωνοποιήσιμου, πίνακα  $A \in \mathbb{R}^{m \times m}$ ). Επειδή στο σύστημα (10') οι Δ.Ε. έχουν αποσυνδεθεί σε  $m$  ανεξάρτητες απλές Δ.Ε. της μορφής  $\psi'_i = \lambda_i \psi_i$ , φαίνεται λογικό ότι για να επεκτείνουμε την έννοια της απόλυτης ευστάθειας έτσι ώστε να έχει εφαρμογή σε συστήματα της μορφής (10), θα πρέπει να μελετήσουμε την συμπεριφορά των λύσεων της αριθμητικής μεθόδου όταν εφαρμοσθεί σε απλή Δ.Ε. της μορφής (1) αλλά με μιγαδικό  $\lambda$  με  $\text{Re} \lambda < 0$ .

Λέμε λοιπόν γενικά, σύμφωνα με τον Dahlquist (1963), ότι μία αριθμητική μέθοδος για την λύση προβλημάτων αρχικών τιμών για συστήματα Δ.Ε. είναι απόλυτα ευσταθής (για κάποιο ορισμένο βήμα  $h$ ) αν, όταν εφαρμοσθεί στο πρόβλημα

$$(11) \quad \begin{cases} y' = \lambda y, & t > 0, \lambda \in \mathbb{C}, \text{Re} \lambda < 0 \\ y(0) = 1. \end{cases}$$

δίνει προσεγγίσεις  $y^n$ ,  $n \geq 0$ , τέτοιες ώστε  $y^n \rightarrow 0$ ,  $n \rightarrow \infty$  (εσταθερό  $h$ ). Οι τιμές του γινομένου  $h\lambda$  (υποεύνολο του αριστερού μιγαδικού ανοιχτού ημιεπιπέδου  $\operatorname{Re} z < 0$ ) για τις οποίες μία μέθοδος έχει αυτήν την ιδιότητα αποτελούν την περιοχή απόλυτης ευστάθειας της μεθόδου. Προφανώς, ο νέος ορισμός της απόλυτης ευστάθειας γενικεύει του προηγούμενου, που υπέθετε ότι  $\lambda \in \mathbb{R}$ ,  $\lambda < 0$ . Στο εξής θα χρησιμοποιούμε τον νέο ορισμό.

Η μέθοδος του Euler εφαρμοζόμενη στο πρόβλημα (11) δίνει την (μιγαδική) ακολουθία  $y^n = (1+h\lambda)^n$  για την οποία  $y^n \rightarrow 0$ ,  $n \rightarrow \infty \Leftrightarrow |1+h\lambda| < 1$ . Βέτοντας  $z=h\lambda$  βλέπουμε ότι η περιοχή απόλυτης ευστάθειας της μεθόδου του Euler είναι ο ανοιχτός δίσκος  $\{z: |1+z| < 1\}$  κέντρου  $-1$  και ακτίνας  $1$  στο μιγαδικό επίπεδο. Προφανώς το διάστημα  $(-2, 0)$  απόλυτης ευστάθειας της μεθόδου (για απλές Δ.Ε.) είναι η τομή του  $|1+z| < 1$  με τον πραγματικό άξονα. Η (άμεση) μέθοδος του μέσου (3.2.6) εφαρμοζόμενη στο πρόβλημα (11) δίνει  $y^n = (1+h\lambda + (h^2\lambda^2)/2)^n$ , δηλ. ότι  $y^n \rightarrow 0$ ,  $n \rightarrow \infty \Leftrightarrow |1+z+z^2/2| < 1$ ,  $z=h\lambda$ .



Σχήμα 3.4.1

Το χωρίο  $\{z \in \mathbb{C}: |1+z+z^2/2| < 1\}$  είναι συμμετρικό ως προς τον πραγματικό άξονα (για κάθε πολυώνυμο με πραγματικούς συντελεστές  $p(z) = \overline{p(\bar{z})}$ ), περιέχεται όλο στο ανοιχτό ημιεπίπεδο  $\operatorname{Re} z < 0$  και περιέχει τον δίσκο  $|1+z| < 1$  της μεθόδου του Euler. Το σύνορό του, δηλ. η καμπύλη  $1+z+z^2/2 = \exp(i\theta)$ , μπορεί να σχεδιασθεί εύκολα αν δίνουμε τιμές στο  $\theta$  και λύσουμε κάθε φορά την εξίσωση ως προς  $z$ . Η καμπύλη ( $\rho=2$  στο σχήμα 1) που προκύπτει, εφάπτεται του φανταστικού άξονα στο 0 και της περιφέρειας  $|1+z|=1$  στο  $(-2, 0)$ .

Η πεπλεγμένη μέθοδος του Euler δίνει την ακολουθία  $y^n = (1-h\lambda)^{-n}$ ,  $n \geq 0$ , δηλ. ικανοποιεί  $y^n \rightarrow 0$ ,  $n \rightarrow \infty$  αν και μόνο αν

$$|(1-z)^{-1}| < 1, \quad z=h\lambda,$$

που ισχύει για κάθε  $z$  με  $\operatorname{Re} z < 0$ . Συνεπώς η περιοχή απόλυτης ευσταθείας της μεθόδου είναι όλο το ημιεπίπεδο  $\operatorname{Re} z < 0$ . Τέτοιες μέθοδοι λέγονται A-ευσταθείς, (Dahlquist, 1963). Η μέθοδος του τραπεζίου, η οποία δίνει  $(1-h\lambda/2)y^{n+1} = (1+h\lambda/2)y^n$ , δηλ.

$$y^n = [(1+h\lambda/2)(1-h\lambda/2)^{-1}]^n$$

είναι απόλυτα ευσταθής για  $|(1+z/2)/(1-z/2)| < 1$ ,  $z=h\lambda$ , δηλ. για κάθε  $z$  με  $\operatorname{Re} z < 0$ . Συμπεραίνουμε ότι και η μέθοδος του τραπεζίου είναι A-ευσταθής

λέμε ότι ένα πρόβλημα της μορφής

$$(12) \quad \begin{cases} y' = Ay + F(t), & t \geq 0 \\ y(0) = y_0 \end{cases}$$

όπου οι ιδιοτιμές  $\lambda_i$  του  $A$  ικανοποιούν  $\operatorname{Re} \lambda_i < 0 \quad \forall i$ , είναι άκαμπτο αν  $\max_i |\operatorname{Re} \lambda_i| \gg \min_i |\operatorname{Re} \lambda_i|$ . Δεν είναι δύσκολο να δούμε ότι ένα τέτοιο σύστημα έχει τα χαρακτηριστικά του άκαμπτου προβλήματος για απλές Δ.Ε. που εξετάσαμε προηγουμένως. Υπάρχουν συνιστώσες της λύσης (αυτές που αντιστοιχούν σε ιδιοτιμές με μεγάλο  $|\operatorname{Re} \lambda_i|$ ) που τείνουν στο μηδέν πολύ γρήγορα καθώς αυξάνει το  $t$ , σε σχέση με άλλες συνιστώσες που τείνουν πιο αργά στο 0 ή μεταβάλλονται σε χρονικές κλίμακες που εξαρτώνται από το  $F(t)$ . Για να είναι μία μέθοδος "πραγματικά ευσταθής" για το πρόβλημα (12) είναι προφανές ότι πρέπει να έχει μία κατάλληλη μη κενή περιοχή απόλυτης ευσταθείας  $S$  στο ημιεπίπεδο  $\operatorname{Re} z < 0$  και το βήμα  $h$  να είναι τέτοιο ώστε  $h\lambda_i \in S, \forall i$ . Αν η περιοχή  $S$  είναι "μικρή" (όπως π.χ. συμβαίνει με την μέθοδο του Euler ή την μέθοδο του μέσου), έπεται ότι πρέπει να πάρουμε το  $h$  πολύ μικρό έτσι ώστε για όλες τις ιδιοτιμές  $\lambda_i$  (ακόμα και εκείνες για τις



οποίες  $|\operatorname{Re} \lambda_j| \gg 1$  και που δεν συμβάλλουν καθόλου εκθεδόν στην λύση!) να έχουμε  $h\lambda_j \in S$ . Άρα οι  $A$ -ευσταθείς μέθοδοι είναι ιδιαίτερα κατάλληλες για άκαμπτα προβλήματα που έχουν ιδιοτιμές  $\lambda_j$  οπουδήποτε στο  $\operatorname{Re} z < 0$ : για τις  $A$ -ευσταθείς μεθόδους έχουμε  $S = \{z : \operatorname{Re} z < 0\}$  και ευσεπώς  $h\lambda_j \in S$  για κάθε  $h > 0$ .

Πολλά άκαμπτα συστήματα προκύπτουν από την αριθμητική επίλυση μερικών διαφορικών εξισώσεων που η λύση τους εξαρτάται από τον χρόνο, όπως π.χ. παραβολικών και υπερβολικών εξισώσεων. θεωρούμε π.χ. το εξής απλό πρόβλημα αρχικών και ευνοριακών τιμών για την εξίσωση της θερμότητας: ζητάμε  $u(x,t)$  πραγματική συνάρτηση ορισμένη για  $x \in [0,1]$ ,  $t \in [0,T]$  τέτοια ώστε

$$(13) \quad \begin{cases} u_t = u_{xx}, & (x,t) \in [0,1] \times [0,T] \\ u(x,0) = u(x), & x \in [0,1] \\ u(0,t) = u(1,t) = 0, & t \in [0,T] \end{cases}$$

όπου  $u(x)$  δεδομένη πραγματική συνάρτηση στο  $[0,1]$  με  $u(0)=u(1)=0$ . Έστω  $x_j = j\Delta x$ ,  $j=0,1,2,\dots,(J+1)$  ένας ομοιόμορφος διαμερισμός του  $[0,1]$  με  $(J+1)\Delta x = x_{J+1} = 1$ . Αντικαθιστώντας στην Μ.Δ.Ε. του (13) την παράγωγο  $u_x(x,t)$  με την κεντρική διαφορά  $(u(x+\Delta x,t) - 2u(x,t) + u(x-\Delta x,t))/(\Delta x)^2$  παίρνουμε το εξής εύστημα Σ.Δ.Ε. για τις προσεγγίσεις  $u_j(t)$ , των τιμών  $u(x_j,t)$  της λύσης του (13):

$$(14) \quad \begin{cases} u_j'(t) = -(u_{j+1}(t) - 2u_j(t) + u_{j-1}(t))/(\Delta x)^2, & 1 \leq j \leq J, t \in [0,T] \\ u_0(t) = u_{J+1}(t) = 0, & t \in [0,T] \\ u_j(0) = u_j \equiv u(x_j), & 0 \leq j \leq J+1 \end{cases}$$

Το  $J \times J$  εύστημα των Δ.Ε. της (14) γράφεται για  $u = [u_1, \dots, u_J]^T$

$$(15) \quad u' = Au$$

όπου  $A$  ο συμμετρικός τριδιαγώνιος πίνακας με  $a_{ii} = -2/(\Delta x)^2$ ,

$a_{i,i+1} = (\Delta x)^{-2}$ . Δεν είναι δύσκολο να δούμε ότι οι ιδιοτιμές του  $A$  είναι όλες αρνητικοί αριθμοί (γιατί ο  $-A$  είναι θετικά ορισμένος) και δίνονται από τους τύπους

$$\lambda_j = (\Delta x)^{-2} [-2 + 2\cos(j\pi/(J+1))], \quad 1 \leq j \leq J$$

Είναι φανερό ότι  $\lambda_J < \lambda_{J-1} < \dots < \lambda_1 < 0$ . Μάλιστα  $\lambda_1 = (-2 + 2\cos(\pi/(J+1)))/(\Delta x)^2 = (-2 + 2\cos(\pi\Delta x))/(\Delta x)^2 = -\pi^2 + O((\Delta x)^2)$  ενώ  $\lambda_J = (-2 + 2\cos(J\pi/(J+1)))/(\Delta x)^2 = -4/(\Delta x)^2 + O(1)$ . Βλέπουμε δηλ. ότι για  $\Delta x$  μικρό

$$\lambda_J \cong -4(\Delta x)^{-2} \ll \lambda_1 \cong -\pi^2.$$

Προφανώς πρόκειται περί ακάμπτου ευστήματος με πραγματικές ιδιοτιμές.

Επειδή στις εφαρμογές συναντάμε συχνά προβλήματα (όπως το (14) π.χ.) με πραγματικές αρνητικές ιδιοτιμές ή ιδιοτιμές με  $\operatorname{Re}\lambda_j < 0$  που βρίσκονται όμως ε' ένα κώνο γωνίας  $0 < \theta < \pi/2$  στο  $\operatorname{Re}z < 0$ , δηλ. στο χωρίο  $S_\theta = \{z: z = re^{i\varphi}, \pi - \theta < \varphi < \pi + \theta, r > 0\}$ , δεν χρειάζεται να καταφεύγουμε πάντα σε μία  $A$ -ευσταθή μέθοδο, αλλά π.χ. σε μία μέθοδο της οποίας η περιοχή απόλυτης ευστάθειας περιλαμβάνει τον αρνητικό πραγματικό ημιάξονα ή το  $S_\theta$  αντίστοιχα. Λέμε ευγκεκριμένα ότι μία μέθοδος είναι  $A_\theta$ -ευσταθής (Cryer, 1973) αν η περιοχή της απόλυτης ευστάθειας της περιέχει τον αρνητικό πραγματικό άξονα  $\{z = x + iy, x < 0, y = 0\}$ . Μία μέθοδος λέγεται  $A(\theta)$ -ευσταθής (Widlund, 1967) αν η περιοχή απόλυτης ευστάθειας της περιέχει τον κώνο  $S_\theta$ . Προφανώς για επιτυχή αριθμητική λύση του ευστήματος (14) ((15)) αρκεί η μέθοδος να είναι  $A_\theta$ -ευσταθής οπότε για κάθε  $h > 0$  οι τιμές  $h\lambda_j$  θα βρεθούν μέσα στην περιοχή απόλυτης ευστάθειας της μεθόδου.

Θα εξετάσουμε τώρα λεπτομερέστερα θέματα απόλυτης ευστάθειας (για το μιγαδικό πρόβλημα (11)) των δύο μεγάλων κατηγοριών μεθόδων που έχουμε ήδη μελετήσει, δηλ. των μεθόδων RK και των (γραμμικών) πολυβηματικών μεθόδων. Αρχίζουμε με τις μεθόδους RK. Θεωρούμε την

γενική (πεπλεγμένη) μέθοδο RK με  $q$ -στάδια (3.2.12a,b): την εφαρμόζουμε στο πρόβλημα (11) και παίρνουμε

$$(16\alpha) \quad y^{n,i} = y^n + h\lambda \sum_{j=1}^q a_{ij} y^{n,j}, \quad 1 \leq i \leq q$$

$$(16\beta) \quad y^{n+1} = y^n + h\lambda \sum_{j=1}^q b_j y^{n,j}$$

Λύοντας ως προς  $y^{n,i}$  την σχέση (16α) (συμβολίζοντας με  $A_q = (a_{ij})$   $1 \leq i, j \leq q$  τον πίνακα των συντελεστών  $a_{ij}$  ώστε να μην γίνεται σύγχυση με τον πίνακα  $A$  του ευστήματος (10)) παίρνουμε, για  $Y^n = (y^{n,1}, \dots, y^{n,q})^T \in \mathbb{C}^q$ ,  $u = (1, \dots, 1)^T \in \mathbb{C}^q$ ,  $I_q$  ταυτότητα στον  $\mathbb{C}^q$ , ότι

$$Y^n = Y^n (I_q - h\lambda A_q)^{-1} u,$$

και, αντικαθιστώντας στην (16β), τελικά ότι

$$y^{n+1} = y^n (1 + h\lambda b^T (I_q - h\lambda A_q)^{-1} u).$$

Θέτοντας  $z = h\lambda \in \mathbb{C}$  ( $\operatorname{Re} z < 0$ ) μπορούμε να γράψουμε την σχέση αυτή στην μορφή

$$(17) \quad y^{n+1} = r(z) y^n, \quad z = h\lambda, \quad n \geq 0,$$

όπου  $r$  είναι η συνάρτηση

$$(18) \quad r(z) = 1 + z b^T (I_q - z A_q)^{-1} u, \quad \operatorname{Re} z < 0.$$

Υπό την προϋπόθεση ότι η  $r(z)$  είναι καλά ορισμένη για  $\operatorname{Re} z < 0$  (ή για κατάλληλο υποσύνολο του  $\operatorname{Re} z < 0$ ) - δηλ. ότι ο πίνακας  $I_q - z A_q$  είναι αντιστρέψιμος - βλέπουμε, γράφοντας τον αντίστροφο  $(I_q - z A_q)^{-1}$  συναρ-

τήσει οριζουμένου ότι η  $r(z)$  είναι ρητή συνάρτηση του  $z$  με βαθμούς αριθμητού και παρονομαστού το πολύ  $q$ .

Η έκθεση (17) είναι το διακριτό ανάλογο της

$$(19) \quad y(t^{n+1}) = e^z y(t^n), \quad z = h\lambda, \quad n \geq 0,$$

που ικανοποιεί η λύση  $y(t)$  του (11). Η συνάρτηση  $r(z)$  είναι ευνεπώς μία ρητή προσέγγιση του εκθετικού  $e^z$  για  $\operatorname{Re}z < 0$  και, ευνεπώς, οι ιδιότητες απόλυτης ευστάθειας της μεθόδου μεταφράζονται σε κατάλληλες ιδιότητες της  $r(z)$ . Είναι π.χ. προφανές από την (17), ότι μία μέθοδος RK είναι απόλυτα ευσταθής για κάποιο  $h$  αν το  $z = h\lambda$  είναι τέτοιο ώστε η  $r(z)$  να είναι καλά ορισμένη και να ικανοποιεί  $|r(z)| < 1$ . Συνεπώς η περιοχή απόλυτης ευστάθειας  $S$  μιάς μεθόδου RK μπορεί να οριστεί ως το σύνολο των  $z = h\lambda$ :

$$S = \{z : \operatorname{Re}z < 0, r(z) \text{ καλά ορισμένη και } |r(z)| < 1\}.$$

Μία μέθοδος είναι  $A$  - ευσταθής αν  $S = \{z : \operatorname{Re}z < 0\}$ ,  $A_0$  - ευσταθής αν  $S \supset \{z : \operatorname{Re}z < 0, \operatorname{Im}z = 0\}$ , κ.ο.κ.

Από τον ορισμό της τάξης ακρίβειας μιάς μεθόδου RK (βλ. Παρ. 3.2) συμπεραίνουμε ότι αν η μέθοδος (3.2.12α,β) έχει τάξη ακρίβειας  $p$ , τότε η ποσότητα  $y(t^{n+1}) - r(z)y(t^n) = (e^z - r(z))y(t^n)$  θα πρέπει να είναι τάξης  $O(z^{p+1})$  για κάθε  $\lambda, n$ . Άρα θα ικανοποιεί την αναγκαία συνθήκη

$$(20) \quad e^z - r(z) = O(z^{p+1}) \text{ για } z \rightarrow 0.$$

(Παρατηρείστε ότι η  $r(z)$  είναι αναλυτική σε μία περιοχή του μηδενός).

Οι ρητές συναρτήσεις που αντιστοιχούν στην μέθοδο του Euler, στην πεπλεγμένη Euler και στην μέθοδο του τραπεζίου είναι, αντίστοιχα,  $r_1(z) = 1+z$ ,  $r_2(z) = (1-z)^{-1}$ ,  $r_3(z) = (1+z/2)/(1-z/2)$ .

\* Δεν είναι όμως γενικά ευσταθής η (20) ευνεπώς ότι η τάξη ακρίβειας μιάς μεθόδου RK που αντιστοιχεί στην ρητή συνάρτηση  $r(z)$  είναι  $p$ . Η (20) είναι ικανή συνθήκη έτσι ώστε η τάξη ακρίβειας της μεθόδου για γραμμικά προβλήματα με σταθερούς συντελεστές να είναι  $p$ .

ικανοποιούν, αντίστοιχα,  $|r_i(z)| < 1$  για  $|z+1| < 1$  (περιοχή απόλυτης ευστάθειας της μεθόδου του Euler) και  $|r_i(z)| < 1$  για  $\operatorname{Re} z < 0$ , αν  $i=2,3$  (A-ευσταθείς μέθοδοι). (Παρατηρείτε ότι  $e^z = r_i(z) + O(z^2)$  αν  $i=1,2$  ενώ  $e^z = r_3(z) + O(z^3)$  σε συμφωνία με την (20) και τις γνωστές μας τάξεις των τριών μεθόδων). Είναι φανερό ότι οι άμεσες μέθοδοι RK έχουν πολυωνυμικές συναρτήσεις  $r(z)$  με  $r(0)=1$ . Αν μία άμεση μέθοδος έχει τάξη  $p$  η (20) δίνει ότι η αντίστοιχη  $r(z)$  της θα είναι αναγκαστικά

$$(21) \quad r(z) = 1 + z + z^2/2! + \dots + z^p/p!$$

Π.χ. όλες οι άμεσες μέθοδοι RK τάξης  $p=2$  έχουν  $r(z) = 1 + z + z^2/2$ . Είναι φανερό ότι δεν υπάρχει άμεση A-ευσταθής μέθοδος RK! (Οι περιοχές όπου  $|r(z)| < 1$  για άμεσες μεθόδους με  $p=1,2,3,4$  είναι τα γραμμοσκιασμένα χωρία του Σχήματος 1)

Από τις μεθόδους που ξεχωρίσαμε στην Παρ. 3.2, η πεπλεγμένη μέθοδος του μέσου (3.2.15') ευπύπτει, για το πρόβλημα (11), με την μέθοδο του τραπέζιου και είναι ευσεπώς A-ευσταθής με  $r(z) = (1+z/2)/(1-z/2)$ . (Δύο διαφορετικές μέθοδοι RK μπορεί να αντιστοιχούν στην ίδια ρητή προσέγγιση του εκθετικού). Οι ημιπεπλεγμένες μέθοδοι (3.2.16) με  $q=2$ ,  $p=3$ , δηλ. με  $\lambda = (1+3^{-1/2})/2$  δίνουν την ρητή συνάρτηση

$$(22) \quad r(z) = (1 + (1-2\lambda)z + (1/2 - 2\lambda + \lambda^2)z^2) / (1 - \lambda z)^2$$

και είναι A-ευσταθείς. A-ευσταθής είναι επίσης και η ημιπεπλεγμένη μέθοδος (3.2.51) με  $q=3$ ,  $p=4$ . Οι μέθοδοι "Gauss-Legendre  $q$  σημείων" (μέθοδος (3.2.19) αν  $q=2$ : βλ. Πρόταση 3.2.1 (γ) για τον ορισμό τους γενικά) είναι πολύ ενδιαφέρουσες όπως ξέρουμε γιατί δίνουν την μέγιστη τάξη ακρίβειας ( $p=2q$ ) για δεδομένο αριθμό σταδίων  $q$ . Η ρητή προσέγγιση που αντιστοιχεί στην (3.2.19) είναι η

$$(23) \quad r(z) = (1 + z/2 + z^2/12) / (1 - z/2 + z^2/12)$$

γιά την οποία ισχύει  $|r(z)| < 1$  για κάθε  $z: \operatorname{Re} z < 0$ , δηλ. ότι η αντίστοιχη μέθοδος είναι  $A$ -ευσταθής. Γενικά η μέθοδος Gauss-Legendre με  $q$  σημεία δίνει ως ρητή προσέγγιση του εκθετικού το  $q$ -στο διαχώνιο στοιχείο του πίνακα Padé για την συνάρτηση  $e^z$  (βλ. Παρατήρηση 2). Από γνωστή ιδιότητα του πίνακα Padé (Birkhoff-Ullmer, 1965) συμπεραίνουμε ότι όλες οι μέθοδοι Gauss-Legendre είναι  $A$ -ευσταθείς.

Εκτός από την περιπτώσιολογία υπάρχει φυσικά και πολλή θεωρία για τις ιδιότητες των ρητών προσεγγίσεων του εκθετικού  $e^z$  και ανάλογη σημαντική θεωρία χαρακτηρισμένων ρητών συναρτήσεων  $r(z)$  διαφόρων μορφών που οδηγούν σε  $A$ -ευσταθείς μεθόδους. Π.χ. ένα απλό αποτέλεσμα είναι το εξής:

**ΠΡΟΤΑΣΗ 1** Αν οι ιδιοτιμές  $\mu_i$  του πίνακα  $A_q$  ικανοποιούν  $\operatorname{Re} \mu_i \geq 0, |\mu_i| \leq q$ , αν η μέθοδος είναι ευγενής και αν για κάθε  $y \in \mathbb{R}$  έχουμε  $|r(iy)| \leq 1$ , τότε  $|r(z)| < 1$  για  $\operatorname{Re} z < 0$ , δηλ. η αντίστοιχη μέθοδος είναι  $A$ -ευσταθής.

Απόδειξη: Επειδή οι ιδιοτιμές  $\mu_i, |\mu_i| \leq q$  του  $A_q$  ικανοποιούν  $\operatorname{Re} \mu_i \geq 0 \quad \forall i$  και επειδή οι ιδιοτιμές του πίνακα  $I_q - zA_q$  είναι οι  $1 - z\mu_i, |\mu_i| \leq q$ , συμπεραίνουμε ότι ο  $I_q - zA_q$  είναι αντιστρέψιμος για κάθε  $z: \operatorname{Re} z < 0$ . (Πράγματι αν για κάποιο  $\mu = \mu_i, \mu \neq 1/z$ , έχουμε ότι  $\operatorname{Re} \mu = \operatorname{Re} z / |z|^2$ , άτοπο). Συνεπώς, η ρητή συνάρτηση  $r(z)$ , βλ. (18), δεν έχει πόλους στο ημιεπίπεδο  $\operatorname{Re} z < 0$ . Από την υπόθεσή μας ότι  $|r(iy)| \leq 1, \forall y \in \mathbb{R}$  συμπεραίνουμε ότι δεν έχει πόλους ούτε στον φανταστικό άξονα. Άρα η  $r(z)$  είναι αναλυτική σε μία περιοχή του κλειστού ημιεπιπέδου  $\operatorname{Re} z \leq 0$ . Επειδή η  $r(z)$  είναι ρητή και ισχύει  $|r(iy)| \leq 1 \quad \forall y \in \mathbb{R}$  συμπεραίνουμε ότι  $\lim_{|z| \rightarrow \infty} |r(z)| \leq 1$ . Άρα από την αρχή του μεγίστου για αναλυτικές συναρτήσεις συμπεραίνουμε ότι  $|r(z)| \leq 1 \quad \forall z: \operatorname{Re} z \leq 0$ . Επειδή η μέθοδος είναι ευγενής ( $p \geq 1$ ) η (20) δίνει ότι  $r(z) \neq 1$  για  $\operatorname{Re} z \leq 0$ . Συμπεραίνουμε, πάλι απ' την αρχή του μεγίστου, ότι, ακριβέστερα,  $|r(z)| < 1$  στο εσωτερικό  $\operatorname{Re} z < 0$ . @

Το αποτέλεσμα αυτό μας επιτρέπει, αν μπορεί να εφαρμοσθεί, να ελέγχουμε αν  $|r(z)| \leq 1$  μόνο για  $z=iy$ ,  $y \in \mathbb{R}$ . Επίσης από την απόδειξη του συμπεραίνουμε ότι αν η μέθοδος είναι ευνηής (αρκεί  $r(z) \neq 1$ ) και η  $r(z)$  είναι καλά ορισμένη για  $\operatorname{Re} z \leq 0$ , τότε  $\{|r(z)| < 1, \operatorname{Re} z < 0\} \Leftrightarrow \{|r(z)| \leq 1, \operatorname{Re} z \leq 0\}$ . Άλλα μία ενδιαφέρουσα ικανή συνθήκη είναι η εξής: Αν οι ιδιοτιμές  $\mu_i$  του  $A_q$  ικανοποιούν  $\operatorname{Re} \mu_i \geq 0 \quad \forall i$ , αν η τάξη  $p$  της μεθόδου ικανοποιεί την ανισότητα  $p \geq 2q-2$  και αν  $\lim_{x \rightarrow \infty} |r(x)| \leq 1$ , τότε η μέθοδος είναι  $A$ -ευσταθής.

Όπως αναφέρθηκε προηγουμένως υπάρχει σημαντική θεωρία για την μελέτη των ιδιοτήτων (ιδίως ευστάθειας, ακρίβειας, εντοπισμού των πόλων κ.λ.π.) ρητών προσεγγίσεων της  $e^z$  που βασίζεται συνήθως σε προχωρημένα θέματα της θεωρίας μιγαδικών συναρτήσεων. Για μία γεύση των μεθόδων που χρησιμοποιούνται σήμερα στις αποδείξεις βλ. π.χ. το άρθρο των G. Wanner, E. Hairer και S.P. Nørsett, "Order stars and stability theorems", BIT 18(1978), 475-89.

Προχωρούμε τώρα στην μελέτη της απόλυτης ευστάθειας (για την μιγαδική απλή Δ.Ε. (11)) των (γραμμικών) πολυβηματικών μεθόδων - που προηγείται ιστορικά (βλ. Dahlquist, 1963) της ανάλογης μελέτης για τις μεθόδους RK -. Εφαρμόζοντας την  $k$ -βηματική μέθοδο (3.3.3) στην εξίσωση (11) παίρνουμε την ομογενή εξίσωση διαφορών

$$(24) \quad \sum_{j=0}^k (a_j - h\lambda\beta_j) y^{n+j} = 0, \quad n \geq 0$$

Για να ισχύει  $y^n \rightarrow 0$ ,  $n \rightarrow \infty$  για κάποιο  $h > 0$  για κάθε λύση της (24) με  $y^0 = 1$  πρέπει ευνηώς (βάσει της θεωρίας των λύσεων ομογενών εξισώσεων διαφορών με σταθερούς συντελεστές, βλ. Παρ. 3.3., ιδίως (3.3.1), (3.3.12)) οι ρίζες  $\zeta_i = \zeta_i(\lambda h)$ ,  $1 \leq i \leq k$ , του πολυωνύμου

$$(25) \quad p(\zeta, \lambda h) = p(\zeta) - \lambda h c(\zeta), \quad \operatorname{Re} \lambda < 0$$

να ικανοποιούν

$$(26) |\zeta_i| < 1, 1 \leq i \leq k.$$

(Τονίζουμε την απειρή ανισότητα στην (26)). Το πρόβλημα ευνηώς του προσδιορισμού της απόλυτης ευσταθείας της μεθόδου (3.3.3) ανάγεται στην μελέτη των ριζών  $\zeta_i$  του  $n$ -ως συναρτήσεων της μιγαδικής παραμέτρου  $h\lambda$ ,  $\operatorname{Re} \lambda < 0$  και στην διατύπωση συνθηκών έτσι ώστε να ιεχύει η (26).

Ας δούμε μερικά παραδείγματα. Για την μέθοδο (3.1.1') έχουμε  $\rho(\zeta) = \zeta^2 - 1$ ,  $\epsilon(\zeta) = 2\zeta$  ευνηώς  $n = \zeta^2 - 2h\lambda\zeta - 1$ , του οποίου οι ρίζες  $\zeta_1, \zeta_2$  ικανοποιούν, για κάθε  $h\lambda \in \mathbb{C}$ , την συνθήκη  $|\zeta_1| |\zeta_2| = 1$ . Συνεπώς δεν είναι δυνατόν να ιεχύει η (26) για καμιά τιμή του  $h\lambda$ : η μέθοδος (3.1.1') - αν και "ευσταθής" σύμφωνα με τον ορισμό της Παρ. 3.3 - δεν είναι ποτέ απόλυτα ευσταθής. Το ίδιο συμβαίνει και με την μέθοδο του Simpson (3.3.2). Έχουμε ήδη δεί ότι οι μέθοδοι πεπλεγμένη Euler και τραπεζίου είναι A-ευσταθείς. Οι μέθοδοι "οπισθοδρομικών διαφορών" με  $k$  βήματα (3.3.6),  $1 \leq k \leq 6$  έχουν τις εξής ιδιότητες: Για  $k=1$  (πεπλεγμένη Euler) και  $k=2$  είναι A-ευσταθείς. Για  $3 \leq k \leq 6$  είναι A-ευσταθείς με γωνίες  $\theta_k$ ,  $3 \leq k \leq 6$  τις εξής:  $\theta_3 \approx 88^\circ$ ,  $\theta_4 \approx 73^\circ$ ,  $\theta_5 \approx 51^\circ$ ,  $\theta_6 \approx 18^\circ$ . Είναι δηλ. κατάλληλες για άκαμπτα ευστήματα με ιδιοτιμές που βρίσκονται μέσα στους αντίστοιχους κώνους  $S_{\theta_k}$ . Θα δούμε στις Ασκήσεις και άλλα παραδείγματα. Βλ. και Παρατήρηση 1.

Υπάρχει σημαντικό ερώμα θεωρίας για τις ιδιότητες απόλυτης ευσταθείας των (γραμμικών) πολυβηματικών μεθόδων, που αρχίζει με την εργασία του G.O. Dahlquist "A special stability problem for linear multistep methods", BIT 3(1963), 27-43. Ο Dahlquist (αφού δίνει τον ορισμό της απόλυτης ευσταθείας) αποδεικνύει ότι δεν υπάρχουν άμεσες ( $\beta_k = 0$ ) A-ευσταθείς πολυβηματικές μέθοδοι και ότι η τάξη ακρίβειας  $p$  μιάς A-ευσταθούς πολυβηματικής μεθόδου δεν μπορεί να είναι μεγαλύτερη του 2. Πάλι στα, η μέθοδος του τραπεζίου είναι εκείνη η A-ευσταθής μέθοδος με  $p=2$  με την μικρότερη "εσταθερά εφάλματος" (βλ. Παρατήρηση 3.3.2)  $c_* = c_3/\epsilon(1) = 1/12$ . Έτσι έχουμε ένα εσοφόρο περιορισμό στην τάξη ακρίβειας μιάς A-ευσταθούς πολυβηματικής μεθό-



δου, πράγμα που αναδεικνύει την σημασία των (πεπλεγμένων) μεθόδων AK μεταξύ των οποίων υπάρχουν A-ευστάθεις μέθοδοι οποιασδήποτε τάξης ακρίβειας (π.χ. οι μέθοδοι Gauss-Legendre με  $q$  ετάδια και  $p=2q$ ).

Η κατάσταση όμως είναι πολύ καλύτερη όταν δεν επιβάλλουμε A-ευστάθεια αλλά αρκεσθούμε σε A(θ)-ευστάθεια για  $0 < \theta < \pi/2$  ή  $A_0$ -ευστάθεια. (Για την αριθμητική λύση με πλήρως διακριτές μεθόδους π.χ. προβλημάτων παραβολικών ΜΔΕ, συνήθως  $A_0$  - ή A(θ) - ευστάθεια είναι υπερφορτωτές. Δεν ισχύει όμως αυτό για υπερβολικά π.χ. προβλήματα όπου οι ιδιοτιμές του A είναι φανταστικές βλ. και AK.7 ). Είναι γνωστό (Widlund, 1967) ότι δεν υπάρχουν άμεσες A(θ)-ευστάθεις μέθοδοι και ότι η μόνη A(θ)-ευστάθης k-βηματική μέθοδος της οποίας η τάξη  $p$  υπερβαίνει το  $k$  είναι η μέθοδος του τραπεζίου ( $k=1, p=2$ ). Αλλά είναι γνωστό ότι για κάθε  $\theta \in (0, \pi/2)$  υπάρχουν A(θ)-ευστάθεις μέθοδοι με  $k=p=3$  και  $k=p=4$  και επιπλέον ότι για τουλάχιστον  $1 \leq k \leq 6$  υπάρχουν A(θ<sub>k</sub>)-ευστάθεις μέθοδοι,  $0 < \theta_k < \pi/2$ , με  $p=k$  (οι μέθοδοι "οπισθοδρομικών διαφορών"). Αν περιοριστούμε σε  $A_0$ -ευστάθεις μεθόδους, είναι γνωστό (Cryer, 1973) ότι ισχύει πάλι ότι οι  $A_0$ -ευστάθεις μέθοδοι είναι πεπλεγμένες και ότι η τάξη τους  $p$  δεν μπορεί να υπερβαίνει το  $k$  με εξαίρεση την μέθοδο του τραπεζίου. Όμως για κάθε αριθμό βημάτων  $k$  υπάρχουν  $A_0$ -ευστάθεις μέθοδοι τάξης  $p=k$ . Για μία πολύ καλή ανασκόπηση των παραπάνω (αλλά και πολλών άλλων παρόμοιων) κλασικών αποτελεσμάτων για την απόλυτη ευστάθεια των πολυβηματικών μεθόδων και για τις αποδείξεις τους με ένα εννοποιημένο και ετοιχειώδη τρόπο βλ. την παράγραφο 3.2. του β' τόμου του βιβλίου [3.3] του Grigorieff. Για μία πρόσφατη ανασκόπηση νεωτέρων αποτελεσμάτων με πιο προχωρημένες μεθόδους της θεωρίας μιγαδικών συναρτήσεων, βλ. π.χ. το άρθρο των Jeltsch και Neunlinna, Numer. Math. 40(1982), 245-296.

Η μελέτη της απόλυτης ευστάθειας των μεθόδων AK και των πολυβηματικών μεθόδων έγινε με βάση την συμπεριφορά αυτών των μεθόδων όταν εφαρμοσθούν στην απλή μιγαδική Δ.Ε. (11). Ας δούμε τώρα την σημασία αυτών των ιδιοτήτων για την αριθμητική λύση του πχππ πραγματικού ευστήματος (10) του οποίου το μη ομογενές ανάλογο γράψουμε ως

$$(27) \begin{cases} y' = My + f(t), t \geq 0, \\ y(0) = y_0, \end{cases}$$

και για το οποίο, γενικεύοντας λίγο, υποθέτουμε ότι για όλες τις ιδιοτιμές  $\lambda_i$  του πίνακα  $M \in \mathbb{R}^{m \times m}$  ισχύει  $\operatorname{Re} \lambda_i \leq 0$  και ότι υπάρχουν ιδιοτιμές τέτοιες ώστε π.χ.  $|\operatorname{Re} \lambda_i| = 0(1)$  ενώ για άλλες  $|\operatorname{Re} \lambda_i| \gg 1$ , δηλ. ότι το (27) είναι άκαμπτα. (Γράφουμε  $M$  αντί  $A$  ώστε να μην υπάρξει σύγχυση μεταξύ των στοιχείων του  $M$  και των συντελεστών  $a_{ij}$  της μεθόδου RK).  
 Ας εξετάσουμε πρώτα την ευστάθεια των μεθόδων RK για το (27). Η μέθοδος (3.2.12a,b) στην περίπτωση του (27) γίνεται (για  $n \geq 0$ ,  $y^0 = y_0$ ):

$$(28) \begin{cases} y^{n,i} = y^n + h \sum_{j=1}^q a_{ij} (My^{n,j} + f(t^{n,j})), 1 \leq i \leq q, \\ y^{n+1} = y^n + h \sum_{j=1}^q b_j (My^{n,j} + f(t^{n,j})). \end{cases}$$

Υποθέτουμε ότι το γραμμικό σύστημα (ως προς  $y^{n,i}$ ) που ορίζεται μέσω της (28) έχει μοναδική λύση για κάθε  $n$ . (Αυτό π.χ. εξασφαλίζεται, χωρίς περιορισμό στο  $h$ , αν οι ιδιοτιμές  $\mu_i$  του πίνακα  $A_q = (a_{ij})$  έχουν

θετικό πραγματικό μέρος, βλ. Πρόταση 1). Σε αναλογία με ό,τι κάναμε στην Παράγραφο 3.2 (Πρόταση 3.2.2) θεωρούμε την "διαταραχή" του

(28) με δεδομένα  $z^0 \in \mathbb{R}^m$ ,  $e^n \in \mathbb{R}^m$ ,  $n \geq 0$ :

$$(29) \begin{cases} z^{n,i} = z^n + h \sum_{j=1}^q a_{ij} (Mz^{n,j} + f(t^{n,j})), 1 \leq i \leq q, \\ z^{n+1} = z^n + h \sum_{j=1}^q a_{ij} (Mz^{n,j} + f(t^{n,j})) + e^n. \end{cases}$$

Για τεχνικούς λόγους θα μελετήσουμε την ευστάθεια της (28) ως προς την ευκλείδεια νόρμα  $\|\cdot\| = \|\cdot\|_2$  στον  $\mathbb{R}^m$ . Σκοπός μας είναι να φράξουμε

## 3.4.21

την διαφορά  $\|y^n - z^n\|$  συναρτήσει των  $\|y^0 - z^0\|$ ,  $\max_{0 \leq j \leq n-1} \|e^j\|$  χωρίς όμως να χρησιμοποιήσουμε το αποτέλεσμα (3.2.25) της Πρότασης 3.2.2., το οποίο φυσικά ισχύει και εδώ αλλά με σταθερά Lipschitz  $L = \|M\| = \max_{1 \leq i \leq m} |\lambda_i(M^*M)|^{1/2}$  που μπορεί να γίνει πολύ μεγάλη (π.χ. θεωρείτε  $M^* = M$  οπότε  $\lambda_i(M^*M) = \lambda_i^2$ ) λόγω της ακαμψίας του (27). Συνεπώς, στην περίπτωση μας, οι σταθερές  $C_1$  και  $C_2$  της (3.2.25) είναι πολύ μεγάλες, η δε (3.2.25) πρακτικά άχρηστη.

Έστω ότι η μέθοδος RK είναι A-ευσταθής. Τότε ισχύει, για κάθε  $n \geq 0$  ότι για κάθε  $h > 0$

$$(30) \quad \|y^n - z^n\| \leq \|y^0 - z^0\| + n \max_{0 \leq j \leq n-1} \|e^j\|.$$

Η απόδειξη της (30) για γενικούς πίνακες  $M$  με  $\operatorname{Re} \lambda_i \leq 0$  είναι τεχνικά όχι τόσο απλή (Crouzeix 1975). Για συμμετρικούς πίνακες  $M$  η απόδειξη είναι πολύ ευκολότερη: έστω λοιπόν  $M$  συμμετρικός  $m \times m$  πίνακας με ιδιοτιμές  $\lambda_i \leq 0$ ,  $1 \leq i \leq m$ . Στην ειδική αυτή περίπτωση θα αποδείξουμε την (30) για A<sub>0</sub>-ευσταθείς μεθόδους, για μεθόδους δηλ. για τις οποίες ισχύει ότι  $|r(x)| < 1$  για  $x < 0$ , όπου  $r(x)$  η ρητή προέκταση του  $e^x$  για  $x < 0$  που δίνεται από την (18). Υποθέτοντας ότι η μέθοδος είναι τουλάχιστον συνεπής (ακριβέστερα ότι  $r(x) \neq 1$ ) αυτό είναι ισοδύναμο με την ασθενέστερη ευαθήκη ότι

$$(31) \quad |r(x)| \leq 1 \text{ για } x \leq 0,$$

την οποία και υποθέτουμε ότι ισχύει. Χρησιμοποιούμε τώρα την φασματική παράσταση του πίνακα  $M$ , δηλ. ότι ο  $M$  έχει  $m$  ορθοκανονικά ιδιοδιανύσματα  $u^i$  τέτοια ώστε  $Mu^i = \lambda_i u^i$ . Τότε, για κάθε πραγματικό

πολυώνυμο  $p$  έχουμε  $\forall u \in \mathbb{R}^m \quad p(M)u = \sum_{i=1}^m p(\lambda_i)(u, u^i)u^i$ , όπου  $(\cdot, \cdot)$  το ευκλείδειο εσωτερικό γινόμενο στον  $\mathbb{R}^m$ . Επίσης, αν ο πίνακας  $\alpha I + \beta M$

είναι αντιετρέφιος, έχουμε ότι  $(\alpha I + \beta \Pi)^{-1} u = \sum_{i=1}^m (\alpha + \beta \lambda_i)^{-1} (u, u^i) u^i$ .

Συμπεραίνουμε ότι για κάθε ρητή συνάρτηση  $r(x)$ , καλά ορισμένη για  $x \in D$  μπορούμε να ορίσουμε για κάθε  $u \in \mathbb{R}^m$

$$r(\Pi)u = \sum_{i=1}^m r(\lambda_i) (u, u^i) u^i.$$

Αφαιρώντας κατά μέλη τις (28) και (29) έχουμε τώρα ότι

$$y^{n,i} - z^{n,i} = y^n - z^n + h \sum_{j=1}^m a_{ij} \Pi(y^{n,j} - z^{n,j})$$

$$y^{n+1} - z^{n+1} = y^n - z^n + h \sum_{j=1}^m b_j \Pi(y^{n,j} - z^{n,j}) - e^n.$$

Εκφράζοντας τα διανύσματα  $y^n, y^{n,i}, z^n, z^{n,i}$  ως αναπτύγματα ιδιοδιανυσμάτων του  $\Pi$  και κάνοντας λίγες πράξεις δεν είναι δύσκολο να δούμε ότι, σε αναλογία με την (17) π.χ., τώρα ισχύει

$$y^{n+1} - z^{n+1} = r(h\Pi)(y^n - z^n) - e^n, \quad n \geq 0,$$

από την οποία προκύπτει με επανάληψη η σχέση

$$(33) \quad \|y^n - z^n\| \leq \|r(h\Pi)\|^n \|y^0 - z^0\| + \sum_{j=0}^{n-1} \|r(h\Pi)\|^{n-1-j} \|e^j\|, \quad n \geq 1.$$

Τώρα, για κάθε συνάρτηση  $\varphi(x)$  ορισμένη για  $x = \lambda_i, 1 \leq i \leq m$ , θεωρούμε τον

αντίστοιχο γραμμικό τελεστή  $\varphi(\Pi)$  που ορίζεται για  $u \in \mathbb{R}^m$  ως

$$(34) \quad \varphi(\Pi)u = \sum_{i=1}^m \varphi(\lambda_i) (u, u^i) u^i.$$

Έχουμε τότε ότι λόγω της ορθογωνιότητας των  $u^i$

$$\|\varphi(\Pi)u\|^2 = \sum_{i=1}^m [\varphi(\lambda_i)]^2 (u, u^i)^2 \leq \max_{1 \leq i \leq m} [\varphi(\lambda_i)]^2 \|u\|^2$$

Συνεπώς ισχύει  $\|\varphi(M)u\| \leq \max_{1 \leq i \leq m} |\varphi(\lambda_i)| \|u\|$  για κάθε  $u \in \mathbb{R}^m$ , με  
 ιδιότητα αν  $u = u^e$  όπου  $|\varphi(\lambda_i)| = \max_i |\varphi(\lambda_i)|$ . Συνεπώς αν η  $\varphi(M)$  ορίζεται  
 από την (34) έχουμε

$$(35) \|\varphi(M)\| = \max_{1 \leq i \leq m} |\varphi(\lambda_i)|.$$

Συμπεραίνουμε, λόγω της (31), επειδή  $\lambda_i \leq 0$ ,

$$(36) \|r(hM)\| = \max_{1 \leq i \leq m} |r(h\lambda_i)| \leq 1 \quad \forall h > 0.$$

Η (36) και η (33) δίνουν συνεπώς την ανισότητα (30), η οποία  
 πραγματικά εκφράζει την ευστάθεια της μεθόδου, κάτω από τις  
 προϋποθέσεις μας, στην περίπτωση του γραμμικού ευστήματος (27).

Τι μπορούμε να πούμε τώρα για μεθόδους που δεν είναι  
 $R$ -ευσταθείς, π.χ. για μεθόδους, όπως οι άμεσες RK, για τις οποίες η  
 συνθήκη  $|r(z)| \leq 1$  ισχύει μόνο για ένα φραγμένο υποσύνολο του  
 ημιεπιπέδου  $\operatorname{Re} z \leq 0$ ; Ένα γενικό αποτέλεσμα δίνεται από τον Crouzeix,  
 op. cit. Εδώ ως περιοριστούμε πάλι σε άκαμπτα ευστήματα με  
 συμμετρικούς πίνακες  $M$  με  $\lambda_i \leq 0$ ,  $1 \leq i \leq m$ . Υποθέτουμε ότι τώρα ισχύει η  
 $|r(x)| \leq 1$  μόνο ε' ένα διάστημα δηλ. ότι αντί της (31) έχουμε για  
 κάποιο  $a < 0$  ότι

$$(37) |r(x)| \leq 1, \quad x \in (a, 0].$$

(Π.χ. για την μέθοδο του Euler ή την μέθοδο του μέσου  $\sigma = -2$  κλπ.).  
 Τότε η (36) δεν ισχύει για κάθε  $h > 0$  αλλά για περιορισμένα  $h$ .  
 Πράγματι, λόγω της (37), μόνο αν  $h\lambda_i \in (a, 0]$ ,  $1 \leq i \leq m$ , δηλ. αν

$$(38) h \max_{1 \leq i \leq m} |\lambda_i| < |a|,$$

θα έχουμε ότι

$$(39) \|r(hM)\| = \max_{1 \leq i \leq m} |r(h\lambda_i)| \leq 1.$$

Συνοπώς αν ισχύει η (37), τότε για μικρό  $h$  (έτσι ώστε να ικανοποιείται η (38), που για άκαμπτο σύστημα αποτελεί ευνήθως σοβαρό περιορισμό στο  $h$ ) προκύπτει πάλι η (30). (Γιαυτόν τον λόγο οι  $A$ -,  $A(\theta)$ - ή  $A_0$ -ευσταθείς μέθοδοι λέγονται μερικές φορές και "απεριόριστα ευσταθείς" ενώ μία μέθοδος για την οποία ισχύει η (37) λέγεται και "ευσταθής υπό ευνήθη" - την (38)).

Η σχέση (3.2.25) ήταν το "κλειδί" στην απόδειξη της εκτίμησης (3.2.31) για την τάξη του σφάλματος  $\|y^n - y(t^n)\|$  των μεθόδων RK υπό την προϋπόθεση ότι είναι γνωστή μία σχέση όπως η (3.2.30) όπου  $\theta^n = e^n$ . Αν υποθέσουμε και εδώ ότι η τάξη ακρίβειας της μεθόδου μας είναι  $p$ , τότε ισχύει εξ ορισμού σχεδόν μία σχέση της μορφής (3.2.30). Όμως η σταθερά  $D$  θα είναι τυπικά πολυωνυμική συνάρτηση νορμών υψηλών παραχάχων  $y^{(j)}$  της λύσης του συστήματος (37), δηλ. υψηλών δυνάμεων του πίνακα  $M$ , δηλ. υψηλών δυνάμεων των ιδιοτιμών  $\lambda_i$  του  $M$ , οπότε η σταθερά  $C$  της (3.2.31) θα είναι πολύ μεγάλη και πρακτικά άχρηστη για άκαμπτο σύστημα. Η απόδειξη φραχμάτων με άριστη τάξη ακρίβειας  $O(h^p)$  αλλά με "λογικές" σταθερές  $C$  είναι πρόβλημα που δεν μπορεί ευνήθως να επιλυθεί στο επίπεδο των ευνήθων διαφορικών εξισώσεων. Χρειάζεται να εισάγουμε ειδικές νόρμες και να χρησιμοποιήσουμε την ομαλότητα της λύσης του προβλήματος (π.χ. της ΜΔΕ) από το οποίο προέρχεται το σύστημά μας. Σχετικά εύκολα όμως μπορεί κανείς να αποδείξει εκτιμήσεις της μορφής

$$\|y^n - y(t^n)\| \leq Ch^p (t^n)^{-p} \|y^0\|, \quad n \geq 1$$

για μία  $A_0$ -ευσταθή μέθοδο π.χ. με  $|r(x)| < 1$  για  $x < 0$ ,  $\lim_{|x| \rightarrow \infty} |r(x)| < 1$ ,  $M$  συμμετρικό με  $\lambda_i \leq 0$  και όπου  $\|\cdot\| = \|\cdot\|_2$ : βλ. π.χ. LeRoux, Math. Comp. 33(1979), 919-931 και Baker, Bramble και Thomée, Math. Comp. 31(1977), 818-847. Εκτιμήσεις τέτοιου τύπου είναι ουσιαστικά εκτιμήσεις του πόσο γρήγορα φθίνει το σφάλμα (εδώ όπως το  $n^{-p}$ ) για μεγάλο  $n$ . Το  $C$  είναι ανεξάρτητο των  $h, n$  και  $y(t)$  (και των  $\lambda_i$ ).

Στην περίπτωση πολυβηματικών μεθόδων υπάρχει επίσης ανάλογη θεωρία (Zlamal 1975, Crouzeix-Raviart, 1978). Για το γενικό πρόβλημα (27) με  $\operatorname{Re} \lambda_i \leq 0$  ισχύει ότι αν η μέθοδος είναι  $A$ -ευσταθής και έχει τάξη

$p$  ( $p \leq 2$ ) τότε, για κάποια σταθερά  $C$  ανεξάρτητη των  $n, h, y(t), \lambda_i$ , έχουμε:

$$\|y(t^n) - y^n\| \leq C(\max_{0 \leq i \leq k-1} \|y(t^i) - y^i\| + h^p \int_0^{t^n} \|y^{(p+1)}(t)\| dt).$$

Παρόμοια εκτίμηση ισχύει και για  $A_0$ -ευσταθείς μεθόδους όταν ο  $M$  είναι συμμετρικός με  $\lambda_i \leq 0$  οπότε μπορούμε να επιτύχουμε και τάξη ακρίβειας  $p > 2$ . Βέβαια πάλι  $\|y^{(p+1)}(t)\| = O(\max_i |\lambda_i|^{p+1})$ .

Τα τελευταία 10-15 χρόνια οι προσπάθειες στην περιοχή της απόλυτης ευστάθειας έχουν στραφεί σε γενικεύσεις της έννοιας για εφαρμογή σε άκαμπτα μη γραμμικά προβλήματα ή σε γραμμικά άκαμπτα προβλήματα με μεταβλητούς συντελεστές. Τα προβλήματα αυτά έχουν γενικά λύσεις που φθίνουν καθώς αυξάνει το  $t$  αλλά έχουν πολύ μεγάλες σταθερές Lipschitz έτσι ώστε η θεωρία ευστάθειας των Παρ. 3.2. και 3.3 να μην είναι πρακτικά εφαρμόσιμη. Η έρευνα έχει αποδειχθεί πολύ καρποφόρα: έχουν ήδη εντοπισθεί και χαρακτηριστεί οικογένειες αποτελεσματικών και πρακτικά ευσταθών μεθόδων ικανοποιητικής ακρίβειας ενώ έχουν κατανοηθεί ε' ένα βαθμό μηχανισμοί "αποσταθεροποίησης" ακόμα και  $A$ -ευσταθών (για γραμμικά προβλήματα) μεθόδων σε άκαμπτα μη γραμμικά (ή γραμμικά με μεταβλητούς συντελεστές) προβλήματα. Η έρευνα συνεχίζεται: μιά, κάπως ασύνδετη, εικόνα της κατάστασης για μεθόδους RK μέχρι το 1984 δίνει το βιβλίο [3.1]. Έδώ θα περιοριστούμε σε μερικές βασικές παρατηρήσεις, μόνο για μεθόδους RK\*.

θα παραλείψουμε το ενδιαμέσο στάδιο των γραμμικών προβλημάτων με μεταβλητούς συντελεστές. θεωρούμε το μιγαδικό σύστημα των μη γραμμικών Δ.Ε.

$$(40) \quad \begin{cases} y' = f(t, y), & t \geq 0, \\ y(0) = y_0. \end{cases}$$

\* Βλ. π.χ. εργασίες των Butcher, BIT 15(1975) 358-361, Crouzeix, Num. Math. 32(1979), 75-82, Burrage-Butcher SIAM J.N.A. 16(1979), 46-57.

όπου η  $y(t)$  και η  $f$  έχουν τιμές στον  $\mathbb{C}^m$ . Υποθέτουμε ότι το (40) έχει μοναδική λύση (τουλάχιστον τοπικά, ε' ένα διάστημα  $[0, T]$ ) και ότι ισχύει η σχέση

$$(41) \operatorname{Re}\{(f(t, y) - f(t, z), y - z)\} \leq 0 \quad \forall t \geq 0, y, z \in \mathbb{C}^m,$$

όπου τώρα με  $(x, y)$  συμβολίζουμε το ευκλείδιο εσωτερικό γινόμενο

$\sum_{i=1}^m x_i \bar{y}_i$  στον  $\mathbb{C}^m$  (και με  $\|\cdot\|$  την αντίστοιχη νόρμα). Για δύο λύσεις  $y(t), z(t)$  του (40) ε' ένα διάστημα  $[t', t'']$  έχουμε

$$y'(t) - z'(t) = f(t, y(t)) - f(t, z(t)), \quad t' \leq t \leq t''.$$

Συνεπώς, παίρνοντας το εσωτερικό γινόμενο με την  $y(t) - z(t)$  έχουμε, επειδή

$$(y' - z', y - z) = (d/dt) \|y - z\|^2 - (y - z, y' - z'),$$

δηλ. επειδή

$$2 \operatorname{Re}\{(y' - z', y - z)\} = (d/dt) \|y - z\|^2,$$

ότι λόγω της (41),

$$(d/dt) \|y - z\|^2 = 2 \operatorname{Re}\{(f(t, y) - f(t, z), y - z)\} \leq 0, \quad t' \leq t \leq t''$$

και συνεπώς ότι

$$(42) \|y(t'') - z(t'')\| \leq \|y(t') - z(t')\|, \quad t' \leq t''.$$

Σημειώστε ότι η σταθερά Lipschitz της  $f$  (αν η  $f$  είναι Lipschitz) δεν εμφανίζεται στην ανισότητα!

Λέμε ότι η μέθοδος RK της μορφής (3.2.12a, b) είναι B-ευσταθής (Butcher, Crouzeix), αν, όταν εφαρμοσθεί στο πρόβλημα (40) (με την ιδιότητα (41)), ικανοποιεί το διακριτό ανάλογο της (42), δηλ. αν για



οποιοδήποτε δύο λύσεις  $y^n, z^n, y^{n+1}, z^{n+1}$  των (3.2.12a,b) ιαχύει

$$(43) \|y^{n+1} - z^{n+1}\| \leq \|y^n - z^n\|, \quad n \geq 0,$$

οπότε πραγματικά έχουμε ευεστάθεια των λύσεων του (3.2.12a,b). Σημειώστε ότι μιά B-ευσταθής μέθοδος εφαρμοζόμενη σε μία μιγαδική εξίσωση της μορφής  $y' = \lambda y$ ,  $-\operatorname{Re} \lambda \leq 0$ , δίνει, επειδή ιαχύει η (41),  $|y^{n+1} - z^{n+1}| \leq |y^n - z^n|$ , δηλ. (με  $z^n = 0$ ),  $|y^{n+1}| \leq |y^n| \Rightarrow |\operatorname{r}(h\lambda)| \leq 1 \quad \forall h > 0$ ,  $\lambda: \operatorname{Re} \lambda \leq 0$ , δηλ. ότι η μέθοδος - υποθέτουμε ότι  $r(z) \neq 1$  - είναι B-ευσταθής. Άρα B-ευστάθεια  $\Rightarrow$  A-ευστάθεια.

Προχωρούμε τώρα σε ικανές και αναγκαίες συνθήκες ώστε η μέθοδος RK που αντιστοιχεί στο μητρώο (3.2.14) να είναι B-ευσταθής. Θεωρούμε τον αχχ συμμετρικό πίνακα  $M = (m_{ij})$  που ορίζεται από

$$(44) m_{ij} = b_i a_{ij} + b_j a_{ji} - b_i b_j, \quad 1 \leq i, j \leq q.$$

**ΠΡΟΤΑΣΗ 2.** (Butcher, Crouzeix). Υποθέτουμε ότι οι αριθμοί  $\tau_i$ ,  $1 \leq i \leq q$  είναι διάφοροι μεταξύ τους. Τότε η μέθοδος RK (3.2.12a,b) είναι B-ευσταθής αν και μόνο αν  $b_i \geq 0$  και ο πίνακας  $M$  ικανοποιεί  $z^* M z \geq 0 \quad \forall z \in \mathbb{C}^q$ ,  $z^* = \bar{z}^T$ . (Οι συνθήκες αυτές στα  $b_i$  και στον  $M$  λέγονται και συνθήκες "αλγεβρικής ευεστάθειας" της μεθόδου). Αν τα  $\tau_i$  δεν είναι όλα διαφορετικά, τότε οι συνθήκες  $b_i \geq 0$  και  $z^* M z \geq 0, \quad \forall z \in \mathbb{C}^q$ , είναι μόνο ικανές για B-ευστάθεια.

Απόδειξη: Θα αποδείξουμε μόνο το ικανό για οποιαδήποτε  $\tau_i$ ,  $1 \leq i \leq q$ . Για απλούστευση της απόδειξης (έτσι ώστε να μην ανησυχούμε για συμμετρία εσωτερικών γινομένων και για να μην παίρνουμε όλο πραγματικά μέρη) υποθέτουμε ότι το πρόβλημα (40) είναι πραγματικό, ότι ιαχύει η (41), ότι τα  $(\cdot, \cdot), \|\cdot\|$  είναι το ευκλείδειο εσωτ. γινόμενο, αντιστχ. νόρμα, στον  $\mathbb{R}^n$  και ότι  $x^T M x \geq 0 \quad \forall x \in \mathbb{R}^q$ . Θεωρούμε δύο λύσεις  $y^n, z^n, y^{n+1}, z^{n+1}$  των (3.2.12a,b). Αφαιρώντας έχουμε

$$(45) y^{n+1} - z^{n+1} = y^n - z^n + \sum_{i=1}^q a_{ij} \varphi^j, \quad 1 \leq i \leq q$$

$$(46) \quad y^{n+1} - z^{n+1} = y^n - z^n + \sum_{i=1}^q b_i \varphi^i,$$

όπου θέσαμε

$$(47) \quad \varphi^i = h(f(t^{n,i}, y^{n,i}) - f(t^{n,i}, z^{n,i})) \in \mathbb{R}^m, \quad 1 \leq i \leq q.$$

Η (46) δίνει

$$(48) \quad \|y^{n+1} - z^{n+1}\|^2 = \|y^n - z^n\|^2 + 2 \sum_{i=1}^q b_i (\varphi^i, y^n - z^n) + \|\sum_{i=1}^q b_i \varphi^i\|^2.$$

Χρησιμοποιώντας τώρα τις σχέσεις (45) έχουμε

$$\sum_{i=1}^q b_i (\varphi^i, y^n - z^n) = \sum_{i=1}^q b_i (\varphi^i, y^{n,i} - z^{n,i}) - \sum_{i,j=1}^q b_i a_{ij} (\varphi^i, \varphi^j).$$

Συνεπώς, η (48) γίνεται

$$(49) \quad \|y^{n+1} - z^{n+1}\|^2 = \|y^n - z^n\|^2 + 2 \sum_{i=1}^q b_i (\varphi^i, y^{n,i} - z^{n,i}) \\ + \sum_{i,j=1}^q b_i b_j (\varphi^i, \varphi^j) - 2 \sum_{i,j=1}^q b_i a_{ij} (\varphi^i, \varphi^j).$$

Τώρα

$$(\varphi^i, y^{n,i} - z^{n,i}) = h(f(t^{n,i}, y^{n,i}) - f(t^{n,i}, z^{n,i}), y^{n,i} - z^{n,i}) \leq 0$$

λόγω της (41). Επίσης λόγω της συμμετρίας του  $(\cdot, \cdot)$  έχουμε

$$2 \sum_{i,j=1}^q b_i a_{ij} (\varphi^i, \varphi^j) = \sum_{i,j=1}^q b_i a_{ij} (\varphi^i, \varphi^j) + \sum_{i,j=1}^q b_j a_{ji} (\varphi^i, \varphi^j).$$

Άρα η (49), λόγω της υπόθεσης  $b_i \geq 0$  δίνει

$$\begin{aligned}
 (50) \quad \|y^{n+1} - z^{n+1}\|^2 &\leq \|y^n - z^n\|^2 - \sum_{i,j=1}^q (b_i a_{ij} + b_j a_{ji} - b_i b_j) / (\varphi^i, \varphi^j) \\
 &= \|y^n - z^n\|^2 - \sum_{i,j=1}^q m_{ij} (\varphi^i, \varphi^j).
 \end{aligned}$$

Έστω  $\{\xi^k\}$ ,  $0 \leq k \leq m$  μία ορθοκανονική βάση του  $\mathbb{R}^m$  ως προς το  $(\cdot, \cdot)$ .

Τότε, αν  $\varphi^i = \sum_{k=1}^m c_k^{(i)} \xi^k$ ,  $\varphi^j = \sum_{k=1}^m c_k^{(j)} \xi^k$  έχουμε  $(\varphi^i, \varphi^j) = \sum_{k=1}^m c_k^{(i)} c_k^{(j)}$ .

Άρα

$$\begin{aligned}
 \sum_{i,j=1}^q m_{ij} (\varphi^i, \varphi^j) &= \sum_{i,j=1}^q m_{ij} \left( \sum_{k=1}^m c_k^{(i)} c_k^{(j)} \right) = \\
 &= \sum_{k=1}^m \left( \sum_{i,j=1}^q c_k^{(i)} m_{ij} c_k^{(j)} \right) \geq 0,
 \end{aligned}$$

επειδή  $x^T M x \geq 0 \quad \forall x \in \mathbb{R}^q$ . Συνεπώς η (50) δίνει την (43). @

Η μέθοδος RK

$$\begin{array}{c|c}
 \lambda & \lambda \\
 \hline
 1 & 
 \end{array}$$

με  $q=1$  είναι προφανώς B-ευσταθής αν  $\lambda \geq 1/2$ . Συνεπώς και η πεπλεγμένη μέθοδος του Euler ( $\lambda=1$ ) αλλά και η (πεπλεγμένη) μέθοδος του μέσου (3.2.15) ( $\lambda=1/2$ ,  $p=2$ ) είναι B-ευσταθείς. Οι διαγώνια πεπλεγμένες μέθοδοι (3.2.18) με  $q=2$  είναι B-ευσταθείς για  $\lambda \geq 1/4$ . (Συνεπώς η μέθοδος με  $\lambda=(1+3^{-1/2})/2$  είναι B-ευσταθής, και έχει τάξη  $p=3$ ). Η διαγώνια πεπλεγμένη μέθοδος με  $q=3$  (3.2.51) είναι B-ευσταθής ετην ενδιαφέρουσα περίπτωση που το  $\beta$  είναι η μεγαλύτερη ρίζα του πολυωνύμου  $\beta^3 - 3\beta^2/2 + \beta/2 - 1/24 = 0$ , δηλ. όταν  $\beta = 2 \cos(\pi/18) / \sqrt{3}$ , οπότε  $p=4$ . Όλες οι μέθοδοι Gauss-Legendre με  $q$  σημεία (τάξης  $p=2q$ ) είναι B-ευσταθείς! (Μάλιστα για αυτές τις μεθόδους  $M=0$ .) Συνεπώς όλες αυτές οι μέθοδοι είναι κατάλληλες για μη γραμμικά άκαμπτα ευστήματα όπως τα (40)-(41).

Απ' την άλλη μεριά η B-ευσταθής μέθοδος του τραπέζιου, που ως γνωστόν δίνεται από το μητρώο (βλ. 3.2.16))

$$\begin{array}{cc|c} 0 & 0 & 0 \\ 1/2 & 1/2 & 1 \\ \hline 1/2 & 1/2 & \end{array}$$

δεν είναι B-ευσταθής' πράγματι έχουμε

$$H = 1/4 \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$$

για το οποίο φυσικά  $\exists x \neq 0$  τέτοιο ώστε  $x^T H x < 0$ . Πολλά αριθμητικά πειράματα για κατάλληλα μη γραμμικά άκαμπτα προβλήματα έχουν όντως δείξει ότι η μέθοδος του τραπέζιου πάσχει από ευδεύρευση εφαλμάτων ετρογχύλευσης για μεγάλο n.

### Παρατηρήσεις

1. Για να βρούμε την περιοχή απόλυτης ευστάθειας μιάς πολυβηματικής μεθόδου στο μιγαδικό επίπεδο, δηλ. για να προσδιορίσουμε την περιοχή  $S \subset \{z: \operatorname{Re} z < 0\}$  για την οποία αν  $z = h\lambda \in S$  τότε οι ρίζες  $\zeta_j$  του πολυωνύμου

$$p(\zeta, z) \equiv p(\zeta) - z \varepsilon(\zeta)$$

ικανοποιούν  $|\zeta_j| < 1$ , μπορούμε να χρησιμοποιήσουμε τα κριτήρια του Schur ή των Routh-Hurwitz (Παρατήρηση 3.3.1). Στην πράξη χρησιμοποιείται και η εξής μέθοδος: Έστω  $\partial S$  το σύνορο της S. Επειδή οι ρίζες  $\zeta_j$  είναι συνεχείς συναρτήσεις της παραμέτρου  $z = h\lambda$ , και  $z \in S \Leftrightarrow |\zeta_j| < 1$ , τότε  $z \in \partial S$  όταν κάποια ρίζα  $\zeta_j$  βρεθεί στην

μοναδιαία περιφέρεια, δηλ. είναι της μορφής  $e^{i\theta}$  για κάποιο  $0 \leq \theta < 2\pi$ . Συνεπώς το εύρος  $\partial D$  δίνεται από την καμπύλη  $z(\theta) = \rho(e^{i\theta})/\varepsilon(e^{i\theta})$ ,  $0 \leq \theta < 2\pi$  την οποία μπορούμε να σχεδιάσουμε στο ημιπίεδο  $\{z: \operatorname{Re} z < 0\}$ . Το διάστημα απόλυτης ευστάθειας θα οριστεί από την τομή της  $z(\theta)$  με τον (πραγματικό) αρνητικό ημιάξονα:

2. Ξέραμε ότι το πολυώνυμο Taylor βαθμού  $\leq m$  μίας συνάρτησης  $f(z)$  αναλυτικής σε μία περιοχή του μηδενός είναι το (μοναδικό) πολυώνυμο βαθμού  $\leq m$  που ικανοποιεί την συνθήκη  $|f(z) - p(z)| = O(|z|^{\nu})$ , καθώς  $|z| \rightarrow 0$ , με  $\nu$  όσο το δυνατόν μεγαλύτερο για αυθαίρετη  $f$  (το  $\nu$  που προκύπτει είναι  $\nu = m+1$ ). Γενικεύοντας, αναζητούμε ρητή προσέγγιση της  $f$  της μορφής

$$r(z) = p(z)/q(z),$$

αναλυτική σε μία περιοχή του μηδενός, με αριθμητή  $p(z)$  βαθμού  $\leq m$  και παρονομαστή  $q(z)$ , βαθμού  $\leq n$ , που, για δεδομένα  $m, n$ , να ικανοποιεί

$$|f(z) - r(z)| = O(|z|^{\nu}),$$

με  $\nu$  όσο το δυνατόν μεγαλύτερο για γενική  $f(z)$ . Εύκολα βλέπουμε ότι υπάρχει μοναδική τέτοια συνάρτηση  $r(z)$ , η λεγόμενη  $(m, n)$  προσέγγιση Padé της  $f$  για την οποία το άριστο γενικό  $\nu$  είναι  $\nu = n+m+1$ . Οι

συντελεστές  $p_i, q_i$  των πολυωνύμων  $p$  και  $q$  ( $p(z) = \sum_{i=0}^m p_i z^i$ )

$q(z) = \sum_{i=0}^n q_i z^i$ ) μπορούν να βρεθούν αναπτύσσοντας κατά Taylor την συνάρτηση  $f(z)q(z) - p(z)$  γύρω απ' το μηδέν και ζητώντας να μηδενίζονται όσο περισσότεροι όροι της σειράς είναι δυνατόν.

Για την συνάρτηση  $f(z) = e^z$  είναι γνωστό (βλ. π.χ. [1.8]) ότι οι  $(m, n)$  προσεγγίσεις Padé δίνονται από τους τύπους

$$r_{m,n}(z) = p(z)/q(z) = \sum_{k=0}^m p_k z^k / \sum_{k=0}^n q_k z^k,$$

όπου

$$p_k = p_k(m, n) = (m+n-k)!m!/(m+n)!k!(m-k)!$$

και

$$q_k = q_k(m, n) = (-1)^k(m+n-k)!n!/(m+n)!k!(n-k)!$$

και διατάσσονται ευθέως ως στοιχεία ενός (άπειρου) πίνακα, του λεγόμενου πίνακα Padé για την  $e^z$ :

$m \backslash n$	0	1	2	...
0	1	$\frac{1}{1-z}$	$\frac{1}{1-z+z^2/2}$	...
1	$1+z$	$\frac{1+z/2}{1-z/2}$	$\frac{1+z/3}{1-2z/3+z^2/6}$	...
2	$1+z+z^2/2$	$\frac{1+\frac{2}{3}z+\frac{z^2}{6}}{1-z/3}$	$\frac{1+z/2+z^2/12}{1-z/2+z^2/12}$	...
⋮	⋮	⋮	⋮	⋮

Πολλές ρητές προσεγγίσεις του εκθετικού που προέρχονται από μεθόδους RK είναι στοιχεία του πίνακα Padé. Π.χ. οι μέθοδοι Gauss-Legendre  $(q, q)$  της διαγωνίου του πίνακα Padé. Η πεπλεγμένη μέθοδος του Euler δίνει  $r(z)=1/(1-z)$  που είναι στοιχείο της 1ης υπερδιαγωνίου του πίνακα Padé, τα στοιχεία της οποίας δίνουν ρητές προσεγγίσεις που αντιστοιχούν σε A-ευσταθείς μεθόδους. Γενικά είναι γνωστό (βλ. το άρθρο των Wanner, Hairer και Hørsett που αναφέραμε προηγουμένως) ότι για τα στοιχεία του πίνακα Padé

$$|r_{m,n}(z)| < 1 \text{ για } \operatorname{Re} z < 0 \Leftrightarrow n-2 \leq m \leq n,$$

δηλ. ότι  $A$ -ευσταθείς προεχθήσεις δίνουν μόνο τα στοιχεία της διαγωνίου και των δύο πρώτων υπερδιαγωνίων του πίνακα.

3. Μεταξύ των  $A$ -ευσταθών μεθόδων RK, δηλ. των μεθόδων που ικανοποιούν την σχέση

$$(51) \quad |r(z)| < 1 \text{ για κάθε } z \in \mathbb{C} \text{ με } \operatorname{Re} z < 0,$$

διακρίνουμε υπο-κατηγορίες ανάλογα με την συμπεριφορά του  $|r(z)|$  καθώς  $|z| \rightarrow \infty$ ,  $\operatorname{Re} z < 0$ , δηλ. της τιμής  $|r(\infty)|$ . Π.χ. για την μέθοδο του τραπεζίου ισχύει η (51) και η  $|r(\infty)| = 1$ . Αυτό μας κάνει να περιμένουμε ότι η μέθοδος του τραπεζίου δεν θα αποσβένει (για οποιοδήποτε  $h > 0$ ) αρκετά καλά συνιετώδες της λύσης ενός άκαμπτου ευστήματος που αντιστοιχούν σε πάρα πολύ μεγάλα  $|\operatorname{Re} \lambda_j|$  - και όπως αυτό φαίνεται στην πράξη -. Αντίθετα, η μέθοδος του Euler (πεπλεγμένη), με  $r(z) = 1/(1-z)$ , ικανοποιεί την (51) αλλά και την

$$(52) \quad |r(\infty)| = 0,$$

που πραγματικά μιμείται την ανάλογη ιδιότητα του εκθετικού  $|e^z| = e^{\operatorname{Re} z} \rightarrow 0$  όταν  $\operatorname{Re} z \rightarrow -\infty$ . Οι  $A$ -ευσταθείς μέθοδοι για τις οποίες ισχύει η (52) (λέγονται και " $L$ -ευσταθείς"), είναι ιδιαίτερα κατάλληλες λοιπόν για την αριθμητική λύση πολύ ακόμπτων ευστημάτων. Κάπου-ευδιάμεσα βρίσκονται  $A$ -ευσταθείς μέθοδοι για τις οποίες ισχύει

$$(53) \quad |r(\infty)| < 1,$$

όπως π.χ. η μέθοδος (3.12.18) με  $\lambda = (1 + 3^{-1/2})/2$  με  $q=2$ ,  $p=3$ . Τέτοιες μέθοδοι λέγονται και "ισχυρά  $A$ -ευσταθείς". (Ανάλογα, στον πραγματικό άξονα ορίζουμε  $L_0$  - και ισχυρά  $A_0$ -ευσταθείς μεθόδους). Για ισχυρά

$A$ -ευσταθείς μεθόδους μπορούμε να βελτιώσουμε την εκτίμηση ευστάθειας (30): π.χ. με  $e^n = 0$  ισχύει η  $\|y^n - z^n\| \leq e^{-|\operatorname{Re} \lambda_s| t} \|y^0 - z^0\|$ , όπου  $\lambda_s$  η

ιδιοτιμή του  $\Pi$  με το μεγαλύτερο πραγματικό μέρος  $\operatorname{Re} \lambda_s$ .

4. Προγράμματα γενικής χρήσης για την αριθμητική λύση ακάμπτων ευστημάτων χρησιμοποιούν ευρήτως μία οικογένεια μεθόδων (όπως π.χ. τις  $A(\beta)$ -ευσταθείς μεθόδους οπισθοδρομικών διαφορών) μεταβλητής τάξης καθώς και μεταβλητό βήμα που ελέγχεται με κάποια στρατηγική εκτίμησης του τοπικού εφάλματος. Προγράμματα τέτοια όπως η "μέθοδος του Gear" ή το πρόγραμμα EPISODE των Byrne et al. υπάρχουν σε πολλές βιβλιοθήκες αλγορίθμων. Δεν μπορούμε όμως να πούμε ακόμα ότι χράφηκε το πρόγραμμα που θα ολοκληρώνει εχθυμένα οποιοδήποτε άκαμπτο πρόβλημα, έστω και "μεσαίου" μεγέθους.

5. Για ανάλογη με την  $\beta$ -ευστάθεια θεωρία για πολυβηματικές μεθόδους βλ. π.χ. την εργασία του Dahlquist στα πρακτικά συνεδρίου Springer LNM v.506 (1976), καθώς και εργασίες των Butcher (SIAM JNA 18 (1981), 37-44), Nevanlinna-Liniger, π.χ. BIT 19 (1979), 53-72) κ.α. Η σημασία της ανάλυσης των μεθόδων για μη γραμμικές εξισώσεις είχε αναγνωριστεί από τον Dahlquist ήδη στην σημαντική εργασία του του 1963, όπου εισήγαγε την  $A$ -ευστάθεια.

#### Ασκήσεις 3.4

1. Βρείτε τα διαστήματα απόλυτης ευστάθειας για άμεσες μεθόδους RK τάξης  $p=2, 3$  και 4.
2. Για την λύση της απλής Δ.Ε.  $y'=f(t,y)$  θεωρείτε την πεπλεγμένη μέθοδο (που δεν ανήκει στην κατηγορία των RK ή των γραμμικών πολυβηματικών μεθόδων)

$$y^{n+1} = y^n + h(f^n + f^{n+1}) + \frac{h^2}{12} ((D_t f)^n + (D_t f)^{n+1}),$$

όπου  $D_t f = \partial_t f + (\partial_y f)f$ . Ορίστε την τάξη ακρίβειας της μεθόδου και υπολογίστε την. Βρείτε την περιοχή απόλυτης ευστάθειας της. (Μέθοδοι όπως η παραπάνω ανήκουν στην κατηγορία των "μεθόδων Obrechhoff" ή των διβηματικών "πολυπαραγωγικών" μεθόδων).



3. Βρείτε την περιοχή απόλυτης ευστάθειας της μεθόδου Adams-Moulton

$$y^{n+3} - y^{n+2} = h(9f^{n+3} + 19f^{n+2} - 5f^{n+1} + f^n)/24$$

(Δείξτε ότι η περιοχή είναι συμμετρική ως προς τον πραγματικό άξονα, ότι το διάστημα απόλυτης ευστάθειας είναι το  $(-3,0)$  και σχεδιάστε την περιοχή στο μιγαδικό επίπεδο. Τι συμβαίνει κοντά στον φανταστικό άξονα;)

4. Δείξτε για την μέθοδο οπισθοδρομικών διαφορών με  $k$  βήματα (3.3.6) ότι

(α) Για  $k=2$  είναι  $A$ -ευσταθής.

(β) Για  $k=3$  είναι  $A_0$ -ευσταθής αλλά όχι  $A$ -ευσταθής.

5(α). Θεωρείστε τις διαγώνια πλεγμένες μεθόδους (3.2.18) με παράμετρο  $\lambda \in \mathbb{R}$ . Για ποιές τιμές του  $\lambda$  είναι οι μέθοδοι  $A$ -ευσταθείς; Εφαρμόστε τις στο γραμμικό εύστημα (27). Τι παρατηρείτε σχετικά με τα γραμμικά ευστήματα που πρέπει να λυθούν σε κάθε στάδιο και σε κάθε βήμα;

(β) Θεωρείστε την μέθοδο του Calahan (3.2.50). Ποιά είναι η περιοχή απόλυτης ευστάθειας της;

6. Βρείτε την περιοχή της απόλυτης ευστάθειας των μεθόδων πρόβλεψης - διόρθωσης (i) και (ii) της Παρατήρησης 3.3.3 για  $m=1$  και  $m=2$  διορθώσεις.

7. Ως πρότυπο για ευστήματα για τα οποία οι ιδιοτιμές του πίνακα  $M$  είναι φανταστικές, (τέτοια ευστήματα προέρχονται π.χ. από υπερβολικές Μ.Δ.Ε.) θεωρούμε το "δυντηρητικό" πρόβλημα

$$(*) \quad \begin{cases} y' = i\lambda y, & t \geq 0, \lambda \in \mathbb{R}, \\ y(0) = 1, \end{cases}$$

του οποίου η λύση  $y(t) = e^{i\lambda t}$  ικανοποιεί  $|y(t)| = 1, t \geq 0$ . Λέμε ότι μία μέθοδος είναι 1-ευσταθής αν, όταν εφαρμοσθεί στο (\*), δίνει, για

κάθε  $h > 0$ , ακολουθία  $\{y^n\}$  τέτοια ώστε  $|y^n| \leq 1$ ,  $n \geq 0$ . Για μεθόδους RK αυτό σημαίνει ότι  $|\rho(iy)| \leq 1 \quad \forall y \in \mathbb{R}$ . Συνεπώς, κάθε A-ευσταθής μέθοδος είναι 1-ευσταθής. Από τις 1-ευσταθείς μεθόδους ξεχωρίζουμε τις λεγόμενες ευντηρητικές μεθόδους, για τις οποίες  $|y^n|=1$ ,  $n \geq 0$ , δηλ. για τις οποίες ισχύει ακριβώς το διακριτό ανάλογο της  $|y(t)|=1$ ,  $t \geq 0$ .

(α). Δείξτε ότι όλες οι μέθοδοι RK που έχουν ρητή συνάρτηση  $\rho(z)$  που δίνεται από οποιοδήποτε στοιχείο της διαγωνίου του πίνακα Padé για την  $e^z$  είναι 1-ευσταθείς και μάλιστα ευντηρητικές.

(β) Χρησιμοποιώντας το (α) και την Πρόταση 1 δείξτε ότι η μέθοδος Gauss-Legendre με 2 σημεία, (23), είναι A-ευσταθής.

(γ) Συγκρίνετε τις δύο μεθόδους Euler και την μέθοδο του τραπεζίου ως προς την καταλληλότητά τους για την αριθμητική ολοκλήρωση του (\*) για μεγάλο  $|z|$ .

8. (α) Να αποδειχθεί ότι οι διαγώνια πεπλεγμένες μέθοδοι (3.2.18) είναι B-ευσταθείς για  $\lambda \geq 1/4$ .

(β) Να αποδειχθεί ότι η μέθοδος Gauss-Legendre με  $q=2$  σημεία (3.2.19) είναι B-ευσταθής.

(γ) Να αποδειχθεί ότι η μέθοδος RK

1/8	1/8	1/4
3/8	3/8	3/4
1/2	1/2	

είναι A-ευσταθής αλλά όχι B-ευσταθής.

9. (Dahlquist) (α) Δείξτε ότι μία  $k$ -βηματική μέθοδος είναι A-ευσταθής αν και μόνο αν η συνάρτηση  $\rho(z)/\sigma(z)$  είναι αναλυτική και έχει μη αρνητικό πραγματικό μέρος για  $|z| > 1$ .

(β) Χρησιμοποιώντας το μέρος (α) δείξτε ότι μία άμεση  $k$ -βηματική μέθοδος δεν μπορεί να είναι A-ευσταθής.

4. ΠΑΡΕΜΒΟΛΗ ΚΑΙ ΠΡΟΣΕΓΓΙΣΗ

## 4.1 ΠΑΡΕΜΒΟΛΗ ΜΕ ΠΟΛΥΝΟΜΟ LAGRANGE

Στο κεφάλαιο αυτό θα κάνουμε μία σύντομη εισαγωγή στην προσέγγιση πραγματικών συνεχών συναρτήσεων μίας μεταβλητής με παρεμβολή με πολύνομα ή τμηματικά πολυωνυμικές συναρτήσεις. Ο λόγος που τα πολύνομα χρησιμοποιούνται κατ' εξοχήν για την προσέγγιση συναρτήσεων είναι βέβαια το γεγονός ότι μπορούν να υπολογισθούν (αλλά και να παραχωχισθούν ή να ολοκληρωθούν) εύκολα με ένα πεπερασμένο πλήθος προσθαφαιρέσεων και πολλαπλασιασμών αλλά και το γεγονός ότι έχουν καλές προσεγγιστικές ιδιότητες, όπως υποδηλώνει το θεώρημα του Weierstrass. Θα συμβολίζουμε με  $P_m$  τον διανυσματικό χώρο όλων των πραγματικών πολυωνύμων βαθμού  $\leq m$ .

Αρχίζουμε εξετάζοντας σύντομα στην παράγραφο αυτή διάφορα ερωτήματα που εχετίζονται με την παρεμβολή Lagrange. Το θέμα μας είναι γνωστό από το μάθημα της Εισαγωγής στην Αριθμητική Ανάλυση' για λεπτομέρειες ή αποδείξεις που θα παραλειφθούν βλ. π.χ. [5.2].

Έστω  $\tau \equiv \{x_i\}_{i=1}^n$  μία ακολουθία  $n \geq 1$  διακριτών σημείων. Το πολύνομο βαθμού  $n-1$

$$(1) L_i(x) = \prod_{j \neq i}^n (x-x_j)/(x_i-x_j),$$

λέγεται το i-στό πολύνομο Lagrange για την  $\tau$ . Προφανώς

$$L_i(x_j) = \delta_{ij} = \begin{cases} 1 & \text{αν } i=j \\ 0 & \text{αν } i \neq j \end{cases}$$

Συνεπώς για μία οποιαδήποτε συνάρτηση  $f(x)$  (ή γενικά για οποιαδήποτε δεδομένα  $f(x_i)$ ), το πολύνομο βαθμού  $\leq n-1$

$$(2) p(x) = \sum_{i=1}^n f(x_i) L_i(x)$$

## 4.1.2

ικανοποιεί  $p(x_i) = f(x_i)$ ,  $1 \leq i \leq n$ , δηλ. παρεμβάλλεται στις τιμές της συνάρτησης (των δεδομένων  $f(x_i)$ ) ετά σημεία της  $\tau$ . (λέγεται και "παρεμβάλλει την  $f$  ετά  $x_i$ "). Προφανώς υπάρχει μόνο ένα τέτοιο  $p \in P_{n-1}$  γιατί αν υπήρχε και άλλο  $q \in P_{n-1}$  με την ίδια ιδιότητα, η διαφορά τους  $r = p - q$  θα ήταν ένα στοιχείο του  $P_{n-1}$  με  $n$  διακριτές ρίζες, δηλ. θα ήταν το μηδενικό πολυώνυμο. Θα ονομάζουμε λοιπόν το πολυώνυμο  $p(x)$ , (2), πολυώνυμο παρεμβολής Lagrange για την συνάρτηση  $f$  ετά σημεία της  $\tau$ .

Η παράσταση (2) του πολυωνύμου παρεμβολής δεν είναι η πιο κατάλληλη για τις εφαρμοχές. Ίσως η πιο εύχρηστη είναι η παράστασή του στη μορφή Newton, που μπορεί να οριστεί με την βοήθεια διαιρεμένων διαφορών.

Η  $k$ -ετή διαιρεμένη διαφορά μίας συνάρτησης  $f$  ετά σημεία  $x_i, \dots, x_{i+k}$  είναι ο συντελεστής του μονωνύμου μεγίστου βαθμού (δηλ. του  $x^k$ ) του πολυωνύμου βαθμού  $\leq k$  που ευφωνεί με τις τιμές της  $f$  ετά σημεία  $x_i, \dots, x_{i+k}$ . Την ευφωλίζουμε με

$$f[x_i, \dots, x_{i+k}].$$

Έχουμε δηλ.

$$f[x_i] = f(x_i),$$

$$f[x_i, x_2] = (f(x_2) - f(x_1)) / (x_2 - x_1), \text{ για } x_1 \neq x_2$$

και γενικά (βλ. Θεκ. 1α)

$$(3) \quad f[x_i, \dots, x_{i+k}] = (f[x_{i+1}, \dots, x_{i+k}] - f[x_i, \dots, x_{i+k-1}]) / (x_{i+k} - x_i),$$

ή χρησιμοποιώντας την μορφή (2) του πολυωνύμου παρεμβολής,

$$(3') \quad f[x_i, \dots, x_{i+k}] = \sum_{j=1}^{i+k} f(x_j) / \left( \prod_{\substack{\mu=1 \\ \mu \neq j}}^{i+k} (x_j - x_\mu) \right).$$

Ο ορισμός της  $k$ -ετής διαιρεμένης διαφοράς μας δίνει το εξής αποτέλεσμα: Αν τα πολυώνυμα  $p_{i-1} \in \mathbb{P}_{i-1}$  συμφωνούν με την  $f$  στα σημεία  $x_1, \dots, x_i$  για  $i=k$  και  $k+1$ , αντίστοιχα, τότε

$$(4) \quad p_k(x) = p_{k-1}(x) + (x-x_1) \dots (x-x_k) f[x_1, \dots, x_{k+1}].$$

Πράγματι, το πολυώνυμο  $p_k - p_{k-1}$  είναι βαθμού  $\leq k$  και μηδενίζεται στα σημεία  $x_1, \dots, x_k$ . Εξ άλλου το μονώνυμο του μεγίστου βαθμού (δηλ. το αντίστοιχο μονώνυμο του  $p_k$ ) έχει εξ ορισμού συντελεστή  $f[x_1, \dots, x_{k+1}]$ . Άρα πρέπει να είναι της μορφής  $p_k(x) - p_{k-1}(x) = f[x_1, \dots, x_{k+1}] (x-x_1) \dots (x-x_k)$  δηλ. ισχύει η (4).

Η σχέση (4) μας δίνει την δυνατότητα να κατασκευάσουμε το πολυώνυμο παρεμβολής  $p(x) \in \mathbb{P}_{n-1}$  στα σημεία  $x_1, \dots, x_{n-1}$  βήμα προς βήμα, προσθέτοντας κάθε φορά ένα νέο σημείο παρεμβολής, πράγμα πολύ χρήσιμο στις εφαρμογές. Με τον παραπάνω συμβολισμό των  $p_i$ ,  $0 \leq i \leq n-1$ , έχουμε από την (4)

$$\begin{aligned} p(x) = p_{n-1}(x) &= p_0(x) + (p_1(x) - p_0(x)) + (p_2(x) - p_1(x)) + \\ &\quad + \dots + (p_{n-1}(x) - p_{n-2}(x)) = \\ (4') \quad &= f[x_1] + (x-x_1)f[x_1, x_2] + (x-x_1)(x-x_2)f[x_1, x_2, x_3] \\ &\quad + \dots + (x-x_1) \dots (x-x_{n-1})f[x_1, \dots, x_n], \end{aligned}$$

η οποία είναι η μορφή Newton του πολυωνύμου παρεμβολής, πολύ χρήσιμη στην θεωρία και στις εφαρμογές. (Ο υπολογισμός του  $p(x)$  από την (41) γίνεται π.χ. εύκολα με τον κανόνα του Horner).

Μία σημαντική ιδιότητα των διαιρεμένων διαφορών (Άσκηση 1β) είναι η εξής γενίκευση του θεωρήματος μέσης τιμής: Αν η συνάρτηση  $f$ , ορισμένη στο  $[a, b]$ , είναι  $k$  φορές παραγωγίσιμη στο  $(a, b)$  και τα  $x_1, \dots, x_{k+1}$  είναι διακριτά σημεία του  $[a, b]$ , τότε υπάρχει  $\xi \in (a, b)$ , τέτοιο ώστε

$$(5) f[x_1, \dots, x_{k+1}] = f^{(k)}(\xi)/k!$$

Η (5) δίνει τώρα την γνωστή μας έκφραση του εφάλματος του πολυωνύμου παρεμβολής. Έστω  $p = p_{n-1}$  το πολυώνυμο παρεμβολής Lagrange για την  $f$ , βαθμού  $\leq n-1$  στα διακριτά σημεία της

$\tau = \{x_i\}_{i=1}^n \subset [a, b]$  και έστω  $e_{n-1}(x) = f(x) - p_{n-1}(x)$  το εφάλμα της παρεμβολής. Έστω σημείο  $\bar{x} \in [a, b]$ , διάφορο των  $x_i$ ,  $1 \leq i \leq n$ . Αν  $p_n$  είναι το πολυώνυμο που παρεμβάλλεται στις τιμές της  $f$  στα  $n+1$  σημεία  $x_i$ ,  $1 \leq i \leq n$ , και  $\bar{x}$ , τότε  $p_n(\bar{x}) = f(\bar{x})$  και λόγω της (4)

$$p_n(x) = p_{n-1}(x) + (x-x_1) \dots (x-x_n) f[x_1, \dots, x_n, \bar{x}].$$

Άρα

$$f(\bar{x}) = p_n(\bar{x}) = p_{n-1}(\bar{x}) + f[x_1, \dots, x_n, \bar{x}] \prod_{i=1}^n (\bar{x} - x_i).$$

δηλ. για κάθε  $\bar{x} \notin \tau$

$$(6) e(\bar{x}) \equiv f(\bar{x}) - p_{n-1}(\bar{x}) = f[x_1, \dots, x_n, \bar{x}] \prod_{i=1}^n (\bar{x} - x_i).$$

Οι (5) και (6) δίνουν λοιπόν το

**ΘΕΩΡΗΜΑ 1.** Έστω ότι η  $f: [a, b] \rightarrow \mathbb{R}^1$  είναι  $n$  φορές παραγωγίσιμη στο  $(a, b)$ . Αν  $p_{n-1}(x)$  είναι το πολυώνυμο παρεμβολής Lagrange στα  $n$  διακριτά σημεία  $x_1, \dots, x_n \in [a, b]$ , τότε για κάθε  $\bar{x} \in [a, b]$  υπάρχει  $\xi = \xi(\bar{x}) \in (a, b)$  τέτοιο ώστε

$$(7) e_{n-1}(\bar{x}) \equiv f(\bar{x}) - p_{n-1}(\bar{x}) = \left( \prod_{i=1}^n (\bar{x} - x_i) \right) f^{(n)}(\xi)/n!$$

θα εξετάσουμε τώρα ορισμένες προβληματικές πλευρές της προέχουσας μιάς συνάρτησης από το πολυώνυμο παρεμβολής της. Είναι

## 4.1.5

γνωστό π.χ. ότι για ισαπέχοντα σημεία, δηλ. ακολουθίες  $\tau$  με  $x_i - x_{i-1} = \text{σταθ.}$ , μπορούμε, ακόμα και για ομαλότερες συναρτήσεις, να οδηγηθούμε σε πολυώνυμα παρεμβολής  $p_{n-1}(x)$  των οποίων το μέγιστο εφάλμα αυξάνει απεριόριστα καθώς αυξάνει το  $n$ . Το κλασικό παράδειγμα είναι το λεγόμενο παράδειγμα του Runge, το οποίο αναφέρεται στην συνάρτηση

$$(8) f(x) = 1 / (1 + 25x^2), \quad x \in [-1, 1],$$

( $f \in C^\infty(\mathbb{R}^1)$ !) και στην παρεμβολή της με πολυώνυμα Lagrange στις ακολουθίες  $\tau_n$ ,  $n \geq 2$ , ισαπέχοντων σημείων στο  $[-1, 1]$  όπου  $\tau_n = \{x_k^n\}$ ,  $k=1, 2, 3, \dots, n$ , με  $x_1^n = -1$  και  $x_k^n = x_{k-1}^n + h_n$ ,  $k=2, 3, \dots, n$ ,  $h_n = 2/(n-1)$ .

Κατ'αρχήν μία προκαταρκτική παρατήρηση. Έστω  $g(x) = 1/(ax+b)$ . Τότε για  $y_i$ ,  $1 \leq i \leq m$ , διακριτά σημεία έχουμε

$$(9) g[y_1, \dots, y_m] = (-a)^{m-1} \prod_{i=1}^m (ay_i + b)^{-1}.$$

Θα αποδείξουμε την (9) με επαγωγή. Για  $m=1$  ισχύει προφανώς. Παρατηρούμε επί τη ευκαιρία ότι η γενική διαιρεμένη διαφορά  $g[x_1, \dots, x_{i+k}]$  είναι συμμετρική συνάρτηση των  $x_1, \dots, x_{i+k}$ , δηλ. εξαρτάται μόνο από τους αριθμούς  $x_1, \dots, x_{i+k}$  και όχι από την σειρά με την οποία εμφανίζονται μέσα στις αγκύλες. Αυτό είναι προφανές από τον ορισμό της επειδή το πολυώνυμο παρεμβολής στα σημεία  $x_1, \dots, x_{i+k}$  εξαρτάται μόνο από τα σημεία και όχι από την σειρά τους. Έστω λοιπόν τώρα ότι η (9) ισχύει για  $m-1$ . Έχουμε, από την (3) και την συμμετρία της  $g$  ως προς τις μεταβλητές  $y_i$ ,

$$\begin{aligned} g[y_1, \dots, y_m] &= g[y_{m-1}, y_1, y_2, \dots, y_{m-2}, y_m] = \\ &= (g[y_1, \dots, y_{m-2}, y_m] - g[y_{m-1}, y_1, \dots, y_{m-2}]) / (y_m - y_{m-1}) \\ &= (g[y_1, \dots, y_{m-2}, y_m] - g[y_1, \dots, y_{m-1}]) / (y_m - y_{m-1}) \end{aligned}$$

## 4.1.6

$$\begin{aligned}
&= \{(-a)^{m-2} (\prod_{i=1}^{m-2} (ay_i + b)^{-1}) (ay_m + b)^{-1} - (-a)^{m-2} (\prod_{i=1}^{m-1} (ay_i + b)^{-1})\} / (y_m - y_{m-1}) \\
&= (-a)^{m-2} \prod_{i=1}^{m-2} (ay_i + b)^{-1} [1 / ((ay_m + b)(y_m - y_{m-1})) - 1 / ((ay_{m-1} + b)(y_m - y_{m-1}))] \\
&= (-a)^{m-1} \prod_{i=1}^m (ay_i + b)^{-1},
\end{aligned}$$

δηλ. ότι η (9) ισχύει και για  $m$ . Εφαρμόζοντας την (9) για τα διακριτά σημεία  $x_1^n, \dots, x_n^n, x$  και παραλείποντας για ευκολία τους άνω δείκτες ( $n$ ) στα σημεία  $x_k^n$  παίρνουμε για  $g(x) = (ax+b)^{-1}$

$$(9') \quad g[x_1, \dots, x_n, x] = ((-a)^n / (ax+b)) \prod_{i=1}^n (ax_i + b)^{-1}.$$

Γράφοντας τώρα την συνάρτηση  $f(x) = (1+25x^2)^{-1}$  ως

$$(10) \quad f(x) = (1/2) \{(5x+1)^{-1} - (5x-1)^{-1}\}$$

και παρατηρώντας (π.χ. από την (3')) ότι γενικά η διαιρεμένη διαφορά  $g[x_i, \dots, x_{i+k}]$  είναι γραμμική ως προς  $g$ , δηλ. ότι

$$(\lambda g_1 + \mu g_2)[x_i, \dots, x_{i+k}] = \lambda g_1[x_i, \dots, x_{i+k}] + \mu g_2[x_i, \dots, x_{i+k}],$$

έχουμε, από τις (9'), (10) ότι

$$\begin{aligned}
(11) \quad f[x_1, \dots, x_n, x] &= \\
&= (-5)^n (1/2) \{(5x+1)^{-1} \prod_{k=1}^n (5x_k + 1)^{-1} - (5x-1)^{-1} \prod_{k=1}^n (5x_k - 1)^{-1}\}.
\end{aligned}$$

Τώρα,  $x_{n-k+1} = -1 + (n-k)h_n = -1 + (n-1)h_n - (k-1)h_n = 1 - (k-1)h_n = -x_k$

Συνεπώς, για  $n$  άρτιο, έχουμε



$$\prod_{k=1}^n (5x_k \pm i)^{-1} = \prod_{k=1}^{n/2} (5x_k \pm i)^{-1} (-5x_k \pm i)^{-1} = - \prod_{k=1}^{n/2} (25x_k^2 + 1)^{-1}$$

Υποθέτοντας από δω και πέρα ότι ο  $n$  είναι άρτιος έχουμε από την (11) ότι

$$\begin{aligned} f[x_1, \dots, x_n, x] &= (-5)^n (1/2) \left\{ \prod_{k=1}^{n/2} (25x_k^2 + 1)^{-1} ((5x_k + i)^{-1} - (5x_k - i)^{-1}) \right\} \\ (12) \quad &= 5^n (-1)^{n/2} (25x^2 + 1)^{-1} \prod_{k=1}^{n/2} (25x_k^2 + 1)^{-1}. \end{aligned}$$

Από την (6) μπορούμε να γράψουμε το εφάλμα  $e_{n-1}(x) = f(x) - p_{n-1}(x)$  της παρεμβολής Lagrange στα σημεία  $x_1, \dots, x_n$  ως

$$e_{n-1}(x) = (x-x_1) \dots (x-x_n) f[x_1, \dots, x_n, x].$$

Συνοψώς η (12) δίνει

$$(13) \quad |(1+25x^2) e_{n-1}(x)| = \left| \prod_{i=1}^n (x-x_i) \right| \left( \prod_{k=1}^{n/2} 25/(25x_k^2+1) \right).$$

θέτουμε τώρα

$$(14) \quad r_n(x) \equiv 2(\log |(1+25x^2) e_{n-1}(x)|) / (n-1).$$

λόγω της (13), για  $x \neq x_k = x_k^n$

$$r_n(x) = 2(n-1)^{-1} \sum_{i=1}^n \log |x-x_i| - 2(n-1)^{-1} \sum_{k=1}^{n/2} \log (x_k^2 + (1/25)).$$

Άρα

$$\begin{aligned} (15) \quad \lim_{\substack{n \rightarrow \infty \\ n \text{ άρτιος}}} r_n(x) &= \int_{-1}^x \log(x-\xi) d\xi + \int_x^1 \log(\xi-x) d\xi \\ &\quad - \int_{-1}^1 \log(\xi^2 + (1/25)) d\xi \equiv r(x). \end{aligned}$$

## 4.1.8

Η  $r(x)$  είναι ευεχής στο  $[-1,1]$  και  $r(1) > 0$  (Άσκηση 2). Συνεπώς υπάρχει ανοιχτό διάστημα στο  $[-1,1]$  στο οποίο η  $r(t)$  είναι θετική. Άρα υπάρχει  $x_0 \neq x_k^n$  στο  $[-1,1]$  για κάθε  $n$  και  $k \leq n$  για το οποίο  $r(x_0) > 0$ . Άρα από την (14)

$$\lim_{\substack{n \rightarrow \infty \\ n \text{ άρτιος}}} (\log |(1+25x_0^2) e_{n-1}(x_0)|) = +\infty,$$

δηλ.

$$\lim_{\substack{n \rightarrow \infty \\ n \text{ άρτιος}}} \log |e_{n-1}(x_0)| = +\infty,$$

δηλ.

$$\lim_{\substack{n \rightarrow \infty \\ n \text{ άρτιος}}} |e_{n-1}(x_0)| = +\infty,$$

που δείχνει ότι υπάρχει ακολουθία ομοιομόρφων διαμερισμών του  $[-1,1]$  για την οποία η αντίστοιχη ακολουθία των μεγίστων εφαλμάτων του πολυώμου παρεμβολής τείνει στο άπειρο.

Ας κοιτάξουμε το πρόβλημα λίγο γενικότερα. Έστω  $P_n f \in P_{n-1}$  το πολυώνυμο παρεμβολής μιας ευεχούς συνάρτησης  $f$  στα (διακριτά) σημεία  $x_i$ ,  $1 \leq i \leq n$ , της πεπερασμένης ακολουθίας  $\tau$ , η οποία υποθέτουμε ότι περιέχεται ε' ένα διάστημα  $[a,b]$ . Για  $f \in C[a,b]$  θεωρούμε την "maximum" νόρμα της ομοιομορφής εύκλεισης

$$(16) \|f\|_\infty = \max_{a \leq x \leq b} |f(x)|.$$

(θεωρούμε γνωστά από την Συναρτησιακή Ανάλυση τα βασικά περί χώρων συναρτήσεων με νόρμα). Από την (2)

$$|(P_n f)(x)| \leq \sum_{i=1}^n |f(x_i)| |L_i(x)| \leq \max_i |f(x_i)| \sum_{i=1}^n |L_i(x)|.$$

Εισάγοντας την λεγόμενη συνάρτηση του Lebesgue

$$(17) \quad \lambda_n(x) = \sum_{i=1}^n |L_i(x)|,$$

έχουμε τελικά ότι

$$(18) \quad \|P_n f\|_\infty \leq \|\lambda_n\|_\infty \|f\|_\infty.$$

Δεν είναι δύσκολο να δούμε (βλκ 3) ότι υπάρχει, για κάθε  $n$  με  $n$  σημεία, συνεχής συνάρτηση  $f \neq 0$  τέτοια ώστε να ισχύει η (18) ως ισότις. (Δηλ. η νόρμα του φραγμένου γραμμικού τελεστή  $P_n: X \rightarrow X$ ,  $X=(C[a,b], \|\cdot\|_\infty)$  είναι ίση με  $\|\lambda_n\|_\infty$ . Είναι γνωστό εξ άλλου ότι υπάρχουν σταθερές  $c_1 > 0$ ,  $c_2 \in \mathbb{R}$ , ανεξάρτητες των  $n, \tau$  τέτοιες ώστε για την λεγόμενη "σταθερά Lebesgue"  $\|\lambda_n\|_\infty$  να ισχύει

$$(19) \quad \|\lambda_n\|_\infty \geq c_1 \log n + c_2.$$

(βλ. π.χ. το βιβλίο του Rivlin [4.9, σελ. 90-91] για μία απόδειξη με  $c_1=4/n^2$ ,  $c_2=-1$ . Είναι γνωστό (Erdős) ότι η καλύτερη σταθερά  $c_1$  είναι  $c_1=2/n$ ). Το αποτέλεσμα αυτό μας οδηγεί στο περίφημο θεώρημα του Faber, βάσει του οποίου για κάθε δεδομένο "τριγωνικό σύστημα παρεμβολής", δηλ. για κάθε ακολουθία  $\tau_1, \tau_2, \dots, \tau_n, \dots, \tau_n = \{x_n^1, \dots, x_n^n\}$ , πεπερασμένων ακολουθιών σημείων παρεμβολής  $x_n^i \in [a,b]$ , υπάρχει συνάρτηση  $f^* \in C[a,b]$  τέτοια ώστε

$$\overline{\lim}_{n \rightarrow \infty} \|P_n f^* - f^*\|_\infty = \infty,$$

βλ. π.χ. [4.9, σελ. 92-3].

Για ισαπέχοντα σημεία μάλιστα είναι γνωστό ότι η σταθερά Lebesgue αυξάνεται εκθετικά υπάρχουν σταθερές  $c_1, c_2 > 0$ ,  $\alpha_1 > \alpha_2 > 1$ , ανεξάρτητες των  $n, h_n = (b-a)/(n-1)$ , τέτοιες ώστε για κάθε  $n > 1$

$$c_2(\alpha_2)^n \leq \|\lambda_n\|_\infty \leq c_1(\alpha_1)^n.$$

βλ. π.χ. [4.9, εελ. 99] για μία απόδειξη με  $n$  άρτιο,  $\alpha_2=(1.5)^{1/2}$ ,  $\alpha_1=2^{1/2}e$ . Συνεπώς το παράδειγμα του Runge δεν πρέπει να μας εκπλήσσει και πολύ.

Απ' την άλλη μεριά είναι ευδαρμυτικό (βλ. Θεω. 4) ότι για κάθε  $f \in C[a,b]$  υπάρχει τριγωνικό σύστημα παρεμβολής τέτοιο ώστε η ακολουθία πολυνομών  $P_n f$  να ευκλίνει ομοιόμορφα στην  $f$  στο  $[a,b]$  καθώς  $n \rightarrow \infty$ . Φυσικά, το τριγωνικό αυτό σύστημα (ακολουθία  $\tau_n$   $n$  διακριτών σημείων) μπορεί να μην είναι εύκολο να κατασκευασθεί και βέβαια θα διαφέρει από συνάρτηση σε συνάρτηση. Τίθεται το ερώτημα: υπάρχει σύστημα παρεμβολής, που να μπορεί να βρεθεί εύκολα, τέτοιο ώστε η  $\|\lambda_n\|_\infty$  να αυξάνεται αργά (λογαριθμικά είναι φυσικά ο αργότερος ασυμπτωτικά δυνατός τρόπος, βλ. (19)) καθώς  $n \rightarrow \infty$ ; Η απάντηση είναι θετική και η λύση δίνεται από τις ρίζες των (πανταχού παρόντων) πολυνομών Chebyshev! Μπορεί ν' αποδειχθεί, βλ. π.χ. [4.9, εελ. 93 et seq.], ότι αν τα σημεία  $x_1, \dots, x_n$  της  $\tau$  είναι οι ρίζες του πολυνομού Chebyshev  $n^{\text{ου}}$  βαθμού για το  $[a,b]$ , δηλ. αν

$$(20) \quad x_i = x_i^c \equiv \{(a+b) - (a-b) \cos((2i-1)\pi/2n)\}/2, \quad 1 \leq i \leq n$$

(απεικονίστε το διάστημα  $[-1,1]$  πάνω στο  $[a,b]$  και χρησιμοποιήστε ότι  $T_n(z) = \cos(n \cos^{-1} z)$ ,  $-1 \leq z \leq 1$ , βλ. Παρ. 1.6), τότε για την αντίστοιχη σταθερά του Lebesgue  $\|\lambda_n\|_\infty \equiv \|\lambda_n^c\|_\infty$  έχουμε

$$(21) \quad \|\lambda_n^c\|_\infty \leq (2/n) \log n + 4,$$

που, σε συνδυασμό με την (19), δείχνει ότι οι ρίζες των πολυνομών Chebyshev δίνουν ένα πολύ καλό τριγωνικό σύστημα παρεμβολής. Το πείραμα δείχνει (βλ. παρακάτω γιατί) ότι η συνάρτηση Runge  $f(x) = 1/(1+25x^2)$  στο  $[-1,1]$  δίνει ακολουθία πολυνομών παρεμβολής στα  $x_i^c$  που ευκλίνει για  $n \rightarrow \infty$  ομοιόμορφα στην  $f$ .

Η (21) μας δίνει επίσης ότι το πολυώνυμο παρεμβολής  $P_n f$  στα σημεία Chebyshev  $x_i^c$ ,  $1 \leq i \leq n$  είναι "εχεδόν" (modulo λογαριθμικό παράγοντα) η καλύτερη προσέγγιση της  $f$  στον  $(P_{n-1}, \|\cdot\|_\infty)$ . Πράγματι, έστω  $p^* \in P_{n-1}$  η βέλτιστη προσέγγιση της  $f \in C[a, b]$  από στοιχεία του  $P_{n-1}$  ως προς την νόρμα  $\|\cdot\|_\infty$ , δηλ. έστω ότι

$$(22) \|f - p^*\|_\infty = \min_{p \in P_{n-1}} \|f - p\|_\infty$$

(Γιά την ύπαρξη, μοναδικότητα, ιδιότητες και κατασκευή του  $p^*$  βλ. π.χ. [5.1]). Επειδή ο τελεστής  $P_n$  της παρεμβολής στα σημεία  $x_i$ ,  $1 \leq i \leq n$  μιάς οποιασδήποτε  $\tau$  είναι ταυτότητα στον  $P_{n-1}$ , έχουμε

$$f - P_n f = f - p - P_n(f - p) \quad \forall p \in P_{n-1}$$

οπότε από την (18) έχουμε  $\forall p \in P_{n-1}$

$$\|f - P_n f\|_\infty \leq \|f - p\|_\infty + \|P_n(f - p)\|_\infty \leq \|f - p\|_\infty + \|\lambda_n\|_\infty \|f - p\|_\infty$$

Συμπεραίνουμε, παίρνοντας  $p = p^*$  στην παραπάνω ανισότητα ότι

$$(23) \min_{p \in P_{n-1}} \|f - p\|_\infty \leq \|f - P_n f\|_\infty \leq (1 + \|\lambda_n\|_\infty) \min_{p \in P_{n-1}} \|f - p\|_\infty$$

Γιά παρεμβολή στις ρίζες των πολυωνύμων Chebyshev,  $\|\lambda_n\|_\infty = \|\lambda_n^c\|_\infty$ , που αυξάνεται πολύ αργά (βλ. (21)) με το  $n$ . Π.χ. για  $n \leq 20$  είναι γνωστό ότι  $1 + \|\lambda_n^c\|_\infty \leq 4$ . Αυτό σημαίνει, λόγω της (23), ότι η καλύτερη δυνατή προσέγγιση της  $f$  στον  $(P_{n-1}, \|\cdot\|_\infty)$  (για  $n \leq 20$ ) δίνει εφάλλα το οποίο

είναι το πολύ 4 φορές μεγαλύτερο από το εφάλμα του πολυωνύμου παρεμβολής στά επρεία Chebyshev  $x_i^c$ . Φυσικά το πολυώνυμο παρεμβολής είναι πολύ απλό να κατασκευασθεί σε αντίθεση με το  $p^*$ !

Η (23) μας δίνει ένα πολύ χρήσιμο άνω φράγμα του εφάλματος  $\|f - P_n^c\|_\infty$  της παρεμβολής ευαρτήσεως της σταθεράς Lebesgue (που

εξαρτάται μόνο από την  $n$ ) και του εφάλματος  $\min_{p \in P_{n-1}} \|f - p\|_\infty = \|f - p^*\|_\infty$

της βέλτιστης προσέγγισης  $p^*$  της  $f$  στον  $(P_{n-1}, \|\cdot\|_\infty)$ , (που εξαρτάται μόνο από την  $f$ , το διάστημα  $[a, b]$  και το  $n$ ). Για το τελευταίο έχουμε πολύ ακριβείς εκτιμήσεις από την θεωρία προσέγγισης για διάφορες κλάσεις ευαρτήσεων  $f$  (θεωρήματα Jackson). Συνήθως υποθέτουμε ότι για κάποιον ακέραιο  $r \geq 0$ ,  $f \in C^r[a, b]$ , δηλ. ότι η  $f$  έχει  $r$  συνεχείς παραγώγους στο  $[a, b]$  ( $C^0[a, b] = C[a, b]$ ), εκφράζουμε δε τα φράγματα των εφαλμάτων ευαρτήσεως του μέτρου συνέχειας της  $f^{(r)}$ .

Για μία ευάρτηση  $g \in C[a, b]$ , το μέτρο συνέχειάς της είναι μία ευάρτηση  $w(g; h)$ , που ορίζεται για  $h \geq 0$  από

$$(24) \quad w(g; h) = \max\{|g(x) - g(y)| : x, y \in [a, b], |x - y| \leq h\}.$$

Είναι προφανές ότι  $0 \leq h_1 \leq h_2 \Rightarrow w(g; h_1) \leq w(g; h_2)$  και ότι για  $h_1, h_2 \geq 0$   $w(g; h_1 + h_2) \leq w(g; h_1) + w(g; h_2)$ . Είναι επίσης προφανές ότι επειδή  $g \in C[a, b]$ ,

$$(25) \quad \lim_{h \downarrow 0} w(g; h) = 0.$$

Ο ρυθμός όμως (εάν ευάρτηση του  $h$ ) με τον οποίο η  $w(g; h)$  τείνει στο μηδέν καθώς  $h \rightarrow 0$ , μεταβάλλεται όταν η  $g$  διατρέχει τις συνεχείς ευαρτήσεις στο  $[a, b]$ . Δεν είναι δύσκολο να δούμε (Άσκηση 5) ότι ο χρησιότερος τρόπος με τον οποίο μπορεί το μέτρο συνέχειας  $w(g; h)$  μιάς  $g \in C[a, b]$  να τείνει στο μηδέν (αν  $g \neq \text{const.}$ ) είναι γραμμικός ως προς  $h$ , δηλ. όταν υπάρχει σταθερά  $c$ , ανεξάρτητη του  $h$ , τέτοια ώστε

$w(g;h) \leq ch$ . Αυτό επιτυγχάνεται από τις συναρτήσεις  $g \in C^1[a,b]$  για τις οποίες  $w(g;h) \leq \|g'\|_{\infty} h$  και, γενικότερα, από τις συναρτήσεις που ικανοποιούν συνθήκη Lipschitz με σταθερά  $L$  στο  $[a,b]$  για τις οποίες  $w(g;h) \leq Lh$ . Μια μεγαλύτερη κατηγορία συνεχών συναρτήσεων είναι οι συναρτήσεις που ικανοποιούν μια συνθήκη Hölder με εκθέτη  $\alpha \in (0,1)$  στο  $[a,b]$  δηλ. οι συναρτήσεις  $g \in C[a,b]$  για τις οποίες υπάρχει σταθερά  $K$  και  $\alpha \in (0,1)$  τέτοιες ώστε

$$(26) \quad w(g;h) \leq Kh^{\alpha} \text{ για } h \geq 0.$$

Π.χ. η συνάρτηση  $g(x) = x^{\alpha}$ ,  $0 < \alpha < 1$  στο διάστημα  $[0,1]$  έχει  $w(g;h) = h^{\alpha}$ . Για την συνάρτηση  $g(x) = |x|^{1/2}$  στο  $[-1,1]$  έχουμε για  $|x-y| \leq h$ ,  $x, y \in [-1,1]$  ότι  $\max |g(x) - g(y)| = w(g;h) = |g(0) - g(h)| = \sqrt{h}$ .

Ένα από τα θεωρήματα του Jackson (βλ. π.χ. [4.9, βελ. 23]) μας λέει ότι αν  $f \in C^r[a,b]$  και  $n > r+1$ , τότε

$$(26') \quad \min_{p \in P_{n-1}} \|f-p\|_{\infty} \leq c(r) \left( (b-a)/(n-1) \right)^r w(f^{(r)}; (b-a)/2(n-1-r)),$$

όπου  $c(r) = 6(3e)^r / (r+1)$ . Π.χ. για την συνάρτηση  $f(x) = \sqrt{x}$  στο  $[0,1]$  για την οποία  $r=0$ ,  $w(f;h) = h^{1/2}$ , έχουμε ότι

$$\min_{p \in P_{n-1}} \|\sqrt{x-p}\|_{\infty} = \|\sqrt{x-p^*}\|_{\infty} \leq 6((b-a)/2(n-1))^{1/2},$$

δηλ. ότι το εφάλμα της βέλτιστης προσέγγισης τείνει στο 0 όπως το  $n^{-1/2}$  όταν  $n \rightarrow \infty$ . Από την (23) βλέπουμε λοιπόν ότι το εφάλμα της παρεμβολής Lagrange για την  $f(x) = \sqrt{x}$  στο  $[0,1]$  επί σημεία Chebyshev  $x_i^C$ ,  $1 \leq i \leq n$ , τείνει στο μηδέν τουλάχιστον όσο γρήγορα τείνει η ακολουθία  $\log n \cdot n^{-1/2}$  όταν  $n \rightarrow \infty$ . Για μία οποιαδήποτε συνάρτηση του  $C^k[a,b]$  (σταθερό  $k \geq 1$ ) - π.χ. για την συνάρτηση του Runge στο  $[-1,1]$  - η παρεμβολή στα σημεία Chebyshev δίνει εφάλμα που τείνει στο μηδέν

- βλ. (26) με  $n=k-1$ ,  $\omega(f^{(k-1)}; h) \leq \|f^{(k)}\|_{\infty} h$  - τουλάχιστον όσο γρήγορα και η ακολουθία  $(\log n)/n^k$  όταν  $n \rightarrow \infty$ .

Από την (26') βλέπουμε ότι το εφάλμα  $\|f-p^*\|_{\infty}$  είναι μικρό αν ο λόγος  $(b-a)/(n-1)$  γίνει μικρός. Αυτό μπορεί να επιτευχθεί είτε αυξάνοντας τον βαθμό του πολυωνύμου  $n-1$  είτε υποδιαιρώντας το διάστημα  $[a,b]$  σε μικρότερα διαστήματα σε κάθε ένα απ' τα οποία προσεγγίζουμε την  $f$  με κατάλληλα πολυώνυμα σταθερού βαθμού, ευνοϊκά δηλ. προσεγγίζοντας την  $f$  με μία τμηματικά πολυωνυμική συνάρτηση. Έστω ότι υποδιαιρούμε το διάστημα  $[a,b]$  σε  $k$  υποδιαστήματα ( $k$  μεγάλο) στο κάθε ένα από τα οποία προσεγγίζουμε την συνάρτησή μας με ένα πολυώνυμο βαθμού  $n-1$  ( $2 \leq n \leq 6$  στην πράξη). Αυτό ισοδυναμεί, όσο αφορά το μέγεθος του φράγματος του εφάλματος, περίπου με την χρήση ενός πολυωνύμου βαθμού  $k(n-1)$  στο  $[a,b]$ . Ο συνολικός αριθμός των σταθερών ("βαθμών ελευθερίας"), δηλ. των συντελεστών των πολυωνύμων που πρέπει να υπολογίσουμε, παραμένει επίσης περίπου ο ίδιος ( $\approx kn$ ) και για τις δύο διαδικασίες. Υπάρχουν όμως στην πράξη σοβαρές διαφορές ανάμεσά τους:

(α) Ο υπολογισμός ενός πολυωνύμου μεγάλου βαθμού χρειάζεται (π.χ. αν πούμε για βαθμό  $\geq 20$ ) πολλή προσοχή στις πράξεις (π.χ. έκφραση του πολυωνύμου συνάρτησει βάσης καταλλήλων ορθογωνίων πολυωνύμων κλπ.), γιατί πολύ εύκολα μπορούμε, λόγω εφάλματος ετρογχύλευσης, να καταλήξουμε σε overflow (αετάθεια στους υπολογισμούς) ή underflow.

(β) Στην πράξη, σε πολλούς υπολογισμούς, εκφράζουμε τις προσεγγίσεις μας σαν γραμμικούς συνδυασμούς καταλλήλων συναρτήσεων βάσης. Για να υπολογίσουμε την τιμή ενός πολυωνύμου βαθμού  $k(n-1)$  στο  $[a,b]$  πρέπει να υπολογίσουμε  $k(n-1)+1$  συντελεστές και τις τιμές  $k(n-1)+1$  συναρτήσεων βάσης που είναι πολυώνυμα, δηλ. έχουν φορές όλο το  $[a,b]$ . Αντίθετα, για τον υπολογισμό ε' ένα σημείο μιας τμηματικά πολυωνυμικής συνάρτησης, υπολογίζουμε το πολύ ένα (μικρό) γραμμικό συνδυασμό (της τάξεως του  $n$ ) συναρτήσεων βάσης που είναι τμηματικά πολυωνυμικές συναρτήσεις βαθμού  $\leq n-1$  με μικρό φάρα (της τάξεως  $n$  υποδιαστημάτων) ανεξάρτητα του αριθμού των υποδιαστημάτων  $k$  που μπορεί να γίνει όσο μεγάλος θέλουμε. Οι συντελεστές, για πολλές εφαρμογές, είναι λύσεις γραμμικών συστημάτων (μεγέθους  $O(kn) \times O(kn)$ )



τα οποία, για προσέγγιση με τμηματικά πολυωνυμικές συναρτήσεις, έχουν πίνακες μεγάλους αλλά αραιούς (ευνήθως πίνακες ζώνης με πλάτη ζώνης  $0(n)$  ενώ για προσέγγιση με πολυώνυμα στο  $[a,b]$  έχουν τυπικά μεγάλους αλλά πυκνούς πίνακες.

(γ) Πολλές φορές η προσέγγιση γίνεται μέσω παρεμβολής. (Ακριβέστερα, το εφάλμα μίας προσέγγισης πολλές φορές είναι δυνατόν να εκτιμηθεί εύκολα συναρτήσει του εφάλματος της παρεμβολής). Είδαμε τα προβλήματα της πολυωνυμικής παρεμβολής για ομοιόμορφο διαμερισμό με βήμα  $h=(b-a)/k$ ,  $k \gg 1$ . Αίτιετα, η παρεμβολή ε' ένα διάστημα μικρού πλάτους  $h=b-a$  με πολυώνυμα μικρού (εταθερού) βαθμού  $n$ , δηλ. η παρεμβολή Lagrange τοπικά, είναι πολύ επιτυχής: η (7) μας δίνει για το εφάλμα  $e_{n-1}(x)$  για σταθερό  $n$  και  $f \in C^n[a,b]$  ότι

$$(27) |e_{n-1}(x)| \leq \prod_{i=1}^n |x-x_i| |f^{(n)}(\xi)/n!| \leq C_n h^n \max_{a \leq x \leq b} |f^{(n)}(x)|$$

#### Παρατηρήσεις

1. Είναι δυνατόν (βλ. π.χ. [4.7]) να βελτιώσουμε λίγο περισσότερο την σταθερά του Lebesgue που δίνουν τα σημεία Chebyshev (20) αν θεωρήσουμε για κάθε  $n$  την λεγόμενη επέκταση των σημείων Chebyshev. Τα αντίστοιχα σημεία  $x_i$ ,  $1 \leq i \leq n$  της  $\tau$  δίδονται τότε από τους τύπους

$$x_i = x_i^\tau = [(a+b) - (a-b) \cos((2i-1)\pi/2n) / \cos(\pi/2n)] / 2, \quad 1 \leq i \leq n.$$

Για την αντίστοιχη σταθερά του Lebesgue ισχύει ότι

$$(2/n) \log n + 0.5 \leq \|L_n^\tau\|_\infty \leq (2/n) \log n + .73.$$

Μπορεί επίσης να αποδειχθεί ότι για κάθε  $n$  η τιμή της  $\|L_n^\tau\|_\infty$  απέχει το πολύ .201 από την μικρότερη δυνατή τιμή  $\|L_n\|_\infty$  που μπορεί να επιτευχθεί για κάθε  $n$  με κατάλληλη επιλογή της  $\tau$ .

2. Από τον τύπο του εφάλματος της παρεμβολής (7) έχουμε για  $f \in C^n[a, b]$  ότι

$$(28) \|f - P_n f\|_\infty \leq \|w\|_\infty \|f^{(n)}\|_\infty / n!$$

όπου  $w(x) = \prod_{i=1}^n (x - x_i)$ ,  $x \in [a, b]$ .

Είναι γνωστό (βλ. Θεκ. 6) ότι

$$(29) \min_{x_i \in [a, b]} \|w\|_\infty = 2(b-a)^n / 4^n$$

και ότι η ελάχιστη αυτή τιμή της  $\|w\|_\infty$  συμβαίνει όταν τα σημεία παρεμβολής  $x_i$  συμπίπτουν με τις ρίζες  $x_i^c$ ,  $1 \leq i \leq n$  (20), του πολυωνύμου Chebyshev  $n^{\text{ου}}$  βαθμού! Ξαναβλέπουμε δηλ. την σημασία των  $x_i^c$  ως κόμβων στην παρεμβολή Lagrange. Συνεπώς η αποτυχία της παρεμβολής Lagrange πρέπει να οφείλεται στον όρο  $\|f^{(n)}\|_\infty / n!$  ο οποίος για πολλές  $C^\infty$  συναρτήσεις δεν είναι φραγμένος καθώς  $n \rightarrow \infty$ . Υπάρχουν όμως βέβαια κλάσεις συναρτήσεων (π.χ. εκείνες για τις οποίες  $\|f^{(n)}\|_\infty \leq n^n$ ,  $n=1, 2, \dots$  για κάποιο  $0 \leq k < \infty$ ) για τις οποίες η ακολουθία των πολυωνύμων παρεμβολής συγκλίνει ομοιόμορφα στην  $f$  καθώς  $n \rightarrow \infty$  για οποιαδήποτε τριγωνικό εύστημα παρεμβολής, (θεκ. 7).

#### Θεκίσεις 4.1

1. (α) Αποδείξτε την (3).
- (β) Αποδείξτε την (5).

2. Αποδείξτε για το παράδειγμα του Runge ότι το όριο  $\lim_{n \rightarrow \infty} r_n(x)$ ,  $n$  άρτιος  $\rightarrow \infty$ , υπάρχει και ορίζει μία συνεχή συνάρτηση  $r(x)$  στο  $[-1, 1]$  με  $r(1) > 0$ .

3. Δείξτε ότι για κάθε  $\tau$  με  $n$  σημεία υπάρχει  $f \in C[a, b]$ ,  $f \neq 0$  τέτοια ώστε να ισχύει η (18) ως ισότητα. (Υπόδειξη: έστω  $\hat{x} \in [a, b]$  τέτοιο ώστε  $\lambda_n(\hat{x}) = \|\lambda_n\|_\infty$ . Ορίστε  $\epsilon_i = \text{sign}(L_i(\hat{x}))$ ,  $1 \leq i \leq n$  και κατασκευάστε συνεχή συνάρτηση  $g$  στο  $[a, b]$  τέτοια ώστε  $g(x_i) = \epsilon_i$ ,  $1 \leq i \leq n$  και  $\|g\|_\infty = 1$  - π.χ. την τμηματικά γραμμική συνεχή συνάρτηση με  $g(x_i) = \epsilon_i$  με επέκταση με σταθερές έξω απ' το διάστημα  $[\min_i x_i, \max_i x_i]$ .  
- Δείξτε τότε ότι  $\|P_n g\|_\infty = \|\lambda_n\|_\infty \|g\|_\infty$ .

4. (α) Δείξτε ότι για κάθε  $f \in C[a, b]$  υπάρχει τριγωνικό εύθετο παρεμβολής τέτοιο ώστε  $\|f - P_n\|_\infty \rightarrow 0$ ,  $n \rightarrow \infty$ . (Χρησιμοποιείστε το γεγονός ότι το πολυώνυμο  $p^* \in P_{n-1}$  της βέλτιστης προσέγγισης της  $f$  από στοιχεία του  $P_{n-1}$  έχει την εξής ιδιότητα: Υπάρχουν  $n+1$  διακριτά σημεία  $x_i$  στο  $[a, b]$  ετά οποία το εφάλμα  $f - p^*$  έχει ίσες απολύτως και εναλλασσόμενες στο πρόσημο τιμές. Μάλιστα  $|(f - p^*)(x_i)| = \|f - p^*\|_\infty$ . Συνεπώς η διαφορά  $f - p^*$  θα μηδενίζεται σε  $n$  διακριτά σημεία  $y_i \in [a, b]$ , δηλ. το  $p^*$  θα είναι το πολυώνυμο παρεμβολής Lagrange της  $f$  ετά σημεία  $y_i$ ,  $1 \leq i \leq n$ . Χρησιμοποιώντας τώρα το θεώρημα του Weierstrass αποδείξτε το ζητούμενο).

(β) Χρησιμοποιώντας την ιδιότητα του  $p^*$  του μέρους (α) δείξτε για δεδομένο  $n$  και  $f \in C^n[a, b]$

$$\min_{p \in P_{n-1}} \|f - p\|_\infty = c(y_1, \dots, y_n) f^{(n)}(\xi)$$

για κάποιο  $\xi \in (a, b)$ , όπου  $c(y_1, \dots, y_n)$  σταθερά που εξαρτάται μόνο από τα σημεία  $y_i$ . Από την παρατήρηση 2 συμπερνάτε ότι

$$\min_{p \in P_{n-1}} \|f - p\|_\infty \geq 2((b-a)^n / 4^n n!) \min_{a \leq x \leq b} |f^{(n)}(x)|$$

5. (α) Δείξτε ότι η  $w(g;h)$  είναι μονοτονική και υπο-αθροιστική, δηλ. ότι  $w(g;h) \leq w(g;h+k) \leq w(g;h) + w(g;k)$  για  $h,k \geq 0$ ,  $g \in C[a,b]$ . Συμπεράνετε ότι  $w(g;\lambda h) \leq H_\lambda w(g;h)$  όπου για  $\lambda \geq 0$ , ο  $H_\lambda$  είναι ο ελάχιστος ακέραιος τέτοιος ώστε  $\lambda \leq H_\lambda$ .

(β) Αν  $w(g;h)/h \rightarrow 0$  όταν  $h \downarrow 0$ , δείξτε ότι η  $g$  είναι σταθερά.

6. Αποδείξτε τις (28) και (29) και ότι η ελάχιστη αυτή τιμή της  $\|w\|_\infty$  εμφανίζεται όταν  $x_i = x_i^c$ ,  $1 \leq i \leq n$ . (θυμηθείτε την παράσταση των πολυωνύμων Chebyshev στο  $[-1,1]$  και υπολογίστε τον ευτελεστικό του μεγιστοβάθμιου όρου του  $T_n$ )

7. Αποδείξτε ότι αν  $f \in C^\infty[a,b]$  τέτοια ώστε  $\|f^{(n)}\|_\infty \leq M^n$ ,  $n=0,1,2,\dots$  για κάποιο  $0 \leq M < \infty$ , τότε  $\|f - P_n f\|_\infty \rightarrow 0$ ,  $n \rightarrow \infty$  για οποιοδήποτε τριγωνικό εύστημα παρεμβολής.

8. Στο θεώρημα 1 δείξτε ότι η  $f^{(n)}(\xi(\bar{x}))$  είναι συνεχής συνάρτηση του  $\bar{x} \in [a,b]$ .

9. Δείξτε ότι για  $-1 \leq x_1 < x_2 < x_3 \leq 1$  η ελάχιστη τιμή της σταθεράς του Lebesgue  $\|\lambda_3\|_\infty$  ( $[a,b]=[-1,1]$ ) είναι  $5/4$  και ότι η τιμή αυτή λαμβάνεται για  $x_2=0$ ,  $-x_1=x_3 > 2\sqrt{2}/3$ . Δείξτε ότι  $\|\lambda_3^c\|_\infty = 5/3$ . Συνεπώς οι ρίζες των πολυωνύμων Chebyshev δεν δίνουν για κάθε  $n$  την ελάχιστη τιμή της σταθεράς του Lebesgue.

10. Στο μιγαδικό επίπεδο έστω  $p_n$  το μιγαδικό πολυώνυμο βαθμού  $\leq n$  που παρεμβάλλει την συνάρτηση  $f(z)=1/z$  στις  $n$ -ετές ρίζες της μονάδας. Δείξτε ότι  $p_n(z)=z^{n-1}$  και ότι

$$\max_{|z|=1} |p_n(z) - f(z)| \not\rightarrow 0, \quad n \rightarrow \infty$$

(11) (α) (Παρεμβολή Hermite). Δείξτε ότι υπάρχει μοναδικό πολυώνυμο  $q(x)$  βαθμού  $\leq 2n-1$  που ικανοποιεί τις σχέσεις  $q(x_i) = f(x_i)$ ,  $q'(x_i) = f'(x_i)$ ,  $1 \leq i \leq n$  για δεδομένη  $\tau = \{x_i\}$ ,  $1 \leq i \leq n$  και συνάρτηση  $f \in C^1[a, b]$ . Το  $q(x)$  λέγεται πολυώνυμο παρεμβολής Hermite για την  $f$  στα σημεία  $x_i$ .

(β) Δείξτε ότι

$$q(x) = \sum_{i=1}^n [f(x_i) A_i(x) + f'(x_i) B_i(x)]$$

όπου τα  $A_i, B_i \in \mathbb{P}_{2n-1}$  είναι τα (μοναδικά) πολυώνυμα βαθμού  $\leq 2n-1$  τέτοια ώστε  $A_i(x_j) = \delta_{ij}$ ,  $A_i'(x_j) = 0$ ,  $B_i(x_j) = 0$ ,  $B_i'(x_j) = \delta_{ij}$ ,  $1 \leq i, j \leq n$  και ότι δίνονται από τους τύπους

$$A_i(x) = (1 - 2(x - x_i) L_i'(x_i)) L_i^2(x)$$

$$B_i(x) = (x - x_i) L_i^2(x).$$

(γ) Αν  $\tau \subset [a, b]$  και  $f \in C^{2n}[a, b]$ , ισχύει το εξής αποτέλεσμα για το εφάλμα της παρεμβολής Hermite: για  $x \in [a, b]$  υπάρχει  $\xi = \xi_x \in (a, b)$  τέτοιο ώστε

$$f(x) - q(x) = \left\{ \prod_{i=1}^n (x - x_i)^2 \right\} f^{(2n)}(\xi_x) / (2n)!$$

## 4.2 ΠΑΡΕΜΒΟΛΗ ΚΑΙ ΠΡΟΣΕΓΓΙΣΗ ΜΕ ΤΗΜΗΜΑΤΙΚΑ ΓΡΑΜΜΙΚΕΣ ΣΥΝΑΡΤΗΣΕΙΣ

Αρχίζουμε την μελέτη της προσέγγισης με τμηματικά πολυωνυμικές συναρτήσεις με την απλούστερη περίπτωση, δηλ. με τις τμηματικά γραμμικές συναρτήσεις, που βέβαια δεν έχουν την πρακτική σημασία π.χ. των τμηματικά κυβικών συναρτήσεων (κυβικών splines). Η μελέτη μας όμως θα εστιασθεί σε προβλήματα (απλής μορφής εδώ) σημαντικά και για την προσέγγιση με τμηματικά πολυωνυμικές συναρτήσεις οποιουδήποτε βαθμού.

Έστω ο διαμερισμός  $\tau: a=x_1 < x_2 < \dots < x_N=b$  του  $[a,b]$  και  $f \in C[a,b]$ . Προφανώς υπάρχει μόνο μία συνάρτηση, γραμμική σε κάθε διάστημα  $[x_i, x_{i+1}]$  που παίρνει τις τιμές  $f(x_i)$  για  $i=1,2,\dots,N$ . Η "τεθλασμένη αυτή γραμμή" λέγεται συνάρτηση παρεμβολής της  $f$  στον χώρο των τμηματικά γραμμικών συναρτήσεων (που ορίζει η  $\tau$ ) και θα παριστάνεται με  $I_2 f$ . Προφανώς η  $I_2 f$  κατασκευάζεται πολύ εύκολα και δίνεται από

$$(1) (I_2 f)(x) = f(x_i) + (x-x_i) f[x_i, x_{i+1}] \text{ για } x_i \leq x \leq x_{i+1}, 1 \leq i \leq N-1.$$

Γιά να εκτιμήσουμε το σφάλμα  $e(x) = f(x) - (I_2 f)(x)$  παρατηρούμε ότι στο διάστημα  $[x_i, x_{i+1}]$ , η  $I_2 f$  είναι το γραμμικό πολυώνυμο παρεμβολής Lagrange της  $f(x)$  για τα σημεία  $x_i, x_{i+1}$ . Συνεπώς, αν  $f \in C^2[a,b]$  η (4.1.7) δίνει ότι για  $x \in [x_i, x_{i+1}]$ , υπάρχει  $\xi = \xi(x) \in (x_i, x_{i+1})$  τέτοιο ώστε

$$(2) e(x) = (x-x_i)(x-x_{i+1}) f^{(2)}(\xi)/2,$$

από την οποία έπεται ότι για  $x \in [x_i, x_{i+1}]$

$$|e(x)| \leq \max_{x_i \leq x \leq x_{i+1}} |(x-x_i)(x_{i+1}-x)| \max_{x_i \leq x \leq x_{i+1}} |f''(x)|/2.$$

### 4.2.2

Η μέγιστη τιμή της συνάρτησης  $(x-x_i)(x_{i+1}-x)$  για  $x \in [x_i, x_{i+1}]$  λαμβάνεται στο σημείο  $x=(x_i+x_{i+1})/2$  και είναι ίση με  $(x_{i+1}-x_i)^2/4$ .

Συνεπώς, η παραπάνω εκέση δίνει για  $f \in C^2[x_i, x_{i+1}]$

$$(3) \max_{x_i \leq x \leq x_{i+1}} |(f-I_2 f)(x)| \leq (x_{i+1}-x_i)^2 \max_{x_i \leq x \leq x_{i+1}} |f''(x)|/8, \quad 1 \leq i \leq N-1$$

που μάλιστα ισχύει ως ιδιότητα για την συνάρτηση  $f(x)=(x-x_i)(x_{i+1}-x)$  για την οποία  $(I_2 f)(x) \equiv 0$  στο  $[x_i, x_{i+1}]$ . Από την (3), για κάθε  $f \in C^2[a, b]$  έχουμε, με  $\|g\|_\infty = \max_{a \leq x \leq b} |g(x)|$ , ότι

$$(4) \|f-I_2 f\|_\infty \leq h^2 \|f''\|_\infty / 8,$$

όπου

$$(5) h = \max_{1 \leq i \leq N-1} (x_{i+1}-x_i)$$

Συνεπώς μπορούμε να κάνουμε το εφάλμα όσο μικρό θέλουμε παίρνοντας όσο και λεπτότερους διαμερισμούς  $\tau$  του  $[a, b]$ . Το πρόβλημα δεν γίνεται πιά πολύπλοκο καθώς  $h \rightarrow 0$  γιατί πάντα σε κάθε υποδιάστημα  $[x_i, x_{i+1}]$  η  $I_2 f$  είναι ευθύγραμμο τμήμα.

Για δεδομένο διαμερισμό  $\tau: a=x_1 < x_2 < \dots < x_N=b$  του  $[a, b]$  ορίζουμε τον διανυσματικό χώρο των (συνεχών) τεθλασμένων γραμμών (χώρο των γραμμικών splines)

$$S_\tau^2 = \{ \varphi: \varphi \in C[a, b], \varphi \in P_1(x_i, x_{i+1}), 1 \leq i \leq N-1 \},$$

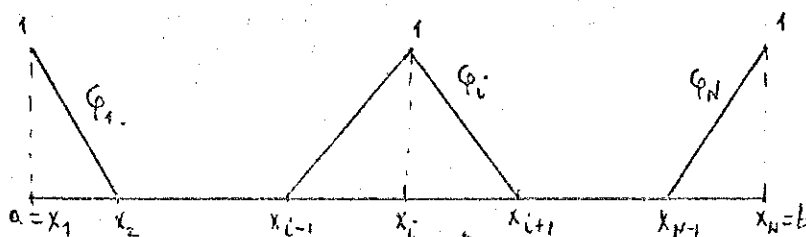
όπου με  $P_k(I)$  εννοούμε τον χώρο των πολυνομίων βαθμού  $\leq k$  στο διάστημα  $I$ . Προφανώς  $\forall f \in C[a, b]$  ορίζεται η  $I_2 f$  ως στοιχείο του  $S_\tau^2$ . Σε πολλές εφαρμογές είναι απαραίτητο να εκφράζουμε τα στοιχεία του

### 4.2.3

$S_T^2$  συναρτήσει μιάς εύχρηστης βάσης του διανυσματικού αυτού χώρου.

**Λήμμα 1.** Ο  $S_T^2$  είναι  $N$ -διάστατος υπόχωρος του  $C[a,b]$ . Οι συναρτήσεις  $\varphi_i \in S_T^2$ ,  $1 \leq i \leq N$  που ορίζονται για  $x \in [a,b]$  από τις σχέσεις  $\varphi_i(x_j) = \delta_{ij}$ , αποτελούν βάση του  $S_T^2$ .

Απόδειξη: Οι συναρτήσεις  $\varphi_i$  δίνονται από



για  $x \in [a,b]$  από τους τύπους

$$\varphi_1(x) = \begin{cases} (x_2 - x) / (x_2 - x_1), & \text{αν } x \in [x_1, x_2] \\ 0 & \text{αλλιώς,} \end{cases}$$

$$\varphi_N(x) = \begin{cases} (x - x_{N-1}) / (x_N - x_{N-1}) & \text{αν } x \in [x_{N-1}, x_N] \\ 0 & \text{αλλιώς} \end{cases}$$

(β) και για  $2 \leq i \leq N-1$ :  $\varphi_i(x) = \begin{cases} (x - x_{i-1}) / (x_i - x_{i-1}) & \text{αν } x \in [x_{i-1}, x_i] \\ (x_{i+1} - x) / (x_{i+1} - x_i) & \text{αν } x \in [x_i, x_{i+1}] \\ 0 & \text{αλλιώς} \end{cases}$

Προφανώς  $\varphi_i \in S_T^2$ ,  $1 \leq i \leq N$ ,  $\varphi_i(x_j) = \delta_{ij}$ ,  $1 \leq i, j \leq N$ . (θα αποδειχθεί σημαντικό αργότερα ότι οι  $\varphi_i$  έχουν μικρά φορέα:  $|\text{supp}(\varphi_i)| \leq 2h$ ). Είναι προφανές

ότι οι  $\varphi_i$  είναι γραμμικά ανεξάρτητες. Πράγματι, αν  $\sum_{i=1}^N c_i \varphi_i(x) = 0$



$\forall x \in [a, b]$ , τότε θέτοντας  $x = x_j$  για  $1 \leq j \leq N$  παίρνουμε  $c_j = 0$ . Είναι επίσης προφανές ότι παράγουν τον  $S_\tau^2$ , γιατί για κάθε  $\psi \in S_\tau^2$  έχουμε την παράσταση

$$(7) \quad \psi(x) = \sum_{i=1}^N \psi(x_i) \varphi_i(x).$$

(Το δεύτερο μέλος της (7) ανήκει  $S_\tau^2$  και για  $x = x_j$  έχει την τιμή  $\psi(x_j)$ .)

Συνοψώς οι δύο τεθλασμένες γραμμές  $\psi(x)$ ,  $\sum_{i=1}^N \psi(x_i) \varphi_i(x)$  ευπρίπτουν παντού στο  $[a, b]$ . @

Συνοψώς, η συνάρτηση παρεμβολής  $I_2 f$  της  $f \in C[a, b]$  στον  $S_\tau^2$  είναι το (μοναδικό) στοιχείο του  $S_\tau^2$  με παράσταση

$$(8) \quad (I_2 f)(x) = \sum_{i=1}^N f(x_i) \varphi_i(x), \quad 1 \leq i \leq N,$$

όπου οι συναρτήσεις βάσης  $\varphi_i$  δίνονται από τις εκθέσεις (6).

Ας ευχαρίσουμε τώρα το εφάλμα  $\|f - I_2 f\|_\infty$  της συνάρτησης παρεμβολής  $I_2 f \in S_\tau^2$  μιάς συνάρτησης  $f \in C[a, b]$  με το ελάχιστο εφάλμα  $\min\{\|f - \varphi\|_\infty : \varphi \in S_\tau^2\}$  που μπορούμε να πετύχουμε προερχόμενος την  $f$  με στοιχεία του  $S_\tau^2$  (ως προς την νόρμα  $\|\cdot\|_\infty$ ), δηλ. με το ελάχιστο εφάλμα βέλτιστης προσέγγισης της  $f$  στον  $(S_\tau^2, \|\cdot\|_\infty)$  (Ότι υπάρχει βέλτιστη προσέγγιση είναι απόρροια της πεπερασμένης διάστασης του  $S_\tau^2$ , βλ. π.χ. [4.9 εελ. 1]). Παρατηρούμε ότι ο τελεστής της παρεμβολής  $I_2: C[a, b] \rightarrow S_\tau^2$  είναι γραμμικός και ευπρίπτει με την ταυτότητα στον  $S_\tau^2$ , δηλ. ότι

$$(9) \quad I_2 \varphi = \varphi \quad \forall \varphi \in S_\tau^2.$$

Επιπλέον, είναι φραγμένος στον  $C[a, b]$  και μάλιστα

$$(10) \|I_2 g\|_\infty = \max_{1 \leq i \leq N} |(I_2 g)(x_i)| = \max_{1 \leq i \leq N} |g(x_i)| \leq \|g\|_\infty \quad \forall g \in C[a, b].$$

Από τις (9) και (10) έχουμε για  $f \in C[a, b]$  ότι

$$\begin{aligned} \|f - I_2 f\|_\infty &= \|(f - \varphi) - I_2(f - \varphi)\|_\infty \leq \|f - \varphi\|_\infty + \|I_2(f - \varphi)\|_\infty \\ &\leq \|f - \varphi\|_\infty + \|f - \varphi\|_\infty = 2\|f - \varphi\|_\infty \quad \forall \varphi \in S_\tau^2. \end{aligned}$$

Συνοπώς, έχουμε τελικά

$$(11) \min_{\varphi \in S_\tau^2} \|f - \varphi\|_\infty \leq \|f - I_2 f\|_\infty \leq 2 \min_{\varphi \in S_\tau^2} \|f - \varphi\|_\infty,$$

δηλ. ότι η  $I_2 f$  είναι "εχεδόν βέλτιστη" - με την έννοια ότι  $\|f - I_2 f\|_\infty \leq 2 \min_{\varphi \in S_\tau^2} \|f - \varphi\|_\infty$ ,  $c$  σταθερά ανεξάρτητη των  $\tau, N$ -προσέγγιση: Αν βρήκαμε μία βέλτιστη προσέγγιση της  $f$  στον  $(S_\tau^2, \|\cdot\|_\infty)$ , το πολύ-πολύ να υποδιπλασιάσαμε το εφάλμα της  $I_2 f$ .

Προχωρούμε τώρα σε μία πιο λεπτομερή μελέτη του εφάλματος της παρεμβολής στον χώρο  $S_\tau^2$ . Θα χρησιμοποιήσουμε το λεγόμενο θεώρημα του πυρήνα Peano σε μία απλή μορφή του. Γενικεύοντας λίγο τον χώρο  $C^k[a, b]$ , θα θεωρήσουμε, για  $k \geq 1$  ακέραιο, τον χώρο  $PC^k[a, b]$ , που ορίζεται ως ο χώρος εκείνων των συναρτήσεων  $f$  του  $C^{k-1}[a, b]$  των οποίων η  $k$ -ετή παράγωγος υπάρχει σε όλα τα σημεία του  $[a, b]$ , εκτός πιθανώς από ένα πεπερασμένο εύλογο σημείων στο  $[a, b]$  (ανάμεσα στα οποία και στα  $a, b$  είναι συνεχής), είναι φραγμένη στο  $[a, b]$  και ορίζεται σε κάθε σημείο του  $[a, b]$  μονοσήμαντα σαν όριο τιμών της από δεξιά ή αριστερά, λόγω συνέχειας.

**ΘΕΩΡΗΜΑ 1** (Πορήια Ρεαπο). Έστω  $E:PC^{n+1}[a,b] \rightarrow \mathbb{R}^1$ ,  $n \geq 0$  ένα γραμμικό συναρτησιακό (δηλ. έστω ότι  $E(\lambda f + \mu g) = \lambda E(f) + \mu E(g)$   $\forall f, g \in PC^{n+1}[a,b]$ ,  $\lambda, \mu \in \mathbb{R}^1$ ), τέτοιο ώστε  $E(p) = 0 \quad \forall p \in \mathbb{P}_n$ . Τότε για κάθε  $f \in PC^{n+1}[a,b]$  έχουμε την παράσταση

$$(12) \quad E(f) = (n!)^{-1} E_x \left[ \int_a^b (x-t)_+^{n-1} f^{(n+1)}(t) dt \right],$$

όπου  $(x-t)_+^n = \begin{cases} (x-t)^n & x \geq t \\ 0 & x < t \end{cases}$

και όπου με  $E_x$  τονίζουμε το γεγονός ότι το  $E$  δρα στον όρο μέσα στις αγκύλες θεωρούμενο ως συνάρτηση του  $x$ .

Απόδειξη: Η συνάρτηση

$$G_n(x) \equiv \int_a^b (x-t)_+^{n-1} f^{(n+1)}(t) dt = \int_a^x (x-t)^{n-1} f^{(n+1)}(t) dt,$$

για  $f \in PC^{n+1}[a,b]$ , ανήκει στον χώρο  $PC^{n+1}[a,b]$  όπως εύκολα μπορούμε να δούμε παραγωγίζοντας το δεύτερο μέλος ως προς  $x$ . Από το θεώρημα του Taylor με την ολοκληρωτική μορφή υπολοίπου, έχουμε για  $x \in [a,b]$  επειδή  $f \in PC^{n+1}[a,b]$

$$f(x) = f(a) + (x-a)f'(a) + \dots + (x-a)^n f^{(n)}(a)/n! + G_n(x)/n!$$

Η (12) τώρα προκύπτει άμεσα επειδή το  $E$  είναι γραμμικό συναρτησιακό και επειδή λόγω της υπόθεσής μας,  $E_x((x-a)^k) = 0$  για  $0 \leq k \leq n$ . @

**ΘΕΩΡΗΜΑ 2.** Έστω  $f \in C^2[a,b]$  και έστω  $e(x) = f(x) - (I_2 f)(x)$  το εφάλμα της παρεμβολής στον  $S_t^2$ . Για  $x_i \leq x \leq x_{i+1}$  με  $h_i = (x_{i+1} - x_i)$  έχουμε

$$(13) \quad e(x) = \int_{x_i}^{x_{i+1}} K(x,t) f^{(2)}(t) dt,$$

όπου

$$(14) \quad K(x,t) = \begin{cases} -(x_{i+1}-x)(t-x_i)/h_i, & x_i \leq t \leq x \leq x_{i+1}, \\ -(x-x_i)(x_{i+1}-t)/h_i, & x_i \leq x < t \leq x_{i+1}. \end{cases}$$

και

$$(15) \quad e'(x) = \int_{x_i}^{x_{i+1}} \Lambda(x,t) f^{(2)}(t) dt,$$

όπου

$$(16) \quad \Lambda(x,t) = \begin{cases} (t-x_i)/h_i, & x_i \leq t \leq x \leq x_{i+1}, \\ -(x_{i+1}-t)/h_i, & x_i \leq x < t \leq x_{i+1}. \end{cases}$$

Απόδειξη: Σταθεροποιούμε ένα  $x$  στο  $[x_i, x_{i+1}]$ . Τότε το  $e(x) = e_x(f) = f(x) - I_2 f(x)$  είναι ένα συναρτησιακό στον χώρο  $C^2[x_i, x_{i+1}]$  για το οποίο  $e_x(p) = 0 \quad \forall p \in P_1[x_i, x_{i+1}]$ . Εφαρμόζοντας λοιπόν το θεώρημα του Πυρήνα του Peano στο  $[x_i, x_{i+1}]$  για  $n=1$  έχουμε

$$(17) \quad e(x) = \epsilon_x \left[ \int_{x_i}^{x_{i+1}} (x-t)_+ f^{(2)}(t) dt \right] = \int_{x_i}^{x_{i+1}} (x-t)_+ f^{(2)}(t) dt - I_{2,x} \left[ \int_{x_i}^x (x-t) f^{(2)}(t) dt \right], \quad x \in [x_i, x_{i+1}]$$

Εξ ορισμού της συνάρτησης παρεμβολής  $I_{2g}$  στο  $[x_i, x_{i+1}]$  (είναι η γραμμική συνάρτηση στο  $[x_i, x_{i+1}]$  με  $(I_{2g})(x_k) = g(x_k)$ ,  $k=i, i+1$ ), βλέπουμε ότι για  $x \in [x_i, x_{i+1}]$

$$I_{2,x} \left[ \int_{x_i}^x (x-t) f^{(2)}(t) dt \right] = ((x-x_i)/h_i) \int_{x_i}^{x_{i+1}} (x_{i+1}-t) f^{(2)}(t) dt.$$

Η (17) και η παραπάνω σχέση δίνουν τις (13)-(14). Γράφοντας τώρα

$$e(x) = \int_{x_i}^x K(x,t) f^{(2)}(t) dt + \int_x^{x_{i+1}} K(x,t) f^{(2)}(t) dt,$$

παρατηρώντας ότι λόγω της συνέχειας της  $K(x,t)$  για  $x=t$

$$\begin{aligned} (d/dx) \left( \int_{x_i}^x K(x,t) f^{(2)}(t) dt \right) &= \int_{x_i}^x K_x(x,t) f^{(2)}(t) dt + K(x,x) f^{(2)}(x-) \\ (d/dx) \left( \int_x^{x_{i+1}} K(x,t) f^{(2)}(t) dt \right) &= \int_x^{x_{i+1}} K_x(x,t) f^{(2)}(t) dt - K(x,x) f^{(2)}(x+) \end{aligned}$$

και χρησιμοποιώντας ότι  $f \in C^2[a,b]$  παίρνουμε τις (15)-(16) γιατί η συνάρτηση  $\Lambda(x,t)$  είναι ίση με  $K_x(x,t)$  τμηματικά, για  $x_i < t < x$  και  $x < t < x_{i+1}$ . @

Με βάση τις παραστάσεις (13), (15) μπορούμε να βρούμε εκτιμήσεις του εφάλματος της παρεμβολής και της παραγώγου του στην νόρμα  $\|\cdot\|_\infty$  αλλά και σε άλλες νόρμες. Ιδιαίτερα μας ενδιαφέρει η  $L^2$ -νόρμα, που για τετραγωνικά ολοκληρώσιμες πραγματικές συναρτήσεις στο  $[a,b]$  θα συμβολίζεται ως

$$(18) \quad \|f\| = \left( \int_a^b f^2(x) dx \right)^{1/2}$$

**ΠΟΡΙΣΜΑ 1.** Έστω  $f \in C^2[a,b]$  και  $h = \max(x_{i+1} - x_i)$ ,  $1 \leq i \leq N-1$ . Για το εφάλμα  $e(x) = f(x) - I_2 f(x)$  της παρεμβολής στον  $S_T^2$  έχουμε

$$(19) \quad \|e\|_\infty \leq h^2 \|f''\|_\infty / 8$$

$$(20) \quad \|e'\|_\infty \leq h \|f''\|_\infty / 2$$

$$(21) \quad \|e\| \leq h^2 \|f''\| / 3\sqrt{10}$$

$$(22) \quad \|e'\| \leq h \|f''\| / \sqrt{6}$$

Απόδειξη:

Οι (13)-(14) δίνουν, για  $x \in [x_i, x_{i+1}]$ ,  $h_i = x_{i+1} - x_i$

$$|e(x)| \leq \left( \int_{x_i}^{x_{i+1}} |K(x,t)| dt \right) \|f''\|_\infty = \|f''\|_\infty h_i^{-1} \left[ \int_{x_i}^x (t-x_i)(x_{i+1}-x) dt \right.$$

$$\left. + \int_x^{x_{i+1}} (x-x_i)(x_{i+1}-t) dt \right]$$

$$= \|f''\|_\infty h_i^{-1} [(x_{i+1}-x)(x-x_i)^2 + (x-x_i)(x-x_{i+1})^2] / 2$$

$$= \|f\|_{\infty} (x_{i+1}-x)(x-x_i)/2 \leq \|f\|_{\infty} h_i^2/8$$

όπως προηγουμένως. Συνεπώς ισχύει η (19) την οποία είχαμε αποδείξει και προηγουμένως χρησιμοποιώντας την "εμφιακή" παράσταση του εφάλματος.

Οι (15)-(16) δίνουν τώρα για  $x \in [x_i, x_{i+1}]$

$$\begin{aligned} |e'(x)| &\leq \|f''\|_{\infty} \int_{x_i}^{x_{i+1}} |A(x,t)| dx = \\ &= \|f''\|_{\infty} h_i^{-1} \left[ \int_{x_i}^x (t-x_i) dt + \int_x^{x_{i+1}} (x_{i+1}-t) dt \right] \\ &= \|f''\|_{\infty} h_i^{-1} [(x-x_i)^2 + (x_{i+1}-x)^2]/2 \leq \|f''\|_{\infty} h_i/2 \end{aligned}$$

από την οποία προκύπτει η (20) με προφανή επέκταση της νόρμας  $\| \cdot \|_{\infty}$  στις συναρτήσεις του χώρου  $PC^0[a,b]$  των τμηματικά συνεχώς φραγμένων συναρτήσεων στον οποίο ανήκει η  $e'(x)$ .

Από τις (13)-(14) τώρα, με χρήση της ανισότητας Cauchy-Schwarz έχουμε, για  $x \in [x_i, x_{i+1}]$  ότι

$$(23) \quad \varepsilon^2(x) \leq \left( \int_{x_i}^{x_{i+1}} K^2(x,t) dt \right) \left( \int_{x_i}^{x_{i+1}} (f''(t))^2 dt \right).$$

Εύκολα βλέπουμε ότι

$$\begin{aligned} \int_{x_i}^{x_{i+1}} K^2(x,t) dt &= h_i^{-2} [(x_{i+1}-x)^2 \int_{x_i}^x (t-x_i)^2 dt + (x-x_i)^2 \int_x^{x_{i+1}} (x_{i+1}-t)^2 dt] \\ &= (3h_i)^{-1} (x_{i+1}-x)^2 (x-x_i)^2 \end{aligned}$$

Άρα, ολοκληρώνοντας και τα δύο μέλη της (23) ως προς  $x$  από  $x_i$  έως  $x_{i+1}$  έχουμε

## 4.2.11

$$\int_{x_i}^{x_{i+1}} e^2(x) dx \leq (3h_i)^{-1} \left[ \int_{x_i}^{x_{i+1}} (x_{i+1}-x)^2 (x-x_i)^2 dx \right] \int_{x_i}^{x_{i+1}} (f''(t))^2 dt$$

$$= (h_i^4/90) \int_{x_i}^{x_{i+1}} (f''(t))^2 dt$$

Συνεπώς,

$$\|e\|^2 = \int_a^b e^2(x) dx = \sum_{i=1}^{N-1} \int_{x_i}^{x_{i+1}} e^2(x) dx$$

$$\leq (h^4/90) \sum_{i=1}^{N-1} \int_{x_i}^{x_{i+1}} (f''(t))^2 dt = (h^4/90) \|f''\|^2$$

από την οποία έπεται η (21).

Τέλος, οι (15) - (16) δίνουν για  $x \in [x_i, x_{i+1}]$

$$(e'(x))^2 \leq \left( \int_{x_i}^{x_{i+1}} \Lambda^2(x,t) dt \right) \left( \int_{x_i}^{x_{i+1}} (f''(t))^2 dt \right),$$

και επειδή

$$\int_{x_i}^{x_{i+1}} \Lambda^2(x,t) dt = h_i^{-2} \left[ \int_{x_i}^x (t-x_i)^2 dt + \int_x^{x_{i+1}} (t-x_{i+1})^2 dt \right]$$

$$= [(x-x_i)^3 + (x_{i+1}-x)^3] / 3h_i^2$$



έχουμε ότι

$$\int_{x_i}^{x_{i+1}} (\varepsilon'(x))^2 dx \leq (3^{-1}h_i^{-2}) \int_{x_i}^{x_{i+1}} [(x-x_i)^3 + (x_{i+1}-x)^3] dx \left( \int_{x_i}^{x_{i+1}} (f''(t))^2 dt \right) \\ \leq (h_i^2/6) \int_{x_i}^{x_{i+1}} (f''(t))^2 dt$$

από την οποία προκύπτει η (22) αν αθροίσουμε ως προς  $i$  και πάρουμε την τετραγωνική ρίζα και των δύο μελών. @

Αν η συνάρτηση  $f$  είναι λιγότερο ομαλή τότε, γενικά, δεν περιμένουμε τάξη ακρίβειας π.χ.  $O(h^2)$  στην (19) αλλά μικρότερη. Θα εξετάσουμε μερικές τέτοιες περιπτώσεις στις Παρατηρήσεις και στις Ηεκκείσεις. Ας προχωρήσουμε όμως τώρα στην μελέτη ενός άλλου προβλήματος, δηλ. στην προσέγγιση μίας συνεχούς συνάρτησης  $f$  από στοιχεία του  $S_T^2$  με την έννοια των ελαχίστων τετραγώνων.

Θεωρούμε στον διανυσματικό χώρο  $C[a,b]$  το εσωτερικό γινόμενο

$$(24) \quad (f, g) = \int_a^b f(x)g(x)dx,$$

το οποίο βέβαια παράγει την  $L^2$ -νόρμα  $\|\cdot\|$ , (18). Επειδή ο χώρος  $S_T^2$  είναι υπόχωρος του  $(C[a,b], (\cdot, \cdot))$  πεπερασμένης διαστάσεως, είναι γνωστό (βλ. Ηεκκείση 1α) ότι το πρόβλημα του προσδιορισμού ενός  $f_h \in S_T^2$  που να ικανοποιεί

$$(25) \quad \|f - f_h\| = \min_{\varphi \in S_T^2} \|f - \varphi\|,$$

δηλ. το πρόβλημα της βέλτιστης προσέγγισης της  $f$  από στοιχεία του  $S_T^2$  για την νόρμα  $\|\cdot\|$  (της λεγόμενης προσέγγισης ελαχίστων τετραγώνων

της  $f$  από στοιχεία του  $S_T^2$ ), έχει μοναδική λύση  $f_h$  που είναι επίσης και η μοναδική λύση του προβλήματος των κανονικών εξισώσεων

$$(26) (f_h, \varphi) = (f, \varphi) \quad \forall \varphi \in S_T^2.$$

Από την (26) π.χ. προκύπτει ότι υπάρχει γραμμικός τελεστής  $P: C[a, b] \rightarrow S_T^2$ , ο λεγόμενος τελεστής της ορθής προβολής (ως προς το εσωτερικό γινόμενο (24)) ή της  $L^2$ -προβολής στον  $S_T^2$  τέτοιος ώστε

$$(27) f_h = Pf.$$

Η (26) γράφεται και ως  $(f - Pf, \varphi) = 0 \quad \forall \varphi \in S_T^2$ . Δηλ. το εφάλμα  $f - Pf$  της  $L^2$ -προβολής της  $f$  στον  $S_T^2$  είναι ορθογώνιο προς του  $S_T^2$ . Από την ορθογωνιότητα αυτή προκύπτει το Πυθαγόρειο θεώρημα, δηλ. ότι  $\|Pf - f\|^2 = \|f\|^2 - \|Pf\|^2$ . Αμα ισχύει

$$(28) \|Pf\| \leq \|f\| \quad \forall f \in C[a, b].$$

Για να εκτιμήσουμε το εφάλμα της  $L^2$ -προβολής ως προς την νόρμα  $\|\cdot\|$ , έχουμε από την (25) ότι

$$(29) \|f - Pf\| = \min_{\varphi \in S_T^2} \|f - \varphi\| \leq \|f - I_2 f\|,$$

από την οποία, χρησιμοποιώντας τα γνωστά μας φράγματα για το εφάλμα της παρεμβολής στην νόρμα  $\|\cdot\|$ , (π.χ. για  $f \in C^2[a, b]$  έχουμε από την (21) ότι  $\|f - I_2 f\| \leq C h^2 \|f''\|$ ) βρίσκουμε ανάλογα φράγματα για το εφάλμα  $\|f - Pf\|$ .

Ο υπολογισμός της  $Pf$  γίνεται από τις κανονικές εξισώσεις (26) που προφανώς είναι ισοδύναμες με τις

$$(30) (Pf, \varphi_i) = (f, \varphi_i), \quad 1 \leq i \leq N,$$

όπου,  $\varphi_i$ ,  $1 \leq i \leq N$ , οι συναρτήσεις βάσης που κατασκευάστηκαν στο Λήμμα

1. Υποθέτοντας ότι

$$(31) \quad Pf = \sum_{i=1}^N c_i \varphi_i,$$

βλέπουμε ότι οι συντελεστές  $c_i$ ,  $1 \leq i \leq N$ , της  $Pf$  ικανοποιούν το  $N \times N$  γραμμικό σύστημα

$$(32) \quad Gc = \beta,$$

όπου  $G \in \mathbb{R}^{N \times N}$ , ο λεγόμενος πίνακας Gram ή πίνακας μάζας της βάσης  $\{\varphi_i\}$ , ορίζεται ως

$$(33) \quad G_{ij} = (\varphi_i, \varphi_j), \quad 1 \leq i, j \leq N,$$

και είναι προφανώς συμμετρικός και θετικά ορισμένος. Στην (32)  $c = [c_1, \dots, c_N]^T$  και  $\beta = [\beta_1, \dots, \beta_N]^T$ , όπου

$$(34) \quad \beta_i = (f, \varphi_i), \quad 1 \leq i \leq N.$$

Η σημασία του μικρού φορέα των συναρτήσεων βάσης  $\varphi_i$  είναι τώρα προφανής: Από την (33) βλέπουμε ότι επειδή  $(\varphi_i, \varphi_j) = 0$  για  $|i-j| > 1$ , ο πίνακας  $G$  είναι αραιός και μάλιστα τριδιαγώνιος. Μερικοί απλοί υπολογισμοί των ολοκληρωμάτων

$$\int_a^b \varphi_i^2 dx, \quad \int_a^b \varphi_i \varphi_{i+1} dx \text{ δίνουν, για } h_j = x_{j+1} - x_j, \quad 1 \leq j \leq N-1,$$



Έστω  $j$  ( $1 \leq j \leq N$ ) δείκτης για τον οποίο  $|c_j| = \max_{0 \leq i \leq N+1} |c_i|$ . Τότε η (37) για  $i=j$  δίνει

$$2|c_j| = |\hat{\beta}_j - (h_j + h_{j-1})^{-1}(h_{j-1}c_{j-1} + h_j c_{j+1})| \leq |\hat{\beta}_j| + (h_j + h_{j-1})^{-1}(h_{j-1} + h_j) \max(|c_{j-1}|, |c_{j+1}|) \leq |\hat{\beta}_j| + |c_j|,$$

δηλ. ότι

$$(38) \quad |c_j| = \max_i |c_i| \leq |\hat{\beta}_j| = \delta(h_j + h_{j-1})^{-1} \left| \int_a^b f \varphi_j dx \right| \\ \leq \delta(h_j + h_{j-1})^{-1} \|f\|_\infty \int_a^b \varphi_j dx = 3 \|f\|_\infty.$$

Επειδή  $Pf \in S_\tau^2$ ,  $\|Pf\|_\infty = \max_{1 \leq i \leq N} |(Pf)(x_i)| = \max_{1 \leq i \leq N} |c_i| = |c_j|$ .

Από την (38) προκύπτει συνεπώς η (36). @

Όπως έχουμε δει και προηγουμένως (βλ. και Άσκηση 5), η "ευετόθεια" (36) της  $L^2$ -προβολής στον  $S_\tau^2$  ως προς την νόρμα  $\|\cdot\|_\infty$  και το γεγονός ότι  $P\varphi = \varphi$  για  $\varphi \in S_\tau^2$  μας οδηγεί στην εξής εκτίμηση για κάθε  $f \in C[a, b]$  και  $\varphi \in S_\tau^2$ :

$$\|Pf - f\|_\infty \leq \|Pf - \varphi\|_\infty + \|\varphi - f\|_\infty = \|P(f - \varphi)\|_\infty + \|f - \varphi\|_\infty \leq 4\|f - \varphi\|_\infty,$$

από την οποία τελικά έχουμε

$$(39) \quad \|Pf - f\|_\infty \leq 4 \min_{\varphi \in S_\tau^2} \|f - \varphi\|_\infty \quad (\text{π.χ. } \leq 4\|f - I_2 f\|_\infty).$$

Έτσι, αν  $f \in C^2[a, b]$  π.χ., παίρνουμε, λόγω της (19), ότι  $\|Pf - f\|_\infty \leq Ch^2 \|f''\|_\infty$ .

### Παρατηρήσεις

1. Έστω ότι  $f \in PC^2[a, b]$ . Από το θεώρημα του πυρήνα του Peano και την απόδειξη του θεωρήματος 2 βλέπουμε ότι εξακολουθεί να ισχύει η παράσταση (13) για το εφάλμα  $e(x) = f(x) - (I_2 f)(x)$  της παρεμβολής, όπου ο πυρήνας  $K(x, t)$  δίνεται πάλι από την (14). Συνεπώς ερμηνεύοντας με επέκταση κατά τον προφανή τρόπο τις νόρμες  $\|f''\|_\infty, \|f''\|$  για  $f \in PC^2[a, b]$  - και ακολουθώντας τα βήματα της απόδειξης του Πορίσματος 1, βλέπουμε ότι ισχύουν οι (19) και (21) και για  $f \in PC^2[a, b]$ . Για να αποδείξουμε φράγματα για το εφάλμα της παραγωγής, βλέπουμε τώρα από την απόδειξη του θεωρήματος 1 ότι για  $f \in PC^2[a, b]$ ,  $x \in [x_i, x_{i+1}]$

$$(40) \quad e'(x) = \int_{x_i}^{x_{i+1}} \lambda(x, t) f^{(2)}(t) dt + K(x, x) [f^{(2)}(x^-) - f^{(2)}(x^+)],$$

όπου ο πυρήνας  $\lambda$  δίνεται από την (16) και ο  $K$  από την (14). (Ο  $K(x, t)$  είναι συνεχής για  $x=t$ ). Χρησιμοποιώντας την (40) εύκολα μπορούμε να δούμε ότι για  $x \in [x_i, x_{i+1}]$

$$|e'(x)| \leq \max(\|f''\|_{L^\infty(x_i, x_{i+1})}, |f^{(2)}(x^+) - f^{(2)}(x^-)|) h_i / 2$$

από την οποία παίρνουμε π.χ. ότι  $\|e'\|_\infty \leq Ch \|f''\|_\infty$ . Πέλοος, μιά ανισότητα της μορφής  $\|e'\| \leq Ch \|f''\|$  (δηλ. το ανάλογο της (22) μπορεί να αποδειχθεί όπως στην Άσκηση 10β με  $C=1$ ).

2. Ας εξετάσουμε τώρα περιπτώσεις και λιγότερο ομαλών  $f$ , υποθέτοντας απλώς ότι  $f \in C[a, b]$ . Από τον τοπικό ορισμό της ευνάρτησης παρεμβολής  $I_2 f$  στο  $[x_i, x_{i+1}]$  εύκολα προκύπτει η παράσταση για  $x \in [x_i, x_{i+1}]$ :

$$e(x) = f(x) - (I_2 f)(x) = h_i^{-1} [(x_{i+1} - x)(f(x) - f(x_i)) + (x - x_i)(f(x) - f(x_{i+1}))].$$

Συμπεραίνουμε λοιπόν ότι για  $x \in [x_i, x_{i+1}]$ ,  $h = \max_i h_i$

$$(41) |e(x)| \leq |h_i^{-1} [(x_{i+1} - x) + (x - x_i)]| \max_{x \in [x_i, x_{i+1}]} (|f(x) - f(x_i)|, |f(x_{i+1}) - f(x)|)$$

$$\leq \omega(f; h_i) \leq \omega(f; h).$$

Συνοπώς επειδή  $f \in C[a, b]$ ,  $\|e\|_\infty \leq \omega(f; h) \rightarrow 0$ ,  $h \rightarrow 0$ , και μάλιστα  $\|e\|_\infty = O(h)$  αν η  $f$  είναι Lipschitz με ειδική περίπτωση το φράγμα  $\|e\|_\infty \leq \|f'\|_\infty h$  για  $f \in C^1[a, b]$ . (Το θεώρημα του Πυρήνα του Peano δίνει όμως το καλύτερο φράγμα  $\|e\|_\infty \leq \|f''\|_\infty h^2/2$ ,  $f \in PC^2[a, b]$ , βλ. Ασκ. 4). Εξ άλλου, το θεώρημα της μέσης τιμής δίνει π.χ.  $\|e'\|_\infty \rightarrow 0$ ,  $h \rightarrow 0$  αν  $f \in C^1[a, b]$ . Η τάξη εύγκλισης όμως μπορεί μερικές φορές να βελτιωθεί με κατάλληλη επιλογή των  $x_i$  (πύκνωση κοντά στις ιδιομορφίες της  $f$ ). Βλ. π.χ. [4.2, σελ. 46].

3. Οι σταθερές  $(1/8$  και  $1/2)$  και οι εκθέτες του  $h$  ( $2$  και  $1$ ) στις σχέσεις (19) και (20) είναι οι καλύτεροι δυνατοί για  $f \in C^2[a, b]$  π.χ.. Αντίθετα οι σταθερές στα φράγματα (21) και (22) (αλλά όχι όμως οι εκθέτες!) μπορούν να βελτιωθούν: Οι καλύτερες σταθερές στα φράγματα  $\|e^{(k)}\| \leq C_k h^{2-k} \|f^{(k)}\|$ ,  $k=0, 1$ ,  $f \in C^2[a, b]$  μπορεί να δείξει ότι είναι  $C_0 = \pi^{-2}$ ,  $C_1 = \pi^{-1}$ , βλ. π.χ. [4.10, παρ. 2, 3].

### Ασκήσεις 4.2

1(α) Επιβεβαιώστε τους ισχυρισμούς περί της  $L^2$ -προβολής στον  $S_T^2$  στο κείμενο και τους τύπους (25)-(28).

(β) Από την ανισότητα  $\max_i |c_i| \leq \max_i |\hat{\beta}_i|$  που προκύπτει από την (38), συμπεράνετε ότι το σύστημα (32) έχει μοναδική λύση.

2. Αποδείξτε τους ισχυρισμούς της Παρατήρησης 1.

3. Αποδείξτε ότι η συνάρτηση παρεμβολής  $I_2 f$  μίας συνεχούς συνάρτησης  $f$  διατήρει ιδιότητες μονοτονίας ή κυρτότητας της  $f$ .

4. Χρησιμοποιώντας το θεώρημα του πυρήνα του Peano δείξτε ότι

$$\|f - I_2 f\|_\infty \leq h \|f'\|_\infty / 2 \text{ για } f \in C^1[a, b].$$

5. Για  $g \in X$ , όπου  $X$  κάποιος διανυσματικός χώρος συναρτήσεων, έστω  $\Pi g$  προσέγγιση του  $g$  ε' ένα υπόχωρο  $S$  του  $X$  με τις ιδιότητες:

$$(i) \Pi \varphi = \varphi, \quad \forall \varphi \in S$$

$$(ii) \Pi(f+g) = \Pi f + \Pi g, \quad \forall f, g \in X$$

$$(iii) \|\Pi g\| \leq \infty, \text{ όπου } \|\Pi g\| = \sup_{\varphi \in S, \varphi \neq 0} \|\Pi g\| / \|\varphi\|,$$

και όπου  $\|\cdot\|$  κάποια νόρμα του  $X$ .

Δείξτε τότε ότι για κάθε  $g \in X$

$$\|g - \Pi g\| \leq (1 + \|\Pi g\|) \inf_{\varphi \in S} \|g - \varphi\|.$$

6. Έστω  $S_T^1$  ο διανυσματικός χώρος των τμηματικά σταθερών συναρτήσεων στο  $[a, b]$  ως προς τον διαμερισμό  $\tau: a = x_1 < x_2 < \dots < x_N = b$ ,

δηλ. έστω



$$S_T^1 = \{\varphi: \varphi(x) = c_i \quad \text{για } x \in [x_i, x_{i+1}], \quad 1 \leq i \leq N-1\}.$$

(Για να αποφύγουμε αριθμούς των  $\varphi \in S^1$  στους κόμβους ας συμβολίσουμε  $\|\varphi\|_\infty = \max \{|\varphi(x)|: x \in [a, b] - \{x_2, \dots, x_{N-1}\}\}$

(α) Δείξτε ότι για  $g \in C[a, b]$

$$\min_{\varphi \in S_T^1} \|g - \varphi\| \leq \omega(g; h/2),$$

όπου κατά τα συνηθιστά  $h = \max_i h_i$ .

(β) Δείξτε ότι  $\|g' - (I_2 g)'\|_\infty \leq 2 \min_{\varphi \in S_T^1} \|g' - \varphi\|$

για κάθε  $g \in C^1[a, b]$  και μάλιστα ότι  $\|g' - (I_2 g)'\|_\infty = O(h)$  για Lipschitz  $g'$ .

(γ) ("Υπερσύγκλιση" στον  $S_T^1$ ). Έστω  $f \in PC^1[a, b]$  και θεωρήστε την συνάρτηση  $I_1 f \in S_T^1$  που ορίζεται ως

$$(I_1 f)(x) = f((x_i + x_{i+1})/2) \quad \text{για } x \in [x_i, x_{i+1}], \quad 1 \leq i \leq N-1. \quad \text{Δείξτε ότι}$$

$$\|f - I_1 f\|_\infty \leq h \|f'\|_\infty / 2.$$

7. ("Υπερσύγκλιση" για παραγώγους στον  $S_T^2$ ). Υποθέστε ότι  $f \in C^3[a, b]$ . Ξεκινώντας από την παράσταση (15), (16) του εφάλατος  $e'(x) = f'(x) - (I_2 f)'(x)$ ,  $x \in [x_i, x_{i+1}]$  και χρησιμοποιώντας ολοκλήρωση κατά μέρη δείξτε ότι

$$|e'((x_i + x_{i+1})/2)| \leq h^2 \|f'''\|_{\omega} / 6, \quad 1 \leq i \leq N-1,$$

δηλ. ότι υπάρχουν σημεία (τα μέσα των διαστημάτων  $[x_i, x_{i+1}]$ ) στο  $[a, b]$  όπου το σημειακό εφάλμα της παραγωγής  $(I_2 f)'$  είναι της τάξης  $O(h^2)$ , (για  $f \in C^3[a, b]$ ), ενώ, όπως ξέρουμε, ισχύει γενικά ότι

$$\|e'\|_{\omega} \leq h \|f''\|_{\omega} / 2 \text{ έστω και για } f \in C^m[a, b].$$

8. Δείξτε ότι για κάθε  $f \in PC^1[a, b]$  ισχύει ότι

$$(\alpha) \quad ((I_2 f)', \varphi') = (f', \varphi') \quad \forall \varphi \in \mathcal{S}_h.$$

$$(\beta) \quad \|(I_2 f)'\|^2 + \|(I_2 f - f)'\|^2 = \|f'\|^2$$

$$(\gamma) \quad \text{Αν } f \in PC^2[a, b] \text{ δείξτε ότι } \|(f - I_2 f)'\|^2 = (f - I_2 f, f'')$$

9. (Αμειωτότητα των Poincaré-Friedrichs). Δείξτε ότι για κάθε  $f \in PC^2[c, d]$  τέτοια ώστε  $f(c) = 0$  ή  $f(d) = 0$  ισχύει

$$\int_c^d f^2(x) dx \leq (d-c)^2 \int_c^d (f'(x))^2 dx.$$

10(α). Χρησιμοποιώντας τα αποτελέσματα του ασκήσεων 8 και 9 δείξτε ότι αν  $f \in PC^1[a, b]$ , τότε

$$\|(f - I_2 f)'\| \leq \|f'\|$$

$$\|(f - I_2 f)\| \leq h \|f'\|$$

## 4.2.22

(β) Αν  $PC^2[a,b]$ -δείξτε, χρησιμοποιώντας τα παραπάνω ότι

$$\|(f - I_2 f)'\| \leq h \|f''\|$$

και

$$\|f - I_2 f\| \leq h^2 \|f''\|$$

(Στην άσκηση αυτή μην χρησιμοποιήσετε το θεώρημα του Πυρήνα του Peano).

### 4.3 ΠΑΡΕΜΒΟΛΗ ΜΕ ΤΗΜΑΤΙΚΑ ΚΥΒΙΚΕΣ ΣΥΝΑΡΤΗΣΕΙΣ ΗΕΡΜΙΤΕ

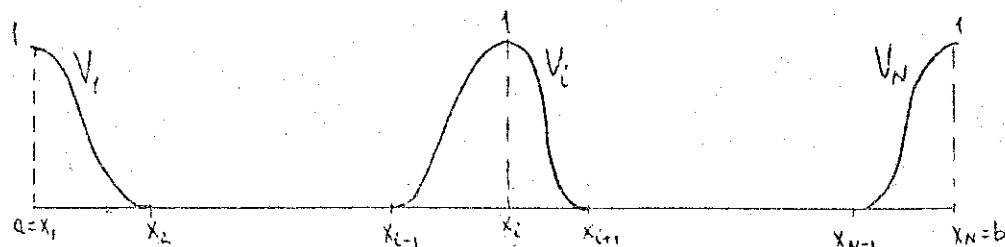
Στην παράγραφο αυτή και στην επόμενη θα ασχοληθούμε με παρεμβολή σε χώρους τμηματικά κυβικών συναρτήσεων. Υπάρχουν αρκετοί τρόποι με τους οποίους μπορεί να γίνει παρεμβολή σε τέτοιους χώρους. Θα αρχίσουμε εδώ με την λεγόμενη παρεμβολή Hermite (που χρησιμοποιεί δηλ. εκτός από τις τιμές  $f(x_i)$  της συνάρτησης στους κόμβους και τις τιμές της παραγώγου της  $f'(x_i)$ ) σε χώρους τμηματικά κυβικών συναρτήσεων που ανήκουν στον χώρο  $C^1[a,b]$ .

Για δεδομένο διαμερισμό  $\tau: a = x_1 < x_2 < \dots < x_N = b$  του  $[a,b]$  θεωρούμε τον διανυσματικό χώρο των (τμηματικά) κυβικών συναρτήσεων Hermite.

$$H_\tau = \{ \varphi : \varphi \in C^1[a,b], \varphi \in P_3(x_i, x_{i+1}), 1 \leq i \leq N-1 \}$$

**Λήμμα 1.** Ο  $H_\tau$  είναι  $2N$ -διάστατος υπόχωρος του  $C^1[a,b]$ . Οι συναρτήσεις  $\{U_i(x), S_i(x)\}, 1 \leq i \leq N$ , που δίνονται από τους τύπους (1), (2), παρακάτω, αποτελούν βάση του  $H_\tau$ .

Βπόδειξη. Προφανώς ο  $H_\tau$  είναι διανυσματικός χώρος, υπόχωρος του  $C^1[a,b]$ . Συμβολίζοντας πάλι  $h_i = x_{i+1} - x_i, 1 \leq i \leq N-1$ , θεωρούμε τις συναρτήσεις  $U_i, S_i \in H_\tau, 1 \leq i \leq N$ , που δίνονται από

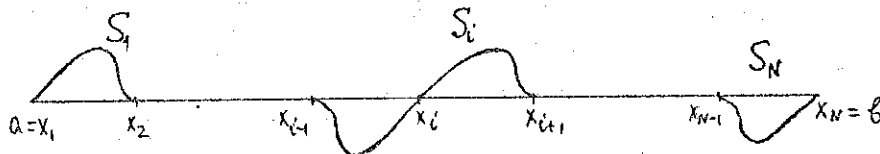


$$(1) \quad U_i(x) = \begin{cases} 2h_i^{-3}(x-x_1)^3 - 3h_i^{-2}(x-x_1)^2 + 1 & \text{αυ } x \in [x_1, x_2] \\ 0 & \text{αλλιώς} \end{cases}$$

Για  $2 \leq i \leq N-1$ 

$$U_i(x) = \begin{cases} -2h_{i-1}^{-3}(x-x_{i-1})^3 + 3h_{i-1}^{-2}(x-x_{i-1})^2 & x \in [x_{i-1}, x_i] \\ 2h_i^{-3}(x-x_i)^3 - 3h_i^{-2}(x-x_i)^2 + 1 & x \in [x_i, x_{i+1}] \\ 0 & \text{αλλοιώς} \end{cases}$$

$$U_N(x) = \begin{cases} -2h_{N-1}^{-3}(x-x_{N-1})^3 + 3h_{N-1}^{-2}(x-x_{N-1})^2 & x \in [x_{N-1}, x_N] \\ 0 & \text{αλλοιώς} \end{cases}$$



και

$$(2) S_1(x) = \begin{cases} h_1^{-2}(x-x_1)(x_2-x)^2 & x \in [x_1, x_2] \\ 0 & \text{αλλοιώς} \end{cases}$$

Για  $2 \leq i \leq N-1$ 

$$S_i(x) = \begin{cases} h_{i-1}^{-2}(x-x_{i-1})^2(x-x_i) & x \in [x_{i-1}, x_i] \\ h_i^{-2}(x-x_i)^2(x_{i+1}-x) & x \in [x_i, x_{i+1}] \\ 0 & \text{αλλοιώς} \end{cases}$$

$$S_N(x) = \begin{cases} h_{N-1}^{-2}(x-x_{N-1})^2(x-x_N) & x \in [x_{N-1}, x_N] \\ 0 & \text{αλλοιώς} \end{cases}$$

## 4.3.3

Είναι προφανές ότι  $S_i, U_i \in H_\tau$ ,  $1 \leq i \leq N$  και ότι ικανοποιούν τις σχέσεις

$$(3) \quad \begin{cases} U_i(x_j) = \delta_{ij}, & U_i'(x_j) = 0 \\ S_i(x_j) = 0, & S_i'(x_j) = \delta_{ij} \end{cases}, \quad 1 \leq i, j \leq N.$$

Επιπλέον οι  $S_i, U_i$  είναι οι συναρτήσεις του  $H_\tau$  με τον μικρότερο δυνατό φορέα που ικανοποιούν τις (3). Επίσης είναι γραμμικά ανεξάρτητες. Πράγματι, αν  $\sum_{j=1}^N c_j U_j(x) + d_j S_j(x) = 0$   $\forall x \in [a, b]$  θα έχουμε και  $\sum_{j=1}^N c_j U_j'(x) + d_j S_j'(x) = 0 \quad \forall x \in [a, b]$ .

Θέτοντας  $x = x_i$ ,  $1 \leq i \leq N$  στην πρώτη σχέση παίρνουμε  $c_i = 0$ ,  $1 \leq i \leq N$ , ενώ από την δεύτερη  $d_i = 0$ ,  $1 \leq i \leq N$ . Είναι εύκολο επίσης να δούμε ότι το εύρος  $\{U_i, S_i\}$  παράγει τον  $H_\tau$ . Πράγματι, για  $\varphi \in H_\tau$  ισχύει

$$(4) \quad \varphi(x) = \sum_{i=1}^N \varphi(x_i) U_i(x) + \varphi'(x_i) S_i(x),$$

όπως εύκολα μπορούμε να δούμε. Για  $x = x_k$ ,  $1 \leq k \leq N$ , οι τιμές και οι τιμές των παραγώγων των δύο μελών της (4) προφανώς ευρύνονται. Σε οποιοδήποτε διάστημα  $[x_k, x_{k+1}]$ ,  $1 \leq k \leq N-1$ , και τα δύο μέλη της (4) είναι λοιπών κυβικά πολυώνυμα των οποίων οι τιμές και οι τιμές των παραγώγων ευρύνονται στα άκρα  $x_k, x_{k+1}$  του διαστήματος. Συνεπώς θα ευρύνονται και για κάθε  $x \in [x_k, x_{k+1}]$  επειδή κάθε κυβικό πολυώνυμο ορίζεται μονοσήμαντα από τις τιμές του και τις τιμές της παραγώγου του σε δύο διακριτά σημεία. @

Η σχέση (4) μας λέει ότι μπορούμε να ορίσουμε μονοσήμαντα οποιαδήποτε συνάρτηση  $\varphi \in H_\tau$  αν δίνονται οι τιμές  $\varphi(x_i), \varphi'(x_i)$ . Για οποιαδήποτε λοιπόν συνάρτηση  $f \in C^1[a, b]$  π.χ., μπορούμε να κατασκευάσουμε την συνάρτηση παρεμβολής της  $I_H f$  στον χώρο  $H_\tau$  ως το (μοναδικό) στοιχείο του  $H_\tau$  που ικανοποιεί τις σχέσεις

## 4.3.4

$$(5) \quad (I_H f)(x_i) = f(x_i), \quad (I_H f)'(x_i) = f'(x_i), \quad 1 \leq i \leq N.$$

Η κατασκευή της  $I_H f$  απαιτεί γνώση των τιμών  $f(x_i)$ ,  $f'(x_i)$ · μάλιστα, λόγω της (4) η  $I_H f$  δίνεται από

$$(6) \quad (I_H f)(x) = \sum_{i=1}^N f(x_i) U_i(x) + f'_i(x) S_i(x)$$

Χρησιμοποιώντας το θεώρημα του Peano είναι δυνατόν να εκτιμήσουμε το εφάλμα της προσέγγισης  $f - I_H f$  και των παραγώγων της ως προς διάφορες νόρμες ευαρτήσεων στο  $[a, b]$ . Θα αποδείξουμε εδώ μέρος του παρακάτω θεωρήματος:

**ΘΕΩΡΗΜΑ 1.** Έστω  $f \in C^4[a, b]$  και  $h = \max_i h_i$ . Τότε

$$(7) \quad \|f - I_H f\|_\infty \leq h^4 \|f^{(4)}\|_\infty / 384$$

$$(8) \quad \|(f - I_H f)'\|_\infty \leq (\sqrt{3}/216) h^3 \|f^{(4)}\|_\infty$$

$$(9) \quad \|(f - I_H f)^{(2)}\|_\infty \leq h^2 \|f^{(4)}\|_\infty / 12$$

$$(10) \quad \|(f - I_H f)^{(3)}\|_\infty \leq h \|f^{(4)}\|_\infty / 2$$

όπου η δεύτερη και η τρίτη "παράγωγος" της  $I_H f$  έχουν υπολογισθεί κατά τμήματα, δηλ. όπου στις (9) και (10) εννοούμε ότι για  $k=2, 3$ :

$$\|(f - I_H f)^{(k)}\|_\infty = \max_{1 \leq i \leq N-1} \left( \max_{x_i \leq x \leq x_{i+1}} |(f - I_H f)^{(k)}(x)| \right)$$

Απόδειξη. Θα αποδείξουμε την (7) και θα εστιάσουμε απλώς τις (8)-(10) παραπέμποντας για την απόδειξη στην βιβλιογραφία.

## 4.3.5

Σταθεροποιούμε ένα  $x \in (x_i, x_{i+1})$ ,  $[i \leq N-1]$ . Το εφάλμα τότε της παρεμβολής Hermite  $e(x) = e_f(x) = f(x) - I_H f(x)$  είναι γραμμικό συναρτησιακό στον χώρο  $C^4[a, b]$ . — Επιπλέον, επειδή  $p = I_H p \quad \forall p \in P_3(x_i, x_{i+1})$ , από το θεώρημα του πυρήνα του Peano 4.2.1 παίρνουμε την παράσταση

$$(11) \quad e(x) = (3!)^{-1} e_x \left[ \int_{x_i}^{x_{i+1}} (x-t)_+^3 f^{(4)}(t) dt \right]$$

όπου ανακαλούμε ότι

$$(x-t)_+^3 = \begin{cases} (x-t)^3 & \text{αν } x \geq t \\ 0 & \text{αν } x < t \end{cases}$$

Από τον ορισμό του τελεστού  $I_H$  της παρεμβολής βλέπουμε ότι οι

$$\text{πράξεις } e_x \text{ και } \int_{x_i}^{x_{i+1}} dt \text{ αντιμετατίθενται, δηλ. ότι}$$

$$(12) \quad e(x) = 6^{-1} \int_{x_i}^{x_{i+1}} e_x [(x-t)_+^3] f^{(4)}(t) dt = 6^{-1} \int_{x_i}^{x_{i+1}} K(x, t) f^{(4)}(t) dt,$$

όπου

$$(13) \quad K(x, t) = (x-t)_+^3 - I_{H, x} [(x-t)_+^3], \quad x_i \leq x, t \leq x_{i+1}.$$

Από τις (12), (13) συμπεραίνουμε ότι

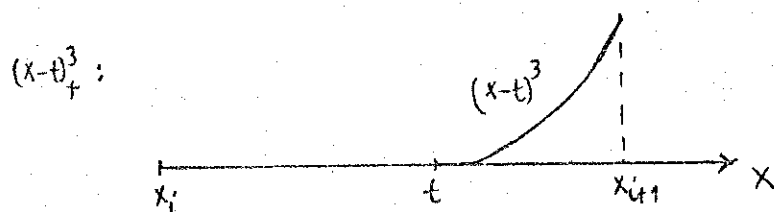
$$(14) \quad |e(x)| \leq \|f^{(4)}\|_\infty \int_{x_i}^{x_{i+1}} |K(x, t)| dt / 6.$$



### 4.3.6

Συνεπώς πρέπει να βρούμε μία όσο το δυνατό καλύτερη εκτίμηση του ολοκληρώματος στο δεύτερο μέλος της (14).

Για να υπολογίσουμε τον πυρήνα  $K(x,t)$  από την (13) πρέπει να βρούμε την συνάρτηση  $I_{H,x}[(x-t)_+^3]$ , για  $x \in [x_i, x_{i+1}]$ , δηλ. το κυβικό πολυώνυμο  $p(x)$  του οποίου η τιμή και η τιμή της παραγώγου ευρπύτουν με τις αντίστοιχες τιμές της συνάρτησης  $(x-t)_+^3$ , θεωρούμενης ως συνάρτησης του  $x$ , στα σημεία  $x=x_i$  και  $x=x_{i+1}$ ,



δηλ. το κυβικό πολυώνυμο  $p(x)$  που ικανοποιεί  $p(x_i) = p'(x_i) = 0$ ,  $p(x_{i+1}) = (x_{i+1} - t)^3$ ,  $p'(x_{i+1}) = 3(x_{i+1} - t)^2$ . Εύκολα βρίσκουμε ότι

$$(15) \quad p(x) = I_{H,x}[(x-t)_+^3] \equiv \Lambda(x,t) = h_i^{-3} (x-x_i)^2 (x_{i+1}-t)^2 [x(x_{i+1}-3x_i+2t) + tx_i - 3tx_{i+1} + 2x_i x_{i+1}], \quad x_i \leq x, t \leq x_{i+1}.$$

Άρα ο πυρήνας  $K(x,t)$ , λόγω των (13), (15) γίνεται

$$(16) \quad K(x,t) = \begin{cases} (x-t)^3 - \Lambda(x,t), & x_i \leq t \leq x \leq x_{i+1} \\ -\Lambda(x,t), & x_i \leq x \leq t \leq x_{i+1} \end{cases}$$

Από τις (15), (16) βλέπουμε ότι  $K(x,t) \geq 0 \quad \forall x, t \in [x_i, x_{i+1}]$ . Ας εξετάσουμε πρώτα την περίπτωση  $x_i < x < t < x_{i+1}$ . Για τέτοια  $x, t$  έχουμε, λόγω των (15), (16) ότι  $\text{sgn } K(x,t) = -\text{sgn } \Lambda(x,t) = -\text{sgn } \Pi(x,t)$ , όπου  $\Pi(x,t) = x(x_{i+1} - 3x_i + 2t) + tx_i - 3tx_{i+1} + 2x_i x_{i+1}$ .

## 4.3.7

Η  $M(x, t)$  είναι γραμμική συνάρτηση του  $t$  και ικανοποιεί  $M(x, x) = -2h_i(x - x_i) < 0$  και  $M(x, x_{i+1}) = -3(x_{i+1} - x)h_i < 0$ . Άρα  $M(x, t) < 0$  για  $x \leq t \leq x_{i+1}$ . Συμπεραίνουμε ότι  $K(x, t) \geq 0$  για  $x_i \leq x \leq t \leq x_{i+1}$ . Μιά σειρά από απλούς υπολογισμούς δίνει επίσης ότι  $K(x, t) \geq 0$  για  $x_i \leq x \leq t \leq x_{i+1}$ .

Συνοπώς έχουμε

$$\begin{aligned} \int_{x_i}^{x_{i+1}} |K(x, t)| dt &= \int_{x_i}^{x_{i+1}} K(x, t) dt = \int_{x_i}^x [(x-t)^3 - \Lambda(x, t)] dt + \int_x^{x_{i+1}} [-\Lambda(x, t)] dt \\ &= \int_{x_i}^x (x-t)^3 dt - \int_{x_i}^{x_{i+1}} \Lambda(x, t) dt = (x-x_i)^4/4 + (x-x_i)^2 [4(xx_{i+1} - 3xx_i + 2x_i x_{i+1}) \\ &\quad + (2x+x_i - 3x_{i+1})(x_{i+1} + 3x_i)]/12 = (x-x_i)^2(x_{i+1}-x)^2/4. \end{aligned}$$

Άρα η (14) δίνει

$$(17) |e(x)| \leq \|f^{(4)}\|_{\infty} (x-x_i)^2 (x_{i+1}-x)^2/24, \quad x \in [x_i, x_{i+1}].$$

Το μέγιστο της συνάρτησης  $(x-x_i)^2 (x_{i+1}-x)^2$ ,  $x \in [x_i, x_{i+1}]$ , προφανώς λαμβάνεται για  $x = (x_i + x_{i+1})/2$  και είναι ίσο με  $h_i^4/16$ . Συνοπώς η (7) έπεται από την (17). Όπως δείχνει η συνάρτηση  $f(x) = (x-x_i)^2 (x_{i+1}-x)^2$  (για την οποία  $I_H f = 0$  για  $x \in [x_i, x_{i+1}]$ ,  $f^{(4)}(x) = 24$ ), η (17) - και συνεπώς και η (7) - δεν είναι δυνατόν να βελτιωθεί.

## 4.3.8

Για να αποδείξουμε τις (8)-(10) παρατηρούμε πρώτα (Ασκ. 1β) ότι για  $x \in (x_i, x_{i+1})$

$$(18) \quad e^{(k)}(x) = 6^{-1} \int_{x_i}^{x_{i+1}} \partial_x^k [K(x,t)] f^{(4)}(t) dt, \quad k=1,2,3,$$

όπου με  $\partial_x^k [K(x,t)]$ ,  $1 \leq k \leq 3$ , συμβολίζουμε την μερική παράγωγο  $k$  τάξεως ως προς  $x$  της συνάρτησης  $K(x,t)$ , ορισμένη τμηματικά λόγω της (16) για  $x < t$  και  $x > t$ . Δεν ισχύει όμως πιά ότι οι συναρτήσεις  $\partial_x^k [K(x,t)]$  έχουν

εταθερό πρόσημο για  $x, t \in (x_i, x_{i+1})$ . Οι Birkhoff και Priver (J. Math. and Phys. 46 (1967), 440-447), βρίσκουν τις ρίζες καθώς και τα σημεία όπου οι  $\partial_x^k [K(x,t)]$  αλλάζουν ασυνεχώς πρόσημο (ως συναρτήσεις

του  $t$ ). Τα ολοκληρώματα  $\int_{x_i}^{x_{i+1}} |\partial_x^k [K(x,t)]| dt$  γράφονται κατόπιν

ως αθροίσματα ολοκληρωμάτων πάνω σε διαστήματα όπου οι  $\partial_x^k [K(x,t)]$

έχουν εταθερό πρόσημο και υπολογίζονται ακριβώς. Προκύπτουν οι (8)-(10). Μπορεί να αποδειχθεί ότι οι εταθερές και οι δυνάμεις του  $h$  στις (7)-(10) είναι οι καλύτερες δυνατές για  $f \in C^k[a,b]$ ,  $k \geq 4$ . @

### Παρατηρήσεις

-1- Μπορούμε να δείξουμε βέβαια και φράγματα εφαλμάτων στην

$L^2$ -νόρμα  $\|f\| = (\int_a^b f^2(x) dx)^{1/2}$ . Π.χ. ισχύει (βλ. Ασκ. 4)

$$\|(f - I_H f)^{(k)}\| \leq C_k h^{2-k} \|f^{(k)}\|, \quad k=0,1,2, \text{ αν } f \in PC^2[a,b]$$

και

$$\|(f - I_H f)^{(k)}\| \leq C_k^* h^{4-k} \|f^{(4)}\|, \quad k=0,1,2, \text{ αν } f \in PC^4[a,b].$$

Οι καλύτερες εταθερές είναι  $C_0 = n^{-2}, C_1 = n^{-1}, C_2 = 1, C_0^* = n^{-4}, C_1^* = n^{-3}, C_2^* = n^{-2}$ .

### 4.3.9

Ανάλογες εκτιμήσεις ισχύουν αν π.χ.  $f \in PC^3[a,b]$  και φυσικά και στον νόρμα  $\| \cdot \|_\infty$  αν η  $f$  είναι λιγότερο ομαλή από  $C^4$ .

2. Κυβική παρεμβολή Bessel. Μία παραλλαγή της παρεμβολής με κυβικές συναρτήσεις Hermite είναι η λεγόμενη κυβική παρεμβολή Bessel, την οποία μπορούμε να χρησιμοποιήσουμε αν γνωρίζουμε μόνο τις τιμές της συνάρτησης  $f$  στα  $x_i$ ,  $1 \leq i \leq N$  και τις τιμές της παραγώγου  $f'$  μόνο στα  $a$  και στα  $b$ . Αντί  $f'(x_i)$ ,  $2 \leq i \leq N-1$  δηλ., στον τύπο (6) χρησιμοποιούμε αριθμούς  $s_i$  που είναι οι κλίσεις στα σημεία  $x_i$  του πολυωνύμου παρεμβολής βαθμού  $\leq 2$  που εμφανεί με τις τιμές της  $f$  στα σημεία  $x_{i-1}, x_i, x_{i+1}$ . Εύκολα βλέπουμε ότι οι κλίσεις αυτές δίνονται από τους τύπους

$$s_i = (h_i f[x_{i-1}, x_i] + h_{i-1} f[x_i, x_{i+1}]) / (h_i + h_{i-1}), \quad 2 \leq i \leq N-1.$$

Η κυβική παρεμβολή Bessel είναι και αυτή τοπική μέθοδος, όπως και η παρεμβολή Hermite. Δηλ. για του προεδιορισμό της τιμής της συνάρτησης παρεμβολής  $(If)(x)$  ε' ένα σημείο  $x \in [x_i, x_{i+1}]$  χρησιμοποιούνται πληροφορίες για τις τιμές της  $f$  (ή και της  $f'$ ) σε λίγα σημεία  $x_k$  κοντά στο  $x$ . Όπως αποδεικνύεται όμως (θεκ. 2) η κυβική παρεμβολή Bessel έχει γενικά εφάλμα της τάξης  $O(h^3)$  και όχι  $O(h^4)$  όπως η παρεμβολή Hermite.

### Θεκήσεις 4.3

1(α) Δείξτε ότι ο πυρήνας  $K(x,t)$  που ορίζεται από την (16) είναι μη αρνητικός και για  $x_i \leq t \leq x_{i+1}$ .

(β) Δείξτε ότι ισχύει η (18) για  $f \in C^4[a,b]$ .

## 4.3.10

2. Έστω  $f$  μία αρκετά ομαλή συνάρτηση στο  $[a, b]$  και έστω  $f_h$  οποιαδήποτε τμηματικά (ως προς οποιοδήποτε διαμερισμό  $\tau$  του  $[a, b]$ ) κυβική συνάρτηση παρεμβολής της  $f$ , δηλ. έστω  $f_h \in P_3(x_i, x_{i+1})$ ,  $1 \leq i \leq N-1$  με  $f_h(x_i) = f(x_i)$ ,  $1 \leq i \leq N$ .

(α) Χρησιμοποιώντας το γεγονός ότι ένα κυβικό πολυώνυμο στο  $[x_i, x_{i+1}]$  προσδιορίζεται μονοσήμαντα από τις τιμές του και τις τιμές της παραγώγου του στα άκρα  $x_i, x_{i+1}$  του διαστήματος, δείξτε ότι το σφάλμα της  $f_h$  γράφεται ως

$$e = f - f_h = f - I_H f + e$$

όπου  $I_H f$  η συνάρτηση παρεμβολής Hermite της  $f$  και όπου

$$e(x) = h_i^2 [e'(x_i)(x_{i+1} - x) - e'(x_{i+1})(x - x_i)](x - x_i)(x_{i+1} - x),$$

$$x_i \leq x \leq x_{i+1}, \quad 1 \leq i \leq N-1.$$

(β) Αποδείξτε ότι για αρκετά ομαλή  $f$  έχουμε για  $h = \max_i (x_{i+1} - x_i)$

$$\|f - f_h\|_\infty \leq \|f - I_H f\|_\infty + h \max_i |e''(x_i)|/4, \text{ δηλ. ότι}$$

$$\|f - f_h\|_\infty = O(h^4) + \max_i |e''(x_i)| O(h).$$

(γ) Έστω  $f_h$  η κυβική συνάρτηση παρεμβολής Bessel για την  $f$  με  $f_h'(a) = f'(a)$ ,  $f_h'(b) = f'(b)$  (βλ. παρατήρηση 2). Δείξτε ότι

$$\|f - f_h\|_\infty = O(h^3) \text{ και ότι } \|f - f_h\|_\infty = O(h^4) \text{ για ομοιόμορφο διαμερισμό,}$$

υποθέτοντας ότι η  $f$  είναι όσο ομαλή χρειάζεται.

## 4.3.11

3. (α) Αν  $f \in PC^2[a, b]$  δείξτε ότι

$$((I_H f - f)'', \varphi'') = 0, \quad \forall \varphi \in H_{\tau, H}$$

(β) Αν  $f \in PC^2[a, b]$  τότε

$$\|(I_H f)''\|^2 + \|(I_H f - f)''\|^2 = \|f''\|^2$$

(γ) Αν  $f \in PC^4[a, b]$  τότε  $\|(f - I_H f)''\|^2 = (f - I_H f, f^{(4)})$ .

4. Χρησιμοποιώντας τις "ολοκληρωτικές ταυτότητες" της Άσκησης 3 και την ανισότητα του Poincaré (Άσκηση 4.2.9) δείξτε,

$$\|(f - I_H f)^{(k)}\| \leq h^{2-k} \|f^{(2)}\|, \quad k=0, 1, 2 \quad \text{αν } f \in PC^2[a, b]$$

και

$$\|(f - I_H f)^{(k)}\| \leq h^{4-k} \|f^{(4)}\|, \quad k=0, 1, 2 \quad \text{αν } f \in PC^4[a, b].$$

## 4.4 ΠΑΡΕΜΒΟΛΗ ΜΕ ΚΥΒΙΚΕΣ SPLINES

Οι τμηματικά κυβικές συναρτήσεις Hermite που εξετάσαμε στην προηγούμενη παράγραφο αποτελούν διανυσματικό χώρο διαστάσεως  $2N$  για δεδομένο διαμερισμό  $\tau$  του  $[a, b]$  με  $N$  διακριτά σημεία  $x_i$ . Συνεπώς ο αριθμός των παραμέτρων που ορίζουν μία συνάρτηση του  $H_\tau$  είναι αρκετά μεγάλος σε σχέση π.χ. με τις  $N$  παραμέτρους που προσδιορίζουν μία συνάρτηση του  $S_\tau^2$ . (Οι δυσκολίες αυξάνουν αν θεωρήσουμε τα πολυδιάστατα ανάλογα αυτών των χώρων για προσέγγιση συναρτήσεων πολλών μεταβλητών). Επιπλέον η κατασκευή της συνάρτησης παρεμβολής  $I_H f$  απαιτεί, εκτός των  $f(x_i)$ , και τις κλίσεις  $f'(x_i)$  που πολλές φορές είναι δύσκολο ή αδύνατο να υπολογισθούν ακριβώς σε πρακτικά προβλήματα.

Τίθεται λοιπόν το ερώτημα αν μπορούμε να κατασκευάσουμε ένα χώρο τμηματικά κυβικών συναρτήσεων με όσο το δυνατόν μικρότερη διάσταση στον οποίο να μπορούμε να λύσουμε αποτελεσματικά το πρόβλημα της παρεμβολής χρησιμοποιώντας εί δυνατόν μόνο τιμές της συνάρτησης  $f(x_i)$  και διατηρώντας εφάλματα παρεμβολής τάξης  $O(h^4)$  για ομαλές  $f$ . Η απάντηση είναι καταφατική και δίνεται από τον χώρο των λεγόμενων κυβικών splines που ορίζεται, για του διαμερισμό  $\tau: a=x_1 < x_2 < \dots < x_N=b$ , ως το εύρος

$$S_\tau^4 = \{ \varphi : \varphi \in C^2[a, b], \varphi \in P_3(x_i, x_{i+1}), 1 \leq i \leq N-1 \},$$

δηλ. ως το εύρος  $S_\tau^4 = H_\tau \cap C^2[a, b]$ .

Αρχίζουμε μελετώντας το πρόβλημα της παρεμβολής στον  $S_\tau^4$ . Εστω λοιπόν μία συνάρτηση  $f \in C[a, b]$ , έστω  $N > 2$  και ας δεχθούμε ότι υπάρχει συνάρτηση  $s(x) = (I_4 f)(x) \in S_\tau^4$  τέτοια ώστε  $s(x_i) = f(x_i)$ ,  $1 \leq i \leq N$ . Θα δείξουμε ότι τότε ισχύουν οι σχέσεις

## 4.4.2

$$(1) \quad h_j s'(x_{j-1}) + 2(h_j + h_{j-1})s'(x_j) + h_{j-1}s'(x_{j+1}) = \\ = 3[h_{j-1}(f(x_{j+1}) - f(x_j))/h_j + h_j(f(x_j) - f(x_{j-1}))/h_{j-1}], \quad 2 \leq j \leq N-1,$$

όπου, όπως και προηγουμένως, χρησιμοποιούμε τον συμβολισμό  $h_j = x_{j+1} - x_j$ ,  $1 \leq j \leq N-1$ .

Πράγματι, όπως εύκολα βλέπουμε, το μοναδικό πολυώνυμο  $p \in P_3[\lambda, \mu]$  που για δεδομένα  $u_1, u_2, u_1', u_2'$  ικανοποιεί τις συνθήκες  $p(\lambda) = u_1, p(\mu) = u_2, p'(\lambda) = u_1', p'(\mu) = u_2'$ , δίνεται, αν  $h = \mu - \lambda$ , από τον τύπο

$$p(x) = u_1 [h^{-2}(x-\mu)^2 + 2h^{-3}(x-\lambda)(x-\mu)^2] + u_2 [h^{-2}(x-\lambda)^2 - 2h^{-3}(x-\lambda)^2(x-\mu)] \\ + u_1' h^{-2}(x-\lambda)(x-\mu)^2 + u_2' h^{-2}(x-\lambda)^2(x-\mu).$$

(Το  $p(x)$  είναι το κυβικό πολυώνυμο Hermite με δύο σημεία στο διάστημα  $[\lambda, \mu]$ ). Συνεπώς, αν γνωρίζαμε τις τιμές  $s'(x_j), s'(x_{j+1})$  θα μπορούσαμε να κατασκευάσουμε την  $s(x)$  από τον τύπο

$$(2) \quad s(x) = f(x_j) [h_j^{-2}(x-x_{j+1})^2 + 2h_j^{-3}(x-x_j)(x-x_{j+1})^2] + f(x_{j+1}) [h_j^{-2}(x-x_j)^2 \\ - 2h_j^{-3}(x-x_{j+1})(x-x_j)^2] + s'(x_j) h_j^{-2}(x-x_j)(x-x_{j+1})^2 + \\ s'(x_{j+1}) h_j^{-2}(x-x_j)^2(x-x_{j+1}), \quad x \in [x_j, x_{j+1}].$$

Παραγωγίζοντας στην (2) παίρνουμε

$$s''(x_j^+) = 2[3h_j^{-2}(f(x_{j+1}) - f(x_j)) - h_j^{-1}(s'(x_{j+1}) + 2s'(x_j))], \quad 1 \leq j \leq N-1.$$





## 4.4.4

ο οποίος έχει αυστηρά κυριαρχική διαχώνιο,, δηλ. ικανοποιεί τις εκθέσεις  $|b_{ii}| > \sum_{j, j \neq i} |b_{ij}|$ ,  $1 \leq i \leq N-2$ . Τέτοιοι πίνακες είναι ως γνωστόν αντιστρέψιμοι (από το θεώρημα του Gerschgorin, βλ. π.χ. [5.4]). Συνεπώς, οι (1) και οι ευνοριακές συνθήκες  $s'(x_1) = f'(x_1)$ ,  $s'(x_N) = f'(x_N)$  ορίζουν μονοσήμαντα τις κλίσεις  $s'(x_j)$ ,  $1 \leq j \leq N$ . Η συνάρτηση  $s(x) = (I_4 f)(x)$ , που ορίζεται τώρα σε κάθε διάστημα  $[x_j, x_{j+1}]$  από την (2), ικανοποιεί λοιπόν τις συνθήκες  $(I_4 f)(x_j) = f(x_j)$ ,  $1 \leq j \leq N$ ,  $(I_4 f)'(x_k) = f'(x_k)$ ,  $k=1, N$ . Αν υπήρχε και άλλη συνάρτηση  $\psi(x)$  με τις ίδιες ιδιότητες στο  $S_c^4$  οι κλίσεις της  $\psi'(x_j)$  θα ικανοποιούσαν τό ίδιο γραμμικό σύστημα με τις  $s'(x_j)$  και συνεπώς θα είχαμε  $\psi'(x_j) = s'(x_j)$ ,  $1 \leq j \leq N$ . Επειδή  $\psi(x_j) = f(x_j)$ ,  $1 \leq j \leq N$ , η παράσταση (2) για την  $\psi(x)$  θα ήταν ταυτόσημη με την παράσταση της  $s(x)$ , δηλ.  $\psi = s$ . @

Δεν είναι δύσκολο να δούμε ότι η κατασκευή του πίνακα  $B$  και του δεύτερου μέλους -βλ. (1)- του συστήματος για τον προσδιορισμό των  $s'(x_j)$  καθώς και η επίλυση του τριδιαχώνιου συστήματος (με τον γνωστό μας "τριδιαχώνιο αλγόριθμο" που είναι ευεταθής στην περίπτωση μας γιατί ο  $B$  έχει αυστηρά κυριαρχική διαχώνιο) απαιτεί  $15N+O(1)$  περίπου πράξεις. Η τιμή  $s(x)$  της κυβικής spline βρίσκεται κατόπιν από την (2) για δεδομένο  $x$ . Βλέπουμε δηλ. ότι ο προσδιορισμός της κυβικής spline παρεμβολής  $I_4 f$  απαιτεί λύση γραμμικού συστήματος, δηλ. ότι η τιμή της  $(I_4 f)(x)$  εξαρτάται τελικά από όλες τις τιμές  $f(x_j)$ ,  $1 \leq j \leq N$ ,  $f'(x_1)$ ,  $f'(x_N)$  των δεδομένων, σε αντίθεση π.χ. με τις συναρτήσεις παρεμβολής  $I_2 f$ ,  $I_H f$  που για  $x \in [x_j, x_{j+1}]$  εξαρτώνται μόνο τοπικά από την συνάρτηση  $f$ . (Η  $(I_2 f)(x)$  από τις τιμές  $f(x_j)$ ,  $f(x_{j+1})$  και η  $(I_H f)(x)$  από τις  $f(x_j)$ ,  $f'(x_j)$ ,  $f(x_{j+1})$ ,  $f'(x_{j+1})$ ).

Προχωρούμε τώρα στην εκτίμηση του εφάλματος της παρεμβολής  $e(x) = f(x) - (I_4 f)(x)$ . Η εκτίμηση στηρίζεται στο εξής αποτέλεσμα:

## 4.4.5

**Λήμμα 1.** Έστω  $f \in C^4[a, b]$  και  $h = \max_j h_j = \max_j (x_{j+1} - x_j)$ . Τότε, αν  $e(x) = f(x) - (I_4 f)(x)$ , έχουμε

$$(4) \quad \max_{1 \leq i \leq N} |e'(x_i)| \leq h^3 \|f^{(4)}\|_{\infty} / 24.$$

Απόδειξη: Εκ κατασκευής της  $I_4 f$  έχουμε  $e'(x_1) = e'(x_N) = 0$ . Στου  $\mathbb{R}^{N-2}$  θεωρούμε τα διανύσματα  $e' = (e'_2, \dots, e'_{N-1})^T$ , όπου  $e'_i = e'(x_i)$ ,  $1 \leq i \leq N$ , και  $r = (r_2, \dots, r_{N-1})^T$  που ορίζεται από την εξίσωση

$$(5) \quad B e' = r,$$

όπου  $B$  ο  $(N-2) \times (N-2)$  πίνακας (3). Χρησιμοποιώντας τις σχέσεις (1) βλέπουμε ότι για  $2 \leq j \leq N-1$ , επειδή  $e'(x_1) = e'(x_N) = 0$

$$\begin{aligned} (6) \quad r_j &= r_j(f) = h_j e'(x_{j-1}) + 2(h_j + h_{j-1}) e'(x_j) + h_{j-1} e'(x_{j+1}) \\ &= h_j f'(x_{j-1}) + 2(h_j + h_{j-1}) f'(x_j) + h_{j-1} f'(x_{j+1}) \\ &\quad - 3[h_{j-1}(f(x_{j+1}) - f(x_j))/h_j + h_j(f(x_j) - f(x_{j-1}))/h_{j-1}], \end{aligned}$$

όπου γράψαμε  $r_j = r_j(f)$  για να συμβολίσουμε το γεγονός ότι για κάθε  $j$ ,  $2 \leq j \leq N-1$ , το  $r_j$  είναι γραμμικό συναρτησιακό στο  $C^4[a, b]$ . Επιπλέον, αν  $p \in P_3$  έχουμε  $r_j(p) = 0$ . Συνεπώς, το θεώρημα του πυρήνα του Peano 4.2.1 μας δίνει:

$$(7) \quad r_j(f) = (1/6) \int_{x_{j-1}}^{x_{j+1}} r_{j,x} [(x-t)_+^3] f^{(4)}(t) dt$$

γιατί η μορφή του συναρτησιακού (6) επιτρέπει την εναλλαγή των πράξεων  $r_{j,x}$  και της ολοκλήρωσης ως προς  $t$ .

θα δείξουμε τώρα, χρησιμοποιώντας τις (6), (7), ότι

$$(8) |r_j(f)| \leq (1/24) \|f^{(4)}\|_{\infty} (h_j h_{j-1}^3 + h_{j-1} h_j^3), \quad 2 \leq j \leq N-1.$$

Μία σειρά απλών πράξεων δίνει, κατ'αρχήν, λόγω της (6), ότι για  $2 \leq j \leq N-1$

$$(9) r_{j,x}[(x-t)_+^3] = \begin{cases} 6(h_j + h_{j-1})(x_j - t)^2 + 3h_{j-1}(x_{j+1} - t)^2 - 3[h_j h_{j-1}^{-1}(x_j - t)^3 \\ + h_{j-1} h_j^{-1}((x_{j+1} - t)^3 - (x_j - t)^3)], & \text{αν } x_{j-1} \leq t \leq x_j, \\ 3h_{j-1}(x_{j+1} - t)^2 - 3[h_{j-1} h_j^{-1}(x_{j+1} - t)^3], & \text{αν } x_j \leq t \leq x_{j+1}. \end{cases}$$

Επίσης μία σειρά απλών πράξεων δίνει άμεσα ότι

$$(10) r_{j,x}[(x-t)_+^3] \geq 0 \quad \forall t \in [x_j, x_{j+1}] \quad \text{και} \quad r_{j,x}[(x-t)_+^3] \leq 0 \quad \forall t \in [x_{j-1}, x_j].$$

(Η απόδειξη της (10) για  $t \in [x_j, x_{j+1}]$  προκύπτει άμεσα από το δεύτερο εκέλος της (9). Για  $t \in [x_{j-1}, x_j]$  το πρώτο εκέλος της (9) δίνει ότι η  $q(t) = r_{j,x}[(x-t)_+^3]$  είναι κυβικό πολυώνυμο στο  $[x_{j-1}, x_j]$  με τις ιδιότητες  $q(x_j) = 0$ ,  $q(x_{j-1}) = 0$ ,  $q'(x_j) = -6h_{j-1}h_j + 9h_{j-1}h_j^{-1}h_j^2 = 3h_{j-1}h_j > 0$ ,  $q'(x_{j-1}) = 0$ .

Συνεπώς το κυβικό πολυώνυμο  $q(t)$  έχει διπλή ρίζα στο σημείο  $x_{j-1}$ , απλή στο  $x_j$  και θετική παράγωγο στο  $x_j$ . Αναγκαστικά λοιπόν  $q(t) \leq 0$  στο διάστημα  $[x_{j-1}, x_j]$ .

Χρησιμοποιώντας τις (10) στην (7) παίρνουμε

$$|r_j(f)| \leq (\|f^{(4)}\|_{\infty}/6) \left\{ \int_{x_{j-1}}^{x_j} -r_{j,x}[(x-t)_+^3] dt + \int_{x_j}^{x_{j+1}} r_{j,x}[(x-t)_+^3] dt \right\}$$

Χρήση τώρα της (9) και απλές ολοκληρώσεις δίνουν την (8).

Το πρόβλημα είναι τώρα να εκτιμήσουμε την λύση  $e' \in \mathbb{R}^{N-2}$  του συστήματος (5) χρησιμοποιώντας τις εκτιμήσεις (8) των συνιστωσών του δευτέρου μέλους. Ας συμβολίσουμε με  $\|\cdot\|_\infty$  και την διανυσματική maximum νόρμα στον  $\mathbb{R}^{N-2}$  καθώς και την παραχόμενη νόρμα πινάκων στον  $\mathbb{R}^{(N-2) \times (N-2)}$ . Πολλαπλασιάζοντας και τα δύο μέλη του συστήματος (5) επί τον  $(N-2) \times (N-2)$  διαγώνιο πίνακα  $D = 2 \text{diag}(h_2+h_1, h_3+h_2, \dots, h_{N-1}+h_{N-2})$  παίρνουμε

$$(11) \quad D^{-1} B e' \equiv (I+M)e' = D^{-1} r.$$

Αμέσως βλέπουμε ότι

$$\|M\|_\infty = \max_i \sum_j |m_{ij}| = \max_{2 \leq i \leq N-1} [(h_i/2(h_i+h_{i-1})) + (h_{i-1}/2(h_i+h_{i-1}))] = 1/2.$$

Συνοψώς από το θεώρημα του Neumann έχουμε ότι ο  $I+M$  είναι αντιστρέψιμος και ότι  $\|(I+M)^{-1}\|_\infty \leq 1/(1-\|M\|_\infty) \leq 2$ . Η (11) λοιπόν δίνει

$$(12) \quad \max_{2 \leq j \leq N-1} |e'(x_j)| = \|e'\|_\infty \leq 2 \|D^{-1} r\|_\infty = \max_{2 \leq j \leq N-1} \{|r_j| / (h_j + h_{j-1})\}.$$

Εξ άλλου από την (8) έχουμε για  $2 \leq j \leq N-1$

$$(13) \quad |r_j| / (h_j + h_{j-1}) \leq (\|f^{(4)}\|_\infty / 24) (h_j h_{j-1}^3 + h_{j-1} h_j^3) / (h_j + h_{j-1})$$

$$\leq \max(h_{j-1}^3, h_j^3) \|f^{(4)}\|_\infty / 24 \leq h^3 \|f^{(4)}\|_\infty / 24.$$

(Οι δύο τελευταίες ανισότητες δεν μπορούν να βελτιωθούν όπως δείχνει η περίπτωση του ομοιόμορφου διαμερισμού). Οι (12) και (13) δίνουν λοιπόν την (4). @

Με το αποτέλεσμα του Λήμματος 1 μπορούμε να αποδείξουμε το βασικό αποτέλεσμα αυτής της παραγράφου:

**ΘΕΩΡΗΜΑ 1.** Αν  $f \in C^4[a, b]$  και  $h = \max_j h_j$ , έχουμε

$$(14) \quad \|f - I_4 f\|_\infty \leq (5/384)h^4 \|f^{(4)}\|_\infty.$$

Απόδειξη: Έστω  $I_H f$  η συνάρτηση παρεμβολής Hermite της  $f$  για τον ίδιο διαμερισμό  $\tau$ . Η τριγωνική ανισότητα μας δίνει

$$(15) \quad \|f - I_4 f\|_\infty \leq \|f - I_H f\|_\infty + \|I_H f - I_4 f\|_\infty.$$

Για τον πρώτο όρο του δεύτερου μέλους έχουμε ήδη την εκτίμηση (4.3.7): απομένει να εκτιμήσουμε την διαφορά  $I_H f - I_4 f$  η οποία, ως τμηματικά κυβική συνάρτηση και στοιχείο του  $C^1[a, b]$ , ανήκει στον χώρο  $H_\tau$  των κυβικών συναρτήσεων Hermite. Εξ ορισμού των  $I_H, I_4$  έχουμε εξ άλλου ότι  $(I_H f - I_4 f)(x_i) = 0$ ,  $1 \leq i \leq N$  και  $(I_H f - I_4 f)'(x_k) = 0$ ,  $k=1, N$ . Συνεπώς ο τύπος (4.3.4) για  $\varphi = I_H f - I_4 f$  δίνει

$$\begin{aligned} (I_H f - I_4 f)(x) &= \sum_{i=2}^{N-1} [(I_H f)'(x_i) - (I_4 f)'(x_i)] S_i(x) \\ &= \sum_{i=2}^{N-1} (f - I_4 f)'(x_i) S(x_i) = \sum_{i=2}^{N-1} e'(x_i) S_i(x), \end{aligned}$$

όπου χρησιμοποίησαμε τον συμβολισμό  $e = f - I_4 f$  του λήμματος 1. Συνεπώς, η (4) δίνει

$$\begin{aligned} (16) \quad \|I_H f - I_4 f\|_\infty &\leq \max_i |e'(x_i)| \max_{a \leq x \leq b} \left( \sum_{i=2}^{N-1} |S_i(x)| \right) \\ &\leq (h^3 \|f^{(4)}\|_\infty / 24) \max_{a \leq x \leq b} \left( \sum_{i=2}^{N-1} |S_i(x)| \right). \end{aligned}$$

Απομένει λοιπόν να εκτιμήσουμε του τελευταίο όρο του δεύτερου μέλους της (16). Σε κάθε υποδιάστημα  $[x_j, x_{j+1}]$ ,  $1 \leq j \leq N-1$  μόνο οι συναρτήσεις  $S_j$  και  $S_{j+1}$  (βλ. (4.3.2)) έχουν μη μηδενικές τιμές. Για  $x \in [x_j, x_{j+1}]$  έχουμε λοιπόν από τις (4.3.2) ότι

$$\begin{aligned} e_j(x) &= |S_j(x)| + |S_{j+1}(x)| = h_j^{-2}(x-x_j)(x_{j+1}-x)^2 + h_j^{-2}(x-x_j)^2(x_{j+1}-x) \\ &= h_j^{-1}(x-x_j)(x_{j+1}-x) \leq h_j/4. \end{aligned}$$

Άρα, αν η συνάρτηση  $\sum_{i=2}^{N-1} |S_i(x)|$  παίρνει το μέγιστό της ε' ένα σημείο  $\bar{x} \in [x_k, x_{k+1}]$  για κάποιο  $k$ ,  $1 \leq k \leq N-1$ , έχουμε

$$(17) \max_{a \leq x \leq b} \sum_{i=2}^{N-1} |S_i(x)| = e_k(\bar{x}) \leq h_k/4 \leq h/4.$$

Οι (16) και (17) δίνουν ευσεπώς

$$(18) \| |f - I_4 f| \|_{\infty} \leq h^4 \| f^{(4)} \|_{\infty} / 96.$$

Οι (4.3.7), (18) και (15) δίνουν τώρα την (14). @

Μπορεί να αποδειχθεί (βλ. Hall και Meyer, J. Approx. Theory, 16(1976), 105-122) ότι η σταθερά  $5/384$  είναι η καλύτερη δυνατή. Στην ίδια εργασία αποδεικνύονται επίσης και φράγματα με άριστες σταθερές για τα εφάλματα των παραχάχων  $(f - I_4 f)^{(k)}$ . Ισχύει το εξής αποτέλεσμα:

**ΘΕΩΡΗΜΑ 2** (Hall και Meyer). Έστω  $f \in C^4[a, b]$ ,  $h = \max_j h_j$ . Τότε έχουμε

$$(19) \quad \|(f-I_4 f)^{(k)}\|_{\infty} \leq c_k h^{4-k} \|f^{(4)}\|_{\infty}, \quad k=0,1,2,3,$$

όπου  $c_0=5/384$ ,  $c_1=1/24$ ,  $c_2=3/8$  και  $c_3=(\lambda+\lambda^{-1})/2$  όπου  $\lambda=h/(\min_j h_j)$ .

Επιπλέον οι σταθερές  $c_0, c_1$  είναι οι καλύτερες δυνατές με την έννοια ότι

$$c_k = \sup_{f \in C^4[a,b], \tau} \{ \|(f-I_4 f)^{(k)}\| / h^{4-k} \|f^{(4)}\|_{\infty} \}, \quad k=0,1. \textcircled{a}$$

Προχωρούμε τώρα στην διερεύνηση των άλλων ερωτημάτων που θέσαμε στην εισαγωγή αυτής της παραγράφου, δηλ. στο να βρούμε την διάσταση του χώρου  $S_{\tau}^4$  και να κατασκευάσουμε μία αποτελεσματική για τις εφαρμογές βάση του. Η διάσταση που περιμένουμε να έχει ο χώρος  $S_{\tau}^4$  είναι  $N+2$ . Πράγματι ο προσδιορισμός ενός στοιχείου του  $S_{\tau}^4$  (δηλ. μιας κυβικής spline) απαιτεί τον προσδιορισμό  $4(N-1)$  σταθερών, δηλ. των τεσσάρων συντελεστών ενός κυβικού πολυωνόμου σε κάθε διάστημα  $[x_j, x_{j+1}]$ ,  $1 \leq j \leq N-1$ . Επιβάλλοντας τον περιορισμό  $S_{\tau}^4 \subset C^2[a,b]$ , δηλ. επιβάλλοντας  $C^2$  ομαλότητα στους εσωτερικούς κόμβους  $x_j$ ,  $2 \leq j \leq N-1$  παίρνουμε  $3(N-2)$  συνθήκες. Άρα οι "βαθμοί ελευθερίας" (συνολικός αριθμός ελευθέρων παραμέτρων = πλήθος "ευντεταχμένων" μίας κυβικής spline = διάσταση του χώρου  $S_{\tau}^4$ ) είναι  $4(N-1) - 3(N-2) = N+2$ . Θα αποδείξουμε αυτό το αποτέλεσμα κατασκευάζοντας ευχρόνως και μια χρήσιμη βάση του  $S_{\tau}^4$ , τις B-splines.

Μια επιθυμητή ιδιότητα για τις συναρτήσεις βάσης είναι να έχουν μικρό φορέα. Είναι εύκολο να δει κανείς ότι δεν υπάρχει άλλη κυβική spline (δηλ.  $C^2$ ) με φορέα  $[x_{j-1}, x_{j+1}]$ , δηλ. δύο διαστήματα, εκτός από την τετριμμένη μηδενική συνάρτηση· πράγματι πρέπει να επιβάλλουμε 9 συνθήκες (τρεις σε κάθε κόμβο  $x_{j-1}, x_j, x_{j+1}$ ) για να προσδιορίσουμε 8 σταθερές. Δεν υπάρχει επίσης μη τετριμμένη κυβική spline με φορέα τριών διαστημάτων. Ο μικρότερος δυνατός φορέας αποτελείται από 4 διαστήματα, είναι δηλ. για  $3 \leq j \leq N-2$  της μορφής  $[x_{j-2}, x_{j+2}]$ .



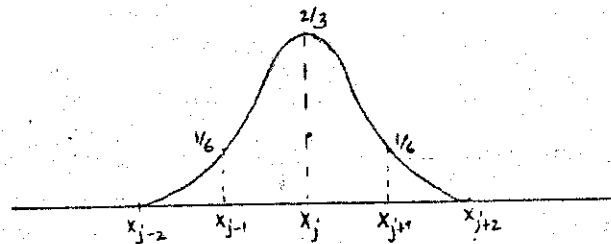
Στην περίπτωση αυτή πρέπει να προσδιορίσουμε 16 παραμέτρους από 15 εξισώσεις που επιτρέπουν το υπολογισμό ενός τμηματικά κυβικού πολυωνύμου στον χώρο  $C^2[a,b]$  με φορέα  $[x_{j-2}, x_{j+2}]$ . όλα τα πολλαπλάσια του επί σταθερά θα έχουν προφανώς τις ίδιες ιδιότητες.

Η συνάρτηση  $s_j(x)$  που προκύπτει λέγεται B-spline (με κέντρο τον κόμβο  $x_j$ ). Χρησιμοποιώντας διαιρεμένες διαφορές (βλ. π.χ. το βιβλίο του De Boor [4.2, κεφ. 9]) μπορούμε να την κατασκευάσουμε αρκετά εύκολα χωρίς να λύσουμε το γραμμικό σύστημα. Η συνάρτηση που προκύπτει είναι, για  $3 \leq j \leq N-2$

$$(20) \quad s_j(x) = (x_{j+2} - x_{j-2}) \sum_{k=-2}^2 [(x_{j+k} - x)_+^3 \prod_{\substack{s=-2 \\ s \neq k}}^2 (x_{j+k} - x_{j+s})^{-1}].$$

Στην περίπτωση του ομοιόμορφου διαμερισμού με  $h = x_{j+1} - x_j = \text{σταθ}$  η B-spline  $s_j$  με κέντρο  $x_j$  δίνεται για  $3 \leq j \leq N-2$  από

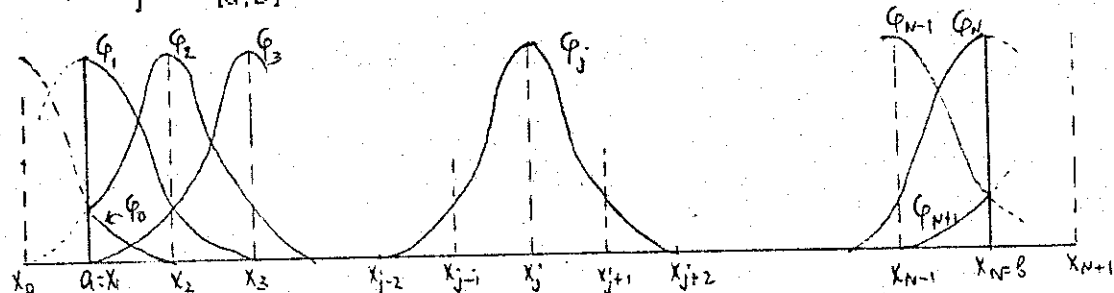
$$(21) \quad s_j(x) = \begin{cases} (x - x_{j-2})^3 / 6h^3, & x_{j-2} \leq x \leq x_{j-1} \\ (2/3) - (x - x_j)^2 / h^2 - (x - x_j)^3 / 2h^3, & x_{j-1} \leq x \leq x_j \\ (1/6) + (x_{j+1} - x) / 2h + (x_{j+1} - x)^2 / 2h^2 - (x_{j+1} - x)^3 / 2h^3, & x_j \leq x \leq x_{j+1} \\ (x_{j+2} - x)^3 / 6h^3, & x_{j+1} \leq x \leq x_{j+2} \\ 0, & x \notin [x_{j-2}, x_{j+2}] \end{cases}$$

$s_j(x)$ :

Μ' αυτόν τόν τρόπο κατασκευάζουμε  $N-4$  B-splines  $s_j$ ,  $3 \leq j \leq N-2$  με κέντρα ετά αντίστοιχα σημεία  $x_j$ . Η βάση του χώρου  $S_C^4$  προσδιορίζεται τώρα ως εξής: θεωρούμε δύο επιπλέον κόμβους  $x_0 < a$  και  $x_{N+1} > b$  - στην περίπτωση ομοιόμορφου διαμερισμού  $x_0 = a-h$ ,  $x_{N+1} = b+h$  - και ορίζουμε  $N+2$  συναρτήσεις  $\varphi_j$  ως

$$(22) \varphi_j(x) = s_j(x)|_{[a,b]}, \quad 0 \leq j \leq N+1,$$

όπου με  $s_j(x)|_{[a,b]}$  συμβολίζουμε τον περιορισμό της  $s_j$  ετο  $[a,b]$ .



Οι συναρτήσεις  $\varphi_j$ ,  $0 \leq j \leq N+1$ , αποτελούν βάση του  $S_C^4$ , ο οποίος ευενώς έχει διάσταση  $N+2$ . Ας αποδείξουμε το αποτέλεσμα αυτό, για να απλοποιήσουμε τα πράγματα, στην περίπτωση του ομοιόμορφου διαμερισμού με  $h = x_{j+1} - x_j$ ,  $0 \leq j \leq N$ . Κατ' αρχήν οι  $\varphi_j$ ,  $0 \leq j \leq N+1$  είναι

γραμμικά ανεξάρτητες: Ας υποθέσουμε ότι  $\sum_{j=0}^{N+1} c_j \varphi_j(x) = 0 \quad \forall x \in [a,b]$ .

Τότε  $\sum_{j=0}^{N+1} c_j \varphi_j(x_i) = 0$ ,  $1 \leq i \leq N$  και  $\sum_{j=0}^{N+1} c_j \varphi_j'(x_k) = 0$ ,  $k=1, N$ . Από την

(21) βλέπουμε ότι οι παραπάνω εξισώσεις οδηγούν στις ευθείες

$$(23) \quad A c = 0, \quad c = [c_1, \dots, c_N]^T, \quad c_0 = c_2, \quad c_{N+1} = c_{N-1},$$

όπου ο  $N \times N$  τριδιαγώνιος πίνακας  $A$  δίνεται από

$$(24) \quad A = \begin{pmatrix} 4 & 2 & & & 0 \\ 1 & 4 & 1 & & \\ & \ddots & \ddots & \ddots & \\ 0 & 1 & 4 & 1 & \\ & & 2 & 4 & \end{pmatrix}$$

και είναι αντιστρέψιμος επειδή έχει αυστηρά κυριαρχική διαγώνιο. Από τις (23) συμπεραίνουμε ότι  $c_j = 0$ ,  $0 \leq j \leq N+1$ , δηλ. ότι οι  $\varphi_j$ ,  $0 \leq j \leq N+1$  είναι γραμμικά ανεξάρτητες. Για να δείξουμε ότι παράχουν του  $S_v^4$ , θεωρούμε οποιοδήποτε διάστημα  $I_j = [x_j, x_{j+1}]$  για κάποιο  $1 \leq j \leq N-2$  και τον τετραδιάστατο χώρο  $P_3(I_j)$ . Επειδή τα κυβικά πολυώνυμα στο  $I_j$   $\varphi_{j-1}, \varphi_j, \varphi_{j+1}, \varphi_{j+2}$  (εκείνες δηλ. οι  $\varphi_j$  των οποίων ο φορέας περιέχει το  $I_j$ ) είναι γραμμικά ανεξάρτητα, θα παράχουν του  $P_3(I_j)$ . Εστω  $\varphi(x)$  τυχόν στοιχείο του  $S_v^4$ . Επειδή  $\varphi|_{I_j} \in P_3(I_j)$  θα έχουμε για  $x \in I_j$

$$(25) \quad \varphi(x) = c_{j-1}^{(j)} \varphi_{j-1}(x) + c_j^{(j)} \varphi_j(x) - c_{j+1}^{(j)} \varphi_{j+1}(x) + c_{j+2}^{(j)} \varphi_{j+2}(x),$$

όπου τονίζουμε με τον άνω δείκτη  $(j)$  των  $c_{j+k}^{(j)}$ ,  $k = -1, 0, 1, 2$  το ότι η τετράδα των συντελεστών αυτών θα εξαρτάται κατ' αρχήν από το διάστημα  $I_j$ . Αν δείξουμε ότι είναι ανεξάρτητοι του  $j$ , δηλ. ότι αν π.χ. για  $x \in I_{j+1}$

$$(26) \quad \varphi(x) = c_j^{(j+1)} \varphi_j(x) + c_{j+1}^{(j+1)} \varphi_{j+1}(x) + c_{j+2}^{(j+1)} \varphi_{j+2}(x) + c_{j+3}^{(j+1)} \varphi_{j+3}(x),$$

ευνεπάγεται  $c_{j+k}^{(j)} = c_{j+k}^{(j+1)}$ ,  $k=0,1,2$ , τότε είναι φανερό ότι κάθε  $\varphi \in S_T^4$  μπορεί να εκφρασθεί για  $x \in [a, b]$  ως γραμμικός συνδυασμός των  $\varphi_j$ ,  $0 \leq j \leq N+1$ . Πράγματι, οι σχέσεις (25) και (26) και η  $C^2$  συνέχεια της  $\varphi$  στο  $x=x_{j+1}$  δίνουν, λόγω της (21),

$$\varphi(x_{j+1}^-) = \varphi(x_{j+1}^+); c_j^{(j)}/6 + 2c_{j+1}^{(j)}/3 + c_{j+2}^{(j)}/6 = c_j^{(j+1)}/6 + 2c_{j+1}^{(j+1)}/3 + c_{j+2}^{(j+1)}/6$$

$$\varphi'(x_{j+1}^-) = \varphi'(x_{j+1}^+); -c_j^{(j)}/2h + c_{j+2}^{(j)}/2h = -c_j^{(j+1)}/2h + c_{j+2}^{(j+1)}/2h$$

$$\varphi''(x_{j+1}^-) = \varphi''(x_{j+1}^+); c_j^{(j)}/h^2 - 2c_{j+1}^{(j)}/h^2 + c_{j+2}^{(j)}/h^2 = c_j^{(j+1)}/h^2$$

$$-2c_{j+1}^{(j+1)}/h^2 + c_{j+2}^{(j+1)}/h^2.$$

Συνοψώς οι αριθμοί  $d_k = c_{j+k}^{(j)} - c_{j+k}^{(j+1)}$ ,  $k=0,1,2$ , ικανοποιούν το  $3 \times 3$  ομογενές σύστημα

$$\begin{pmatrix} 1 & 4 & 1 \\ -1 & 0 & 1 \\ 1 & -2 & 1 \end{pmatrix} \begin{pmatrix} d_0 \\ d_1 \\ d_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

που έχει την μοναδική λύση  $d_k = 0$ ,  $k=0,1,2$ .

Η χρήση της βάσης  $\{\varphi_j\}$ ,  $0 \leq j \leq N-1$  του  $S_T^4$  μας επιτρέπει να λύσουμε αποτελεσματικά, δηλ. με λύση απλών γραμμικών συστημάτων με καλούς δείκτες κατάστασης, προβλήματα όπως π.χ. του υπολογισμού της  $L^2$ -προβολής μιάς συνάρτησης  $f$  στον  $S_T^4$  (βλ. παρ. 4.2) της αριθμητικής λύσης διαφορικών εξισώσεων με μεθόδους τύπου Galerkin κ.ά. Σαν παράδειγμα αναφέρουμε το πρόβλημα της παρεμβολής στον  $S_T^4$ : Γράφοντας για  $f \in C^1[a, b]$

$$(27) \quad (I_4 f)(x) = \sum_{j=0}^{N+1} c_j \varphi_j(x),$$

και χρησιμοποιώντας τις ευθείες παρεμβολές  $(I_4 f)(x_i) = f(x_i)$ ,  $1 \leq i \leq N$ ,  
 $(I_4 f)'(x_k) = f'(x_k)$ ,  $k=1, N$  οδηγούμαστε στις εξισώσεις (π.χ. για  
 ομοιόμορφο διαμερισμό):

$$\sum_{j=0}^{N+1} c_j \varphi_j(x_1) = f(x_1), \quad 1 \leq i \leq N, \quad c_0 \varphi'_0(x_1) + c_2 \varphi'_2(x_1) = f'(x_1),$$

$$c_{N-1} \varphi'_{N-1}(x_N) + c_{N+1} \varphi'_{N+1}(x_N) = f'(x_N)$$

που δίνουν το  $N \times N$  γραμμικό σύστημα

$$(28) \quad A c = b,$$

όπου  $c = [c_1, \dots, c_N]^T$ ,  $A$  ο τριδιαγώνιος πίνακας (24) και  
 $b = [b_1, \dots, b_N]^T$ , όπου  $b_1 = 6f(x_1) + 2hf'(x_1)$ ,  $b_i = 6f(x_i)$ ,  $2 \leq i \leq N-1$ ,  
 $b_N = 6f(x_N) - 2hf'(x_N)$ . Το (28) προφανώς έχει μοναδική λύση. Τέλος  
 προσδιορίζοντας τα  $c_0, c_{N+1}$  από τις εξισώσεις  $c_0 = c_2 - 2hf'(x_1)$ ,  
 $c_{N+1} = c_{N-1} + 2hf'(x_N)$  υπολογίζουμε την συνάρτηση παρεμβολής  $I_4 f$  από την  
 (27). Η λύση του (28) βρίσκεται με  $5N$  περίπου πράξεις  
 χρησιμοποιώντας τον (ευσταθή στην περίπτωση μας) τριδιαγώνιο  
 αλγόριθμο.

#### Παρατηρήσεις

1. Η πρώτη απόδειξη ότι  $\|f - I_4 f\|_{\infty} = O(h^4)$  για αρκετά ομαλή  $f$   
 οφείλεται στους Birkhoff και De Boor (J. Math. Mech., 13(1964),  
 827-836) για διαμερισμούς για τους οποίους ο λόγος  $\lambda = \max_j h_j / \min_j h_j$   
 μένει φραγμένος. Η απόδειξη του θεωρήματος 1 χωρίς περιορισμό στον  
 διαμερισμό οφείλεται στον Hall (J. Approx. Theory 1(1968), 209-218)  
 ενώ του θεωρήματος 2 στους Hall και Meyer (op. cit.).

2. Μία σημαντική ιδιότητα της κυβικής spline παρεμβολής  $I_4 f$  (βλ. θεο. 8γ) είναι ότι ελαχιστοποιεί το ολοκλήρωμα  $\int_a^b (g''(t))^2 dt$  μεταξύ όλων των  $C^2$ -ευναρτήσεων  $g$  που ικανοποιούν τις εκθέσεις παρεμβολής  $g(x_i) = f(x_i)$ ,  $1 \leq i \leq N$ . Η ιδιότητα αυτή έχει κάποια σημασία στην μηχανική (θεωρία εύκαμπτων λεπτών δοκών) και ε' αυτήν οφείλεται η ονομασία "spline", που είναι το σχήμα που παίρνει μία εύκαμπτη λεπτή οδηγός καμπύλη όταν διέρχεται από δεδομένα σημεία  $(x_i, y_i)$ ,  $1 \leq i \leq N$  στο επίπεδο. Τέτοιες καμπύλες χρησιμοποιούνται για την σχεδίαση πλοίων, αυτοκινήτων κ.λ.π.

3. Θα δούμε στην Άσκηση 10 ανάλογα αποτελέσματα του θεωρήματος 2 αλλά ως προς την  $L^2$ -νόρμα  $\| \cdot \|$ . Προφανώς υπάρχουν και τα αναμενόμενα ανάλογα αποτελέσματα ε' όλης τις ευνήθεις νόρμες για φράγματα εφαλμάτων μικρότερης τάξης ακρίβειας αν η  $f$  δεν είναι  $C^4$ .  
βλ. π.χ. θεο. 2, 6, 9.

4. Η κυβική spline παρεμβολής  $f_h = I_4 f$  που μελετήσαμε είναι η λεγόμενη πλήρης συναρτησιακή παρεμβολής της  $f$  στον χώρο  $S_c^4$  και, όπως είδαμε, αντιστοιχεί σε ευνοριακές ευνήθεις δεδομένης κλίσεως στα άκρα, δηλ. στις ευνήθεις  $(I_4 f)'(a) = f'(a)$ ,  $(I_4 f)'(b) = f'(b)$ . Είναι βέβαια δυνατόν να επιβάλουμε, πέρα από τις ευνήθεις παρεμβολής στις τιμές της  $f$  στα σημεία  $x_i$ ,  $1 \leq i \leq N$ , και άλλα ζευγάρια βοηθητικών επιπλέον ευνήθων, π.χ. άλλου τύπου ευνοριακές ευνήθεις. Αν π.χ. οι τιμές  $f''(a)$ ,  $f''(b)$  είναι γνωστές, μπορούμε να επιβάλουμε τις ευνήθεις  $f''_h(a) = f''(a)$ ,  $f''_h(b) = f''(b)$  που οδηγούν στις εξισώσεις:

$$2s'(x_1) + s'(x_2) = 3h_1^{-1}(f(x_2) - f(x_1)) - h_1 f''(x_1)/2$$

και

$$s'(x_{N-1}) + 2s'(x_N) = 3h_{N-1}^{-1}(f(x_N) - f(x_{N-1})) + h_{N-1} f''(x_N)/2,$$

οι οποίες μαζί με τις  $N-2$  εξισώσεις (1) επιτρέπουν τον υπολογισμό των κλίσεων  $s'(x_i)$ ,  $1 \leq i \leq N$  της spline  $f_h$ .

Αν τώρα είναι γνωστές μόνο οι τιμές  $f(x_i)$ ,  $1 \leq i \leq N$ , μπορούμε να προεχθίσουμε π.χ. τις τιμές  $f'(x_1), f'(x_N)$  που απαιτούνται για τον προσδιορισμό του  $I_4 f$  π.χ. αντικαθιστώντας την  $f'(x_1)$  με  $p'(x_1)$  όπου  $p$  το κυβικό πολυώνυμο παρεμβολής Lagrange της  $f$  για τα σημεία  $x_1, x_2, x_3, x_4$  κ.ο.κ..

5. Αξιοσημείωτη είναι η ιδιότητα "υπερέυγκλισης" (που ισχύει για ομοιόμορφο διαμερισμό και  $f \in C^5[a, b]$ ):

$$(29) \quad \max_{1 \leq i \leq N} |(f - I_4 f)'(x_i)| \leq h^4 \|f^{(5)}\|_{\infty} / 60$$

(βλ. Αεκ. 3) και που εξηγεί την μεγάλη ακρίβεια των τιμών και της παραχώχου  $(I_4 f)'$  στους κόμβους για ομοιόμορφο διαμερισμό.

6. Παραπέμπουμε τον αναγνώστη π.χ. στο βιβλίο [4.2] του De Boor για την μελέτη τμηματικά πολυωνυμικών συναρτήσεων ("splines") βαθμού 2 καθώς και οποιουδήποτε βαθμού  $k > 3$ .

#### Αεκήσεις 4.4

1. Μία άλλη κατασκευή της συναρτήσεως  $I_4 f$  είναι η εξής:

Υπολογίστε το κυβικό πολυώνυμο  $s_i(x) = (I_4 f)(x) |_{[x_i, x_{i+1}]}$ ,  $1 \leq i \leq N-1$ ,

συναρτήσεως των τιμών του στα σημεία  $x_i$  και  $x_{i+1}$  και των τιμών  $m_i, m_{i+1}$  των δευτέρων παραχώχων του στα  $x_i, x_{i+1}$ . Οι σταθερές  $m_i, 1 \leq i \leq N$  προσδιορίζονται κατόπι από τις συνθήκες  $s'_{i-1}(m_i) = s'_i(x_i)$ ,  $2 \leq i \leq N-1$  ( $C^1$  στους εσωτερικούς κόμβους) και  $s'_1(x_1) = f'(x_1)$ ,  $s'_{N-1}(x_N) = f'(x_N)$ .

Βρείτε το γραμμικό σύστημα που προκύπτει για τα  $m_i$  και αποδείξτε ότι έχει μοναδική λύση.

2. (Σφάλμα της  $I_4 f$  όταν  $f \in C^1[a, b]$ .)

(α) Χρησιμοποιώντας τις τεχνικές της απόδειξης του Λήμματος 4.2.2, δείξτε ότι η λύση  $s'(x_j)$ ,  $2 \leq j \leq N-1$  του ευστήματος με πίνακα  $B$ , (3), στην απόδειξη της πρότασης 1 ικανοποιεί, για δεδομένα  $s'(x_1)$ ,  $s'(x_N)$ ,

$$\max_j |s'(x_j)| \leq 3 \max(|s'(x_1)|, \max_{2 \leq j \leq N-1} |b_j|, |s'(x_N)|),$$

όπου  $b_j$ ,  $2 \leq j \leq N-1$  είναι το δεύτερο μέλος του ευστήματος, δηλ. το δεύτερο μέλος των εξισώσεων (1).

(β). Συμπεράνετε ότι  $\|(I_4 f)'\|_{\infty} \leq 4 \max_{0 \leq i \leq N} |h_i^{-1}(f(x_{i+1}) - f(x_i))|$  όπου ορίζουμε  $h_0 = h_N = 0$ . (Υπόδειξη: Δείξτε ότι

$$(I_4 f)'((x_i + x_{i+1})/2) = 3h_i^{-1}(f(x_{i+1}) - f(x_i))/2 - (s'(x_i) + s'(x_{i+1}))/4$$

και κατόπιν χρησιμοποιείτε την ανισότητα

$$\max_{a \leq x \leq b} |p(x)| \leq (5/4) \max_{a \leq x \leq b} (|p(a)|, |p((a+b)/2)|, |p(b)|),$$

που ισχύει για κάθε  $p \in P_3$ ).

(γ) Συμπεράνετε από το (β) και την άσκηση 4.2.5 ότι

$$\|f' - (I_4 f)'\|_{\infty} \leq 5 \inf_{\varphi \in H_4} \|f' - \varphi\|_{\infty}$$

(δ) Τέλος αποδείξτε ότι

$$\|f' - (I_4 f)'\|_{\infty} \leq 5h \|f'\|_{\infty} / 2.$$

3. Υποθέστε ότι έχουμε ομοιόμορφο διαμερισμό και  $f \in C^5[a, b]$ . Τότε, από την (δ) έχουμε



4.4.19

$$r_j(f)/3 = h(f'(x_{j-1}) + 2f'(x_j) + f'(x_{j+1}))/3 - \int_{x_{j-1}}^{x_{j+1}} f'(x) dx, \quad 2 \leq j \leq N-1$$

(α) Χρησιμοποιώντας ένα γινόμενο από την Εισαγωγή στην Αριθμητική Ανάλυση αποτέλεσμα για το σφάλμα του κανόνα του Simpson για αριθμητική ολοκλήρωση δείξτε ότι

$$|r_j(f)| \leq h^5 \|f^{(5)}\|_{\infty} / 30, \quad 2 \leq j \leq N-1.$$

(β) Συμπεράνετε, όπως στην απόδειξη του Λήμματος 1, ότι ισχύει η (29).

6. Μιμηθείτε την απόδειξη του Θεωρήματος 1 και χρησιμοποιώντας το θεώρημα του Peano δείξτε ότι

$$\|f - I_4 f\|_{\infty} \leq 2h^2 \|f''\|_{\infty} / 3 \quad \text{αν } f \in C^2[a, b].$$

7. Με χρήση Peano δείξτε επίσης τα άξι άριστα ως προς τις σταθερές φράγματα (Hall, 1968) για ομοιόμορφο διαμερισμό:

$$\|(f - I_4 f)''\|_{\infty} \leq (\sqrt{3/216} + 1/24) h^3 \|f^{(4)}\|_{\infty}$$

$$\|(f - I_4 f)'''\|_{\infty} \leq 5h^2 \|f^{(4)}\|_{\infty} / 12$$

$$\|(f - I_4 f)''''\|_{\infty} \leq h \|f^{(4)}\|_{\infty}$$

8. (α) Αν  $f \in PC^2[a, b]$  τότε  $((I_4 f - f)'', \varphi) = 0, \quad \forall \varphi \in S_c^4$

(β) Αν  $f \in PC^2[a, b]$ , τότε  $\|(I_4 f)''\|^2 + \|(I_4 f - f)''\|^2 = \|f''\|^2$

(γ) Αποδείξτε τον ισχυρισμό της Παρατήρησης 2.

(δ) Αν  $f \in PC^4[a, b]$ , τότε  $\|(f - I_4 f)^{(k)}\|^2 = (f - I_4 f, f^{(4)})$

9. Χωρίς χρήση του θεωρήματος του Peano αλλά με χρήση των "ολοκληρωτικών ταυτοτήτων" της Άσκησης 8 και της αμειωότητας του Poincaré (Άσκ. 4.2.9) δείξτε ότι αν  $f \in PC^2[a, b]$  τότε

$$\|(f - I_4 f)^{(k)}\| \leq C_k h^{2-k} \|f''\|, \quad k=0, 1, 2 \text{ με } C_0=2, C_1=2, C_2=1.$$

10. Με παρόμοιο τρόπο δείξτε, αν  $f \in PC^4[a, b]$  ότι

$$\|(f - I_4 f)^{(k)}\| \leq C_k h^{4-k} \|f^{(4)}\|, \quad k=0, 1, 2, \text{ με } C_0=4, C_1=4, C_2=2.$$