

$$(3.33) (1+\|A\|)^{-1} \leq \|(1-A)^{-1}\| \leq (1-\|A\|)^{-1}.$$

ΑΠΟΔΕΙΞΗ Από την εξίσωση $\|A\| < 1$ και το γεγονός ότι $\rho(A) \leq \|A\|$ έπεται ότι για κάθε ιδιοτιμή λ_i του A ισχύει ότι $|\lambda_i| < 1$. Μ' άλλα λόγια για κάθε ιδιοτιμή έχουμε ότι $\lambda_i \neq 1$. Συνεπώς ο πίνακας $I-A$ είναι αντιστρέψιμος γιατί, αν δεν ήταν, θα υπήρχε τουλάχιστον μια μη μηδενική λύση $u \neq 0$ του συστήματος $(I-A)u = 0 \Leftrightarrow Au = u$ με $u \neq 0$, δηλ. ο A θα έχει μια ιδιοτιμή ίση με 1, άτοπο. Επίσης, επειδή η $\|\cdot\|$ είναι φυσική νόρμα, έχουμε ότι $\|I\| = 1$ (βλ. κλην 3.8). Άρα επειδή υπάρχει ο $(I-A)^{-1}$ έχουμε:

$$\begin{aligned} 1 = \|I\| &= \|(I-A)^{-1} (I-A)\| \leq \|(I-A)^{-1}\| \|I-A\| \leq \\ &\leq \|(I-A)^{-1}\| (\|I\| + \|A\|) = \|(I-A)^{-1}\| (1+\|A\|). \end{aligned}$$

Άρα $(1+\|A\|)^{-1} \leq \|(I-A)^{-1}\|$. (αριστερή ανισότητα της (3.33)).

Αν' την άλλη μεριά, επειδή $(I-A)^{-1} (I-A) = I$ έχουμε $(I-A)^{-1} - A(I-A)^{-1} = I$ δηλ. $(I-A)^{-1} = I + A(I-A)^{-1}$. Άρα $\|(I-A)^{-1}\| = \|I + A(I-A)^{-1}\| \leq \|I\| + \|A(I-A)^{-1}\| \leq 1 + \|A\| \|(I-A)^{-1}\|$. Συνεπώς $(1-\|A\|) \|(I-A)^{-1}\| \leq 1$ και επειδή $\|A\| < 1$ έπεται ότι $\|(I-A)^{-1}\| \leq (1-\|A\|)^{-1}$ (δεξιά ανισότητα της (3.33)).

4. ΔΕΙΚΤΗΣ ΚΑΤΑΣΤΑΣΗΣ ΠΙΝΑΚΑ

4.1 Δείκτης κατάστασης και ευστάθεια συστημάτων

Με τα εργαλεία του προηγούμενου κεφαλαίου ξαναχρυσίζουμε στο πρόβλημα της ευστάθειας γραμμικών συστημάτων και αλγορίθμου για τη λύση τους που θίξαμε στην παρ. 2.5. Εκεί είχαμε ορίσει άτυπα την ευστάθεια (ενός προβλήματος ή ενός αλγορίθμου) σαν την ιδιότητα βάσει της οποίας "μικρά" εφάλματα (στα δεδομένα, στους υπολογισμούς) προκαλούν "μικρή" μεταβολή στο τελικό αποτέλεσμα. Θα δούμε ότι καθοριστικό ρόλο στην ευστάθεια τόσο του συστήματος $Ax=b$ όσο και στην ευστάθεια του αλγορίθμου της απαλοιφής του Gauss για την λύση του παίζει ο λεγόμενος δείκτης κατάστασης του πίνακα.

Αρχίζουμε από το εξής απλό παράδειγμα. Έστω x η λύση του συστήματος

$$(4.1) Ax = b,$$

όπου A είναι ένας $n \times n$ αντιστρέψιμος πίνακας και b ένα $n \times 1$ διάνυσμα. Αν μεταβάλλουμε σε $b+\delta b$ το δεύτερο μέλος της (4.1), πόσο θα μεταβληθεί η λύση; Έστω $x+\delta x$ η λύση του νέου συστήματος

$$(4.2) A(x+\delta x) = b+\delta b.$$

Αφαιρώντας κατά μέλη έχουμε ότι $A\delta x = \delta b$ ή

$$(4.3) \delta x = A^{-1}\delta b.$$

Έστω τώρα $\|\cdot\|$ μια διανυσματική νόρμα και η αντίστοιχη φυσική νόρμα πινάκων που παράχεται από αυτήν. Η (4.3) δίνει ότι

$$(4.4) \|\delta x\| = \|A^{-1}\delta b\| \leq \|A^{-1}\| \|\delta b\|.$$

Αν υποθέσουμε ότι $b \neq 0$ τότε $x \neq 0$ και συνεπώς $\|\delta x\|/\|x\| \leq \|A^{-1}\| \|\delta b\|/\|x\|$. Αλλά, από την (4.1), $\|A\| \|x\| \geq \|b\|$. Συνεπώς

$$(4.5) \|\delta x\|/\|x\| \leq \|A\| \|A^{-1}\| (\|\delta b\|/\|b\|).$$

Η ποσότητα $\|b\|/\|b\|$ είναι η εχρητική μεταβολή του δευτέρου μέλους της (4.1) και η $\|b\|/\|x\|$ η αντίστοιχη εχρητική μεταβολή της λύσης (ή εχρητικό "εφάλμα", αν θεωρήσουμε π.χ. ότι το b παριστάνει ένα μικρό εφάλμα στην παράσταση των δεδομένων του δευτέρου μέλους της (4.1)). Βλέπουμε λοιπόν ότι η ποσότητα $\|A\| \|A^{-1}\|$ είναι ένας συντελεστής που προσδιορίζει πόσο μεγάλο μπορεί να γίνει το εχρητικό εφάλμα $\|b\|/\|x\|$, όταν το εχρητικό εφάλμα του δευτέρου μέλους είναι $\|b\|/\|b\|$.

Ο αριθμός

$$(4.6) \kappa(A) = \|A\| \|A^{-1}\|,$$

που εξαρτάται φυσικά από τον A και την χρησιμοποιούμενη νόρμα λέγεται δείκτης κατάστασης του πίνακα A ως προς την νόρμα $\|\cdot\|$. Αν ο $\kappa(A)$ είναι πολύ μεγάλος τότε περιμένουμε ότι μια μικρή μεταβολή του b μπορεί να προκαλέσει μεγάλη μεταβολή στη λύση.

ΑΣΚΗΣΗ 4.1 Να δείχθει ότι ο δείκτης κατάστασης ως προς την νόρμα $\|\cdot\|_1$ του συστήματος (2.28) είναι περίπου 2.6×10^6 . @

Ο δείκτης κατάστασης ενός πίνακα ως προς οποιαδήποτε φυσική νόρμα $\|\cdot\|$ είναι πάντα μεγαλύτερος ή ίσος της μονάδας, γιατί $1 = \|I\| = \|AA^{-1}\| \leq \|A\| \|A^{-1}\| = \kappa(A)$. Λέμε ότι ο πίνακας έχει κακή κατάσταση αν $\kappa(A) \gg 1$ και καλή κατάσταση αν ο $\kappa(A)$ είναι μικρός, π.χ. αν $1 \leq \kappa(A) \leq 100$. Εκτός από μέτρο της "ευαισθησίας" της λύσης x σε μεταβολές του b , μπορούμε να αποδείξουμε ότι ο $\kappa(A)$ είναι και (αντίστροφο) μέτρο του πόσο κοντά βρίσκεται ο A στο σύνολο των μη αντιστρέψιμων πινάκων. Πράγματι έστω B ένας οποιοσδήποτε μη αντιστρέψιμος $n \times n$ πίνακας. Τότε υπάρχει $x \neq 0$ τέτοιο ώστε $Bx=0$. Συνεπώς $\|A-B\| \|x\| \geq \|(A-B)x\| = \|Ax-Bx\| = \|Ax\|$. Επίσης $\|x\| = \|A^{-1}Ax\| \leq \|A^{-1}\| \|Ax\| \Rightarrow \|Ax\| \geq \|x\|/\|A^{-1}\|$. Συνεπώς έχουμε, επειδή $x \neq 0$, ότι $\|A^{-1}\|^{-1} \leq \|A-B\|$, για οποιοδήποτε μη αντιστρέψιμο πίνακα B , δηλ. ότι

$$(\kappa(A))^{-1} \leq \inf \{ \|A-B\|/\|A\|, B \text{ μη αντιστρέψιμος} \}.$$

Μπορεί να αποδειχθεί ακριβέστερα ότι η παραπάνω ανισότητα \leq είναι στην πραγματικότητα ισότητα, δηλ. ότι ο αντίστροφος $(\kappa(A))^{-1}$ του δείκτη κατάστασης του A μετράει την (εχρητική ως προς το μέγεθος του $\|A\|$) απόσταση του πίνακα A από το σύνολο των $n \times n$ μη αντιστρέψιμων πινάκων.

ΑΣΚΗΣΗ 4.2 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). @

ΑΣΚΗΣΗ 4.3 Έστω A αντιστρέψιμος άνω ή κάτω τριγωνικός πίνακας. Να δείχθει ότι

$$\kappa_\infty(A) \geq \|A\|_\infty / \min_i |a_{ii}|. @$$

ΑΣΚΗΣΗ 4.4 Αν ο A είναι αντιστρέψιμος και ο πίνακας B ικανοποιεί την σχέση

$$\|A-B\| < 1/\|A^{-1}\|,$$

τότε ο B είναι αντιστρέψιμος. @

Επιστρέφουμε τώρα στην (4.5). Θα έλεγε κανείς ότι ο δείκτης κατάστασης $\kappa(A)$ είναι μόνο ένα άνω όριο του λόγου $(\|b\|/\|x\|)/(\|b\|/\|b\|)$. Μπορεί όμως να δείχθει ότι για κάθε πίνακα A μπορούμε να βρούμε ένα b και κατάλληλη μεταβολή δb τέτοια ώστε η (4.5) να ισχύει με το σημείο της ισότητας, δηλ. ότι πάντα υπάρχει ένα δεύτερο μέλος b και μια κατάλληλη μεταβολή δb που να μεταβάλλουν την λύση x κατά το μέγιστο δυνατό. Έστω π.χ.:

$$A = \begin{pmatrix} 4.1 & 2.8 \\ 9.7 & 6.6 \end{pmatrix}, b = \begin{pmatrix} 4.1 \\ 9.7 \end{pmatrix}, \text{ οπότε } x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

και $\|b\|_1 = 13.8, \|x\|_1 = 1$. Αν το δεύτερο μέλος υποστεί την μεταβολή $\delta b = (0.01, 0)^T$, τότε η λύση $x+\delta x$ του νέου συστήματος $A(x+\delta x) = b+\delta b$ είναι $x+\delta x = (0.34, 0.97)^T$. Συνεπώς $\|b\|_1 = .01, \|b\|_1 = 1.63$ και $\|b\|_1/\|b\|_1 = .0007246, \|b\|_1/\|x\|_1 = 1.63$. Συνεπώς τα εχρητικά εφάλματα έχουν λόγο $1.63/.0007246 = 2249.4$. Αν την άλλη μεριά μπορούμε να υπολογίσουμε για τον πίνακα A ότι $\kappa_1(A) = \|A\|_1 \|A^{-1}\|_1 = 2249.4$. Δηλ. η (4.5) ισχύει εαν ισότητα.

Τα δεδομένα του προβλήματος (4.1) δεν είναι μόνο τα στοιχεία του b αλλά και τα στοιχεία του A . Ο δείκτης κατάστασης του πίνακα καθορίζει επίσης και το πώς μεταβολές του A επηρεάζουν την λύση του (4.1). Πραγματικά, αν όταν μεταβάλλουμε τον πίνακα A σε $A+\delta A$ (όπου δA ένας "μικρός" $n \times n$ πίνακας) η λύση x του (4.1) γίνει $x+\delta x$, έχουμε ότι

$$(4.7) \quad (A+\delta A)(x+\delta x) = b.$$

Εξ' υποθέσεως έχουμε ότι ο A είναι αντιστρέψιμος, δηλ. ότι το σύστημα (4.1) έχει λύση. Πρέπει να υποθέσουμε επίσης ότι και το (4.7) έχει λύση, δηλ. ότι και ο $A+\delta A$ είναι αντιστρέψιμος. Για συνθήκη η οποία το εξασφαλίζει αυτό είναι η

$$(4.8) \quad \|\delta A\| \|A^{-1}\| < 1,$$

για οποιαδήποτε φυσική νόρμα πινάκων $\|\cdot\|$. Πραγματικά, η (4.8) συνεπάγεται ότι $\|\delta A A^{-1}\| < 1$. Η πρόταση 3.1 δίνει τώρα ότι ο πίνακας $I+(\delta A)A^{-1}$ είναι αντιστρέψιμος \Rightarrow ο πίνακας $(I+(\delta A)A^{-1})A = A + \delta A$ είναι αντιστρέψιμος. Δηλ. η (4.8), που μπορεί να γραφτεί σαν $\|\delta A\| < 1/\|A^{-1}\|$, είναι μια συνθήκη που εξασφαλίζει ότι η μεταβολή δA είναι αρκετά μικρή ώστε όχι μόνο ο A αλλά και ο $A+\delta A$ να είναι αντιστρέψιμος. (Πρβλ. άσκηση 4.4 με $B=A+\delta A$). Υπό την προϋπόθεση ότι η (4.8) ισχύει έχουμε, αφαιρώντας κατά μέλη τις (4.1) και (4.7) ότι $(A+\delta A)\delta x + (\delta A)x = 0 \Rightarrow \delta x = -(A+\delta A)^{-1}(\delta A)x = -[(I+(\delta A)A^{-1})A]^{-1}(\delta A)x = -A^{-1}(I+(\delta A)A^{-1})^{-1}(\delta A)x$.

Άρα, παίρνοντας νόρμες,

$$(4.9) \quad \|\delta x\| = \|A^{-1}(I+(\delta A)A^{-1})^{-1}(\delta A)x\| \leq \|A^{-1}\| \|(I+(\delta A)A^{-1})^{-1}\| \|\delta A\| \|x\|.$$

Η (3.32) δίνει τώρα λόγω της (4.8) ότι

$$(4.10) \quad \|(I+(\delta A)A^{-1})^{-1}\| \leq 1/(1-\|\delta A\| \|A^{-1}\|).$$

Άρα, από τις (4.9) και (4.10) έχουμε

$$\|\delta x\| \leq \|A^{-1}\| \|\delta A\| \|x\| / (1 - \|\delta A\| \|A^{-1}\|),$$

δηλ. για $x \neq 0$

$$(4.11) \quad \|\delta x\|/\|x\| \leq (\kappa(A)/1-\|A^{-1}\| \|\delta A\|) \cdot (\|\delta A\|/\|A\|).$$

Ξανά δηλ. βλέπουμε τον καθοριστικό ρόλο του $\kappa(A)$ ως "συντελεστή μεγέθυνσης" εσφαλμάτων ή μεταβολών των στοιχείων του πίνακα A , δηλ. ως μέτρου του πόσο το σχετικό εφάλμα $\|\delta A\|/\|A\|$ επηρεάζει το $\|\delta x\|/\|x\|$. Το γενικό βέβαια πρόβλημα είναι το

$$(4.12) \quad (A+\delta A)(x+\delta x) = (b+\delta b),$$

η περίπτωση δηλ. που και ο A και το b μεταβάλλονται. Δεν είναι δύσκολο να δείξουμε ότι στην περίπτωση (4.12), αν ισχύει η (4.8) πάλι, έχουμε ότι

$$(4.13) \quad \|\delta x\|/\|x\| \leq \kappa(A)/(1-\|\delta A\| \|A^{-1}\|) (\|\delta b\|/\|b\| + \|\delta A\|/\|A\|),$$

δηλ. μια ανισότητα της οποίας οι (4.5) και (4.11) αποτελούν ειδικές περιπτώσεις, αν $\delta A=0$ και $\delta b=0$, αντίστοιχα.

ΑΣΚΗΣΗ 4.5 Αν οι πίνακες A, B είναι αντιστρέψιμοι, ναδειχθεί ότι $\|B^{-1} - A^{-1}\|/\|B^{-1}\| \leq \kappa(A) (\|A-B\|/\|A\|)$. @

ΑΣΚΗΣΗ 4.6 Αν ο A είναι αντιστρέψιμος αυτοσυζυγής πίνακας (δηλ. αν $A = A^*$) ναδειχθεί ότι ο δείκτης κατάστασής του $\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2$ δίνεται από τον τύπο

$$(4.14) \quad \kappa_2(A) = \max_i |\lambda_i| / \min_i |\lambda_i|,$$

όπου $\lambda_i, 1 \leq i \leq n$ οι ιδιοτιμές του A . @

ΑΣΚΗΣΗ 4.7 Δίνεται ο πίνακας A και ο A^{-1} :

$$A = \begin{pmatrix} 6 & 13 & -17 \\ 13 & 29 & -38 \\ -17 & -38 & 50 \end{pmatrix}, \quad A^{-1} = \begin{pmatrix} 6 & -4 & -1 \\ -4 & 11 & 7 \\ -1 & 7 & 5 \end{pmatrix}.$$

Ο πίνακας A έχει ιδιοτιμές $\lambda_1 = .0568$, $\lambda_2 = .2007$, $\lambda_3 = 84.74$.
Χρησιμοποιώντας τις εκθέσεις (3.26), (4.14) να βρείτε τα $\|A\|_2$, $\|A^{-1}\|_2$,
 $\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2$. Βρείτε επίσης τους $\kappa_\infty(A)$, $\kappa_1(A)$ @

ΑΣΚΗΣΗ 4.8 Έστω $x \neq 0$ η (ακριβής) λύση του συστήματος $Ax=b$, όπου A αντιστρέψιμος. Έστω \tilde{x} μια προσέγγιση της x . Ορίζουμε όπως στην άσκηση 2.16, το υπόλοιπο r της \tilde{x} εάν $r = A\tilde{x} - b$. Ναδειχθεί ότι για κάθε διανυσματική νόρμα (και αντίστοιχη φυσική νόρμα πινάκων) $\|\cdot\|$:

$$(4.15) \|\tilde{x} - x\| / \|x\| \leq \kappa(A) (\|r\| / \|b\|), \quad \kappa(A) = \|A\| \|A^{-1}\|.$$

(Η (4.15) συμπληρώνει την άσκηση 2.16: Έστω ότι το υπόλοιπο $\|r\|$ είναι μικρό. Αν όμως ο A έχει κακή κατάσταση, δηλ. αν $\kappa(A) \gg 1$, τότε η (4.15) δείχνει ότι εχτικό εφάλμα $\|\tilde{x} - x\| / \|x\|$ μπορεί να είναι μεγάλο. Άρα μικρό υπόλοιπο \neq μικρό εφάλμα). @

ΑΣΚΗΣΗ 4.9 Για τον πίνακα A της άσκησης 4.7, υποθέστε ότι έχουμε μια προσέγγιση \tilde{x} στη λύση x του συστήματος $Ax=b$, $b=(1,1)^T$, που ικανοποιεί: $\|A\tilde{x} - b\|_2 \leq 10^{-2}$. Βρείτε άνω φράγματα

για το απόλυτο ($\|\tilde{x} - x\|_2$) και το εχτικό εφάλμα ($\|\tilde{x} - x\|_2 / \|x\|_2$) της \tilde{x} : @

ΑΣΚΗΣΗ 4.10 Να δείξετε ότι για οποιαδήποτε νόρμα πινάκων και οποιοδήποτε μιγαδικό αριθμό $\lambda \neq 0$

$$(4.16) \kappa(\lambda A) = \kappa(A).$$

Επίσης, αν ο D είναι διαγώνιος πίνακας με στοιχεία $d_{ii} \neq 0$, $1 \leq i \leq n$, να δείξετε ότι

$$(4.17) \kappa_\infty(D) = \kappa_1(D) = \kappa_2(D) = \max_i |d_{ii}| / \min_i |d_{ii}|. @$$

ΑΣΚΗΣΗ 4.11 Θα έλεγε κανείς ότι το μέγεθος της ορίζουσας ενός πίνακα καθορίζει την ευστάθεια του προβλήματος $Ax=b$ σύμφωνα με τον συλλογισμό: μικρή ορίζουσα (δηλ. κοντά στο 0) $\Rightarrow A$ "εχεδόν" μη αντιστρέψιμος \Rightarrow σύστημα ασταθές.

Αυτό όμως δεν είναι αλήθεια: θεωρήστε τον 100×100 διαγώνιο πίνακα με $d_{ii} = 0.1$, $1 \leq i \leq 100$. Βρείτε τη ορίζουσα $\det(D)$ και τους δείκτες κατάστασης $\kappa_\infty(D)$, $\kappa_1(D)$, $\kappa_2(D)$. (Προφανώς το σύστημα $Dx=b$ δεν είναι καθόλου ασταθές. Λύνεται άμεσα: $x=10b$ και για μια μεταβολή δb έχουμε ότι αν $D(x+\delta x) = b + \delta b \Rightarrow \delta x = 10\delta b \Rightarrow \|\delta x\| / \|x\| = \|\delta b\| / \|b\|$). Συνεπώς το μέγεθος της ορίζουσας δεν έχει καμιά σημασία για την ευστάθεια του συστήματος. Το σύστημα αυτό είναι όσο ευσταθές θα μπορούσε να γίνει. @

ΑΣΚΗΣΗ 4.12 Είδαμε ότι πάντα $\kappa(A) \geq 1$. Τι είδους πίνακες έχουν μικρό $\kappa(A)$; Να αποδειχθεί ότι:

(i) Αν ο A είναι διαγώνιος με ίσα στοιχεία ($\neq 0$) τότε $\kappa_\infty(A) = \kappa_1(A) = \kappa_2(A) = 1$.

(ii) Αν ο A ικανοποιεί $\|Ax\| = \|x\|$ για κάθε x , να αποδειχθεί ότι $\kappa(A) = \|A\| \|A^{-1}\| = 1$.

(iii) Για ένα πίνακα U τέτοιου ώστε $UU^* = I$, ναδειχθεί ότι $\kappa_2(U) = 1$. @

4.2 Επιρροή του δείκτη κατάστασης στην απαλοιγή.

Ο δείκτης κατάστασης παίζει μεγάλο ρόλο και στη μελέτη της διάδοσης των εσφαλμάτων ετρογχύλευσης κατά την απαλοιγή Gauss (ή την ανάλυση $PA=LU$) με αριθμητική πεπερασμένης ακρίβειας. Ας θεωρήσουμε πρώτα την "ιδανική" περίπτωση όπου εσφάλματα ετρογχύλευσης συμβαίνουν μόνο κατά την παράσταση των στοιχείων του A και του b στον υπολογιστή και όπου οι υπόλοιπες αριθμητικές πράξεις της απαλοιγής γίνονται ακριβώς. Με τον συμβολισμό της παρ. 2.2 (βλ. και Σημειώσεις Αριθμητικής Ανάλυσης I) κάθε πραγματικός a (μέσα στο εύρος των αριθμών της μηχανής) παριστάνεται στον υπολογιστή από τον $f(a)$ όπου

$$(4.18) f(a) = a(1+\epsilon), \quad |\epsilon| \leq \epsilon,$$

και όπου η "μονάδα του εσφαλματος ετρογχύλευσης" u είναι b^{1-t} για αποκοπή και $b^{1-t}/2$ για ετρογχύλευση. Στον υπολογιστή επομένως το δεύτερο μέλος b θα παραταθεί από ένα διάνυσμα \tilde{b} τέτοιο ώστε $\tilde{b}_i = b_i(1+\epsilon_i)$, $|\epsilon_i| \leq \epsilon$, $1 \leq i \leq n$. Συνεπώς,

$$(4.19) \tilde{b} = b + \delta b, \quad \|\delta b\|_\infty \leq \epsilon \|b\|_\infty.$$

(Η χρήση της νόρμας $\|\cdot\|_\infty$ είναι ευδεδειγμένη για την εκτίμηση των εσφαλμάτων ετρογχύλευσης). Ανάλογα, ο A θα παραταθεί από ένα $n \times n$ πίνακα \tilde{A} , όπου

$$(4.20) \tilde{A} = A + \delta A, \quad \|\delta A\|_\infty \leq \epsilon \|A\|_\infty.$$

Έστω x η θεωρητική λύση του συστήματος $Ax=b$. Βάσει των υποθέσεων μας η υπολογιστική λύση \tilde{x} είναι η ακριβής λύση του συστήματος $\tilde{A}\tilde{x}=\tilde{b}$, δηλ. του συστήματος

$$(4.21) (A + \delta A)\tilde{x} = (b + \delta b).$$

Για να βρούμε μια εκτίμηση του σχετικού εσφαλματος $\|\tilde{x} - x\|_\infty / \|x\|_\infty$ ανατρέχουμε τώρα στη θεωρία της προηγούμενης παραγράφου, δηλ. στην περίπτωση (4.12), όπου $x + \delta x = \tilde{x}$. Ας υποθέσουμε ότι ισχύει η συνθήκη

$$(4.22) \kappa_\infty(A) = r < 1.$$

Συνεπώς ισχύει η (4.8), δηλ. ο πίνακας $A + \delta A$ είναι αντιστρέψιμος και εύκολα, η (4.13) δίνει

$$\|\delta x\|_\infty / \|x\|_\infty = \|x - \tilde{x}\|_\infty / \|x\|_\infty \leq (\kappa_\infty(A) / (1-r)) (\|\delta b\|_\infty / \|b\|_\infty + \|\delta A\|_\infty / \|A\|_\infty)$$

δηλ., λόγω των (4.19) και (4.20), ότι

$$(4.23) \|x - \tilde{x}\|_\infty / \|x\|_\infty \leq 2\kappa_\infty(A) / (1-r).$$

Είναι λοιπόν φανερό ότι αν το χιόνιμο $\kappa_\infty(A)$ πλησιάζει την μονάδα, θα πρέπει να περιμένουμε μεγάλο σχετικό εσάγμα για την υπολογιστική λύση \tilde{x} .

ΑΣΚΗΣΗ 4.13 Δείξτε ότι στην περίπτωση που τα στοιχεία του πίνακα A παριστάνονται ακριβώς, (δηλ. όταν το εσάγμα ετρογχύλευσης προέρχεται μόνο από την παράσταση του b), ισχύει η εκτίμηση $\|x - \tilde{x}\|_\infty / \|x\|_\infty \leq \kappa_\infty(A)$, χωρίς την υπόθεση (4.22). Εφαρμογή: Έστω ότι τα στοιχεία του A και b είναι ακέραιοι με την εξαίρεση του b_i που είναι ίσο με 10^{-1} . Παριστάνουμε τα στοιχεία του A και b σε απλή ακρίβεια στο VAX 11 και υποθέτουμε ότι οι υπόλοιπες πράξεις της απαλοιγής γίνονται ακριβώς. Ποιό είναι το αναμενόμενο σχετικό εσάγμα $\|x - \tilde{x}\|_\infty / \|x\|_\infty$ αν $\kappa_\infty(A) = 10^4$; Το ίδιο ερώτημα για απλή και διπλή ακρίβεια στο IBM 4361. (Η αριθμητική του IBM 4361 είναι ίδια με την αριθμητική του συστήματος IBM 370). @

Προχωρούμε τώρα στην περίπτωση που η απαλοιγή (με οδήγηση ή χωρίς) γίνεται με αριθμητική πεπερασμένης ακρίβειας και που τα εσάγματα ετρογχύλευσης ευεωρεούνται κατά τις πράξεις. Θέλουμε να εκτιμήσουμε το (σχετικό) εσάγμα της υπολογιστικής λύσης x . Χρησιμοποιώντας την σχέση (4.18) και το γεγονός ότι $f(a \oplus b) = (a \oplus b)(1+\epsilon)$, $|\epsilon| \leq \epsilon$, όπου \oplus οποιαδήποτε αριθμητική πράξη $+, -, \cdot, /$, είναι δυνατόν να παρακολουθήσουμε την επίδραση και την ευεωρευση των εσφαλμάτων ετρογχύλευσης σε κάθε πράξη της απαλοιγής (δηλ. στον υπολογισμό των

α_{ij} , των στοιχείων σ_{ij} κλπ.) και της οπισθοδρόμησης. Η "λογιστική" βεβαία ενός τέτοιου εγχειρήματος είναι πολύπλοκη και η τεχνική αυτή (που λέγεται "απ' ευθείας" (forward) εκτίμηση του εσφαλματος) μάλλον θα δώσει ένα εντελώς μη ρεαλιστικό (δηλ. ένα πολύ μεγάλο) άνω φράγμα για το εσάγμα $\|x - \tilde{x}\|_\infty / \|x\|_\infty$. Για άλλη τεχνική, η

λεχόμενη "αντίστροφη" (inverse) ανάλυση του εφάλματος του προβλήματος σφείζεται στον άγγλο μαθηματικό J.H. Wilkinson (δεκαετία '60) που απέδειξε ότι η υπολογιστική (προσεγγιστική) λύση \tilde{x} που τελικά παίρνουμε, μπορεί να θεωρηθεί ως ακριβής λύση (δηλ. ως λύση με αριθμητική απεριόριστη ακρίβεια) όχι του συστήματος $Ax=b$ βέβαια, αλλά του "παραπλήσιου" συστήματος

$$(4.24) (A+\delta A)\tilde{x} = b,$$

όπου ο πίνακας μεταβολής δA έχει γενικά "μικρά" στοιχεία. Αν έχουμε μια καλή εκτίμηση του $\| \delta A \|_{\infty}$ π.χ., τότε η θεωρία της παρ. 4.1 (βλ. (4.7) - (4.11)) μας δίνει μια καλή εκτίμηση του εφάλματος $\|x - \tilde{x}\|_{\infty} / \|x\|_{\infty}$.

ΠΣΚΗΣΗ 4.14 (Για φανατικούς). Η "αντίστροφη" ανάλυση του εφάλματος έχει εφαρμογή και σε άλλα παραδείγματα υπολογισμών. Π.χ. θεωρείστε δυο 2×2 άνω τριγωνικούς πίνακες A και B και έστω $f(A, B)$ ο πίνακας που υπολογίζεται από τον υπολογιστή ως το γινόμενο τους (δηλ. με εφάλματα στρογγύλευσης στην παράσταση των A και B και στις πράξεις του πολλαπλασιασμού των δυο πινάκων). Δείξτε ότι υπάρχουν δυο 2×2 πίνακες \tilde{A}, \tilde{B} , παραπλήσιοι στους A και B αντίστοιχα, τέτοιοι ώστε $f(A, B) = \tilde{A}\tilde{B}$, δηλ. τέτοιοι ώστε το ακριβές γινόμενο τους $\tilde{A}\tilde{B}$ να είναι ίσο με το προσεγγιστικό "γινόμενο" $f(A, B)$. @

Υποθέτοντας για απλούστευση ότι τα στοιχεία των A, B παριστάνονται ακριβώς αρχικά στον υπολογιστή, ο Wilkinson έδειξε (για την απόδειξη, που παραλείπεται, βλ. π.χ.: τα βιβλία [A4], [B1], [B3] ή την "βίβλο" [B5] στη βιβλιογραφία του μαθήματος) ότι αν $n^2 u < 1$, τότε για τον πίνακα δA του (4.24) ισχύει

$$(4.25) \| \delta A \|_{\infty} \leq 1.01(n^3 + 3n^2) \epsilon \| A \|_{\infty} u$$

όπου

$$(4.26) \epsilon = \max_{i,j,k}^{(k)} |a_{ij}| / \| A \|_{\infty},$$

(k)

και όπου τα a_{ij} είναι οι μετασχηματισμοί των a_{ij} κατά την απαλοιφή.

Η σταθερά $1.01(n^3 + 3n^2)$ στην (4.25) είναι ένα πολύ "απαισιόδοξο" φράγμα και μπορεί να αγνοηθεί στην πράξη.

Με μια συγκεκριμένη στρατηγική οδήγησης μπορούμε τώρα να ελέγξουμε το μέγεθος του ϵ . Ο Wilkinson έδειξε ότι για την απαλοιφή Gauss με ολική οδήγηση το ϵ αυξάνεται πολύ αρχά ως συνάρτηση του n ενώ υπάρχουν ακραία παραδείγματα όπου για μερική οδήγηση το ϵ αυξάνεται εκθετικά με το n . Όπως όμως έχει παρατηρηθεί στην πράξη και όπως ο ίδιος ο Wilkinson είχε υπολογίσει, είναι εξαιρετικά σπάνιο φαινόμενο να πάρουμε, με μερική οδήγηση π.χ., ϵ μεγαλύτερο του 10.

Για να κάνουμε τώρα ορισμένες ακόμα παρατηρήσεις σχετικά με το εφάλμα της προσεγγιστικής λύσης \tilde{x} , υποθέτουμε ότι ισχύει η (4.24) και κάνουμε την υπόθεση ερχασίας ότι για κάποια φυσική νόρμα $\| \cdot \|$ ισχύει

$$(4.25) \| \delta A \| / \| A \| = \rho \beta^{-t}.$$

(Στην πράξη εκτός από σπανιότητες περιπτώσεις, έχει παρατηρηθεί ότι με μερική οδήγηση, το ρ είναι της τάξης μεγέθους του β). Από τις (4.24) και (4.25) μπορούμε να εκτιμήσουμε το υπόλοιπο της προσεγγιστικής λύσης \tilde{x} και το σχετικό της εφάλμα. Για το υπόλοιπο r έχουμε, από την (4.24),

$$r = A\tilde{x} - b = -(\delta A)\tilde{x}.$$

Συνεπώς $\| r \| \leq \| \delta A \| \| \tilde{x} \|$, οπότε η (4.25) (αν $A, \tilde{x} \neq 0$) δίνει

$$(4.26) \| r \| / (\| A \| \| \tilde{x} \|) \leq \rho \beta^{-t}.$$

Η άλλα λόγια το "σχετικό" υπόλοιπο (δηλ. το υπόλοιπο δια του γινόμενου $\| A \| \| \tilde{x} \|$, που μετράει την κλίμακα του προβλήματος, δηλ. το μέγεθος των δεδομένων A και της υπολογιστικής λύσης \tilde{x}) είναι της τάξης $\rho \beta^{-t}$. Επειδή συνήθως το ρ είναι της τάξης του β , συμπεραίνουμε ότι το υπόλοιπο της υπολογιστικής (προσεγγιστικής) λύσης που παράγει η απαλοιφή Gauss είναι, ανεξάρτητα απ' την κατάσταση του πίνακα A , μικρό, και μάλιστα της τάξης β^{1-t} , δηλ. του εφάλματος στρογγύλευσης! (Πρβλ. Άσκηση 2.15). Συνεπώς, αν σε κάποια εφαρμογή μας ενδιαφέρει να υπολογίσουμε ένα \tilde{x} που να έχει μικρό υπόλοιπο $r = A\tilde{x} - b$, μπορούμε άφοβα να κάνουμε απαλοιφή με μερική οδήγηση, ανεξάρτητα από την κατάσταση του A . (Μάλιστα, η ανισότητα (4.26) ισχύει και για μη αντιστρέψιμο A - αρκεί βέβαια να υπάρχει λύση \tilde{x} του (4.24)-).

Προχωρούμε τώρα στην μελέτη του εσφάλματος της υπολογιστικής λύσης \tilde{x} . Αν ο A είναι αντιστρέψιμος, οι (4.24), (4.25) και (4.11) δίνουν

$$(4.27) \quad \|x-\tilde{x}\|/\|x\| \leq [k(A)/(1-k(A)\rho\beta^{-t})] \rho\beta^{-t},$$

υπό την προϋπόθεση ότι ισχύει η $\|δA\| \|A^{-1}\| = k(A) \rho\beta^{-t} < 1$. Δηλ. βλέπουμε τώρα ότι το εσχετικό εσφάλμα της \tilde{x} θα είναι μικρό, αν ο δείκτης κατάστασης $k(A)$ είναι μικρός, αλλά ότι μπορεί να είναι αρκετά μεγάλο, αν $k(A) \gg 1$.

Η ανισότητα (4.27) δεν είναι και πολύ χρήσιμη στην πράξη γιατί το εσχετικό εσφάλμα $\|x-\tilde{x}\|/\|x\|$ υπολογίζεται ως προς το $\|x\|$, δηλ. ως προς μια άγνωστη ποσότητα στον παρανομαστή, ενώ ο αλγόριθμος υπολογίζει το \tilde{x} . Αλλά μπορούμε να δούμε από τον ορισμό του υπόλοιπου $r=A\tilde{x}-b$ ότι

$$x-\tilde{x} = A^{-1}b-\tilde{x} = A^{-1}(b-A\tilde{x}) = -A^{-1}r.$$

Συνεπώς, από την (4.26),

$$\|x-\tilde{x}\| \leq \|A^{-1}\| \|r\| \leq \|A^{-1}\| \|A\| \|\tilde{x}\| \rho\beta^{-t},$$

δηλ. ισχύει η ανισότητα

$$(4.28) \quad \|x-\tilde{x}\|/\|\tilde{x}\| \leq k(A) \rho\beta^{-t},$$

που ήταν η ζητούμενη. Επειδή το $\|\tilde{x}\|$ είναι γνωστό, η (4.28) δίνει μια εκτίμηση του απόλυτου εσφάλματος $\|x-\tilde{x}\| \leq \|\tilde{x}\| k(A) \rho\beta^{-t}$ συναρτήσει της υπολογιστικής λύσης \tilde{x} - ένα παράδειγμα εκτίμησης εκ των υστέρων (a posteriori), δηλ. με γνώση του αλγοριθμικού αποτελέσματος \tilde{x} , σε αντίθεση με την εκτίμηση του $\|x-\tilde{x}\|$ που δίνει η (4.27) που είναι εκτίμηση εκ των προτέρων (a priori), δηλ. διατυπώνεται συναρτήσει της ακριβούς λύσης x.

Κλείνουμε αυτό το κεφάλαιο με ορισμένες παρατηρήσεις

1. Σε ορισμένες περιπτώσεις είναι δυνατόν, πολλαπλασιάζοντας πριν λύσουμε το σύστημα από αριστερά και δεξιά τον πίνακα A επί διαγώνιους πίνακες (που αντιστοιχεί σε πολλαπλασιασμό των γραμμών και

των στηλών του A με σταθερές, δηλ. σε μια εκ των προτέρων εστίαση-αλλαγή κλίμακας (scaling) - των στοιχείων του A) να πάρουμε ένα σύστημα με μικρότερο δείκτη κατάστασης, βλ. [B1],[B3]. Τέτοιοι κατάλληλοι διαγώνιοι πίνακες μπορούν να βρεθούν για ειδικές περιπτώσεις πινάκων A αλλά το γενικό πρόβλημα του υπολογισμού τους δεν έχει διελευκασθεί ακόμα αρκετά.

2. Αν ο δείκτης κατάστασης ενός πίνακα δεν είναι πολύ μεγάλος, εσχετικά με την χρησιμοποιούμενη ακρίβεια, δηλ. αν $k(A)u \ll 1$, είναι δυνατόν, χρησιμοποιώντας τεχνικές επαναληπτικής βελτίωσης (iterative improvement) να βελτιώσουμε την ακρίβεια της υπολογιστικής λύσης \tilde{x} . Για τέτοια τεχνική υπολογίζει το υπόλοιπο της \tilde{x} και λύνει ένα νέο σύστημα με πίνακα A και δεύτερο μέλος το υπόλοιπο r, υπολογίζοντας μια "διόρθωση", η οποία, προστιθέμενη στη λύση \tilde{x} , δίνει μια νέα προσέγγιση, που ευχά είναι ακριβέστερη της \tilde{x} εφ' όσον το υπόλοιπο r έχει υπολογιστεί με μεγαλύτερη ακρίβεια (π.χ. διπλή) απ' ό,τι οι υπόλοιπες πράξεις. Βλ. π.χ. [A1],[B1],[B3],[B6].

3. Αν χωρίζουμε τον δείκτη κατάστασης ενός πίνακα, μπορούμε, χρησιμοποιώντας π.χ. την ανισότητα (4.28), να εκτιμήσουμε το εσφάλμα $\|x-\tilde{x}\|$ της υπολογιστικής λύσης \tilde{x} . Βεβαίως $k(A)=\|A\| \|A^{-1}\|$, αλλά θέλουμε ν' αποφύγουμε τον υπολογισμό του A^{-1} , δηλ. μια πράξη κατά πολύ ακριβέτερη από την λύση του συστήματος. Στην πράξη λοιπόν δεν υπολογίζουμε ακριβώς τον δείκτη κατάστασης αλλά κάνουμε μια προσεγγιστική του εκτίμηση αφού εξ' άλλου ευθύως μόνο η γνώση της τάξης του $k(A)$ μας χρειάζεται.

Για τεχνική για την εκτίμηση του $k(A)$ είναι η εξής: Λύνουμε μερικά συστήματα $Aw^i=y^i, i=1,2,\dots,k$ (k συνήθως 2 ή 3) με τυχαία δεύτερα μέλη y^i . (Υπενθυμίζεται ότι ο A έχει ήδη τριγωνοποιηθεί για τη λύση του συστήματος $Ax=b$. Έτσι, ο υπολογισμός κάθε w^i απαιτεί μόνο $O(n^2)$ πράξεις). Υπολογίζουμε κατόπιν τα πηλίκια $\|w^i\|/\|y^i\|$ και προσεγγίζουμε κατόπιν $\|A^{-1}\| \approx \max_{1 \leq i \leq k} (\|w^i\|/\|y^i\|)$. (Βέβαια, η πραγματική τιμή του $\|A^{-1}\|$ είναι ίση με το $\sup_{y \neq 0} (\|A^{-1}y\|/\|y\|) = \sup_{y \neq 0} (\|w\|/\|y\|, w=A^{-1}y \Leftrightarrow w$ λύση συστήματος $Aw=y)$).

Στο υποπρόγραμμα DCOMR του εργαστηρίου του μαθήματος ο δείκτης κατάστασης $k_1(A)$ προσεγγίζεται ως $k_1(A) \approx \|A\|_1 \|z\|_1 / \|y\|_1$, όπου $\|A^{-1}\|_1 \approx \|z\|_1 / \|y\|_1$. Τα διανύσματα z,y υπολογίζονται ως λύσεις δυο συστημάτων $A^T y=e, Az=y$, όπου e ένα κατάλληλα επιλεγμένο διάνυσμα του οποίου οι συνιστώσες είναι ίσες με ± 1 . Για την θεωρητική απόδειξη του ότι αυτό είναι μια καλή εκτίμηση του $k_1(A)$, βλ. [B3],[B6].

Προχωρούμε τώρα στην μελέτη του εσφάλματος της υπολογιστικής λύσης \tilde{x} . Αν ο A είναι αντιστρέψιμος, οι (4.24), (4.25) και (4.11) δίνουν

$$(4.27) \quad \|x-\tilde{x}\|/\|x\| \leq [k(A)/(1-k(A)\rho\beta^{-l})] \rho\beta^{-l},$$

υπό την προϋπόθεση ότι ισχύει η $\|bA\| \|A^{-1}\| = k(A) \rho\beta^{-l} < 1$. Δηλ. βλέπουμε τώρα ότι το εσχετικό εσφάλμα της \tilde{x} θα είναι μικρό, αν ο δείκτης κατάστασης $k(A)$ είναι μικρός, αλλά ότι μπορεί να είναι αρκετά μεγάλο, αν $k(A) \gg 1$.

Η ανισότητα (4.27) δεν είναι και πολύ χρήσιμη στην πράξη γιατί το εσχετικό εσφάλμα $\|x-\tilde{x}\|/\|x\|$ υπολογίζεται ως προς το $\|x\|$, δηλ. ως προς μια άγνωστη ποσότητα στον παρανομαστή, ενώ ο αλγόριθμος υπολογίζει το \tilde{x} . Αλλά μπορούμε να δούμε από τον ορισμό του υπολοίπου $r=A\tilde{x}-b$ ότι

$$x-\tilde{x} = A^{-1}b-\tilde{x} = A^{-1}(b-A\tilde{x}) = -A^{-1}r.$$

Συνοπώς, από την (4.26),

$$\|x-\tilde{x}\| \leq \|A^{-1}\| \|r\| \leq \|A^{-1}\| \|A\| \|\tilde{x}\| \rho\beta^{-l},$$

δηλ. ισχύει η ανισότητα

$$(4.28) \quad \|x-\tilde{x}\|/\|\tilde{x}\| \leq k(A) \rho\beta^{-l},$$

που ήταν η ζητούμενη. Επειδή το $\|\tilde{x}\|$ είναι γνωστό, η (4.28) δίνει μια εκτίμηση του απολύτου εσφάλματος $\|x-\tilde{x}\| \leq \|\tilde{x}\| k(A) \rho\beta^{-l}$ συναρτήσει της υπολογιστικής λύσης \tilde{x} - ένα παράδειγμα εκτίμησης εκ των υστέρων (a posteriori), δηλ. με γνώση του αλγοριθμικού αποτελέσματος \tilde{x} , σε αντίθεση με την εκτίμηση του $\|x-\tilde{x}\|$ που δίνει η (4.27) που είναι εκτίμηση εκ των προτέρων (a priori), δηλ. διατυπώνεται συναρτήσει της ακριβούς λύσης x .

Κλείνουμε αυτό το κεφάλαιο με ορισμένες παρατηρήσεις

1. Σε ορισμένες περιπτώσεις είναι δυνατόν, πολλαπλασιάζοντας πρίν λύσουμε το σύστημα από αριστερά και δεξιά τον πίνακα A επί διαγώνιους πίνακες (που αντιστοιχεί σε πολλαπλασιασμό των γραμμών και

των στηλών του A με σταθερές, δηλ. σε μια εκ των προτέρων εστάθμιση-αλλαγή κλίμακας (scaling) - των στοιχείων του A) να πάρουμε ένα σύστημα με μικρότερο δείκτη κατάστασης, βλ. [B1],[B3]. Τέτοιοι κατάλληλοι διαγώνιοι πίνακες μπορούν να βρεθούν για ειδικές περιπτώσεις πινάκων A αλλά το γενικό πρόβλημα του υπολογισμού τους δεν έχει διελευκανθεί ακόμα αρκετά.

2. Αν ο δείκτης κατάστασης ενός πίνακα δεν είναι πολύ μεγάλος, εσχετικά με την χρησιμοποιούμενη ακρίβεια, δηλ. αν $k(A)u \ll 1$, είναι δυνατόν, χρησιμοποιώντας τεχνικές επαναληπτικής βελτίωσης (iterative improvement) να βελτιώσουμε την ακρίβεια της υπολογιστικής λύσης \tilde{x} . Για τέτοια τεχνική υπολογίζει το υπόλοιπο της \tilde{x} και λύνει ένα νέο σύστημα με πίνακα A και δεύτερο μέλος το υπόλοιπο r , υπολογίζοντας μια "διόρθωση", η οποία, προστιθέμενη στη λύση \tilde{x} , δίνει μια νέα προσέγγιση, που ευχά είναι ακριβέστερη της \tilde{x} εφ' όσον το υπόλοιπο r έχει υπολογιστεί με μεγαλύτερη ακρίβεια (π.χ. διπλή) απ' ότι οι υπόλοιπες πράξεις. Βλ. π.χ. [A1],[B1],[B3],[B6].

3. Αν χυπρίζουμε τον δείκτη κατάστασης ενός πίνακα, μπορούμε, χρησιμοποιώντας π.χ. την ανισότητα (4.28), να εκτιμήσουμε το εσφάλμα $\|x-\tilde{x}\|$ της υπολογιστικής λύσης \tilde{x} . Βεβαίως $k(A)=\|A\| \|A^{-1}\|$, αλλά θέλουμε ν' αποφύγουμε τον υπολογισμό του A^{-1} , δηλ. μια πράξη κατά πολύ ακριβύτερη από την λύση του συστήματος. Στην πράξη λοιπόν δεν υπολογίζουμε ακριβώς τον δείκτη κατάστασης αλλά κάνουμε μια προσεγγιστική του εκτίμηση αφού εξ' άλλου ευρήτως μόνο η γνώση της τάξης του $k(A)$ μας χρειάζεται.

Για τεχνική για την εκτίμηση του $k(A)$ είναι η εξής: Λύνουμε μερικά συστήματα $Aw^i=y^i, i=1,2,\dots,k$ (k συνήθως 2 ή 3) με τυχαία δεύτερα μέλη y^i . (Υπενθυμίζεται ότι ο A έχει ήδη τριγωνοποιηθεί για τη λύση του συστήματος $Ax=b$. Έτσι, ο υπολογισμός κάθε w^i απαιτεί μόνο $O(n^2)$ πράξεις). Υπολογίζουμε κατόπιν τα ηνλικά $\|w^i\|/\|y^i\|$ και προσεγγίζουμε κατόπιν $\|A^{-1}\| \approx \max_{1 \leq i \leq k} (\|w^i\|/\|y^i\|)$. (Βέβαια, η πραγματική τιμή του $\|A^{-1}\|$ είναι ίση με το $\sup_{y \neq 0} (\|A^{-1}y\|/\|y\|) = \sup_{y \neq 0} (\|w\|/\|y\|, w=A^{-1}y \Leftrightarrow w$ λύση συστήματος $Aw=y)$).

Στο υποπρόγραμμα DECOMP του εργατηρίου του μαθήματος ο δείκτης κατάστασης $k_1(A)$ προσεγγίζεται ως $k_1(A) \approx \|A\|_1 \|z\|_1 / \|y\|_1$, όπου $\|A^{-1}\|_1 \approx \|z\|_1 / \|y\|_1$. Τα διανύσματα z, y υπολογίζονται ως λύσεις δύο συστημάτων $A^T y = e, Az = y$, όπου e ένα κατάλληλα επιλεγμένο διάνυσμα του οποίου οι συνιστώσες είναι ίσες με ± 1 . Για την θεωρητική απόδειξη του ότι αυτό είναι μια καλή εκτίμηση του $k_1(A)$, βλ. [B3],[B6].