

RA

**stichting  
mathematisch  
centrum**



---

REKENAFDELING

RA

NR 19/71

JULI

W. HOFFMANN EN J.P. HOLLENBERG  
FOUTENBEREKENING VOOR LINEAIRE KLEINSTE-  
KWADRATENPROBLEMEN

---

**2e boerhaavestraat 49 amsterdam**

## Inhoud

	blz.
Inleiding	1
§1. Relatie met het oorspronkelijke probleem	2
§2. Voorwaartse foutenanalyse van de "klassieke" methode	4
§3. Voorwaartse foutenanalyse van de QR-methode	6
§4. Foutenanalyse met residuberekening	8
Literatuur	10

### Inleiding

Laten gegeven zijn een reële  $m \times n$  matrix  $A$ ,  $m \geq n$  en een reële  $m$ -vector  $b$ .

Voor een willekeurige  $n$ -vector  $x$  definiëren we  $r = Ax - b$ . Onder de kleinste-kwadratenoplossing van  $[A|b]$  verstaan we die vector  $x$  waarvoor  $\|r\|_2$  minimaal is;  $x$  voldoet aan  $A^T Ax = A^T b$  (zie [1], p.57).

De 2-norm van een vector  $v$  (notatie:  $\|v\|_2$ ) wordt gedefinieerd door:

$$\|v\|_2 = \sqrt{v^T v} = \sqrt{\sum_i (v_i)^2}.$$

De 2-norm van een matrix  $M$  (ook spectraal norm genoemd) wordt gedefinieerd door:

$$\|M\|_2 = \sqrt{\lambda_{\max}(M^T M)}.$$

Van belang is ook de Euclidische norm (oftewel Frobenius-norm) van een matrix:

$$\|M\|_E = \sqrt{\sum_i \sum_j (M_{ij})^2};$$

deze is namelijk makkelijker te berekenen dan de 2-norm.

Voor een  $n \times n$  matrix  $M$  weten we:

$$\|M\|_2 \leq \|M\|_E \leq \sqrt{n} \|M\|_2.$$

Zie voor een verdere beschouwing over normen bijv. ([1], p. 7 e.v.)

In het vervolg zal met  $\|\tilde{M}\|$ , resp.  $\|\tilde{v}\|$  altijd bedoeld worden  $\|M\|_2$ , resp.  $\|v\|_2$ .

Matrices en vectoren die in de oorspronkelijke probleemstelling voorkomen zullen we aangeven met  $\tilde{M}$ , resp.  $\tilde{v}$ .

De representatie van deze matrices en vectoren in de computer zullen we aangeven met  $M$ , resp.  $v$ .

Matrices en vectoren die een resultaat van bewerkingen in de computer zijn, worden aangegeven met  $\bar{M}$ , resp.  $\bar{v}$ .

In §1 geven we aan in hoeverre  $x$  verschilt van  $\tilde{x}$ ; het is duidelijk dat dit onafhankelijk is van de gevolgde methode om  $\bar{x}$  te berekenen. In §2 geven we een foutenanalyse voor de "klassieke" methode om  $\bar{x}$  te bepalen.

In §3 doen we dit voor de QR-methode volgens Householder.

In §4 tenslotte geven we aan welke verfijning aangebracht kan worden indien  $\bar{r} = \text{fl}(A\bar{x}-b)$  berekend wordt.

### §1. Relatie met het oorspronkelijke probleem

Laat  $x$  de kleinste-kwadratenoplossing zijn van  $[A;b]$  en  $\tilde{x}$  die van  $[\tilde{A};\tilde{b}]$ . Zij  $\tilde{A}^T \tilde{A} = \tilde{M}$  en zij verder  $A = \tilde{A} + \delta A$ ,  $b = \tilde{b} + \delta b$ .  $\delta A$  en  $\delta b$  stellen voor, de som van de waarnemingsfouten en de representatiefouten. De representatiefouten ontstaan door afronding tot in de computer representeerbare getallen en zijn gewoonlijk veel kleiner dan de waarnemingsfouten.

We weten dat geldt:

$$A^T A x = A^T b,$$

hieruit volgt:

$$\tilde{A}^T \tilde{A} x + \tilde{A}^T \delta A x + \delta A^T A x = \tilde{A}^T \tilde{b} + \tilde{A}^T \delta b + \delta A^T b.$$

Dit levert ons:

$$1.1 \quad \tilde{M} x = \tilde{A}^T \tilde{b} - \delta A^T r - \tilde{A}^T (\delta A x - \delta b).$$

Eveneens geldt:

$$1.2 \quad \tilde{M} \tilde{x} = \tilde{A}^T \tilde{b}.$$

Uit combinatie van 1.1 en 1.2 volgt:

$$1.3 \quad x - \tilde{x} = -\tilde{M}^{-1} \delta A^T r - \tilde{M}^{-1} \tilde{A}^T (\delta A x - \delta b).$$

We weten dat  $\tilde{A}$  geschreven kan worden als het product van een  $m \times n$ -orthonormale matrix  $\tilde{Q}$ , dwz.  $\tilde{Q}^T \tilde{Q} = I$ , en een  $n \times n$  bovendriehoeks-matrix  $\tilde{R}$ . Hieruit volgt dat  $\tilde{M} = \tilde{R}^T \tilde{R}$ .

Toepassing hiervan in 1.3 levert ons:

$$1.4 \quad x - \tilde{x} = -\tilde{M}^{-1} \delta A^T r - \tilde{R}^{-1} \tilde{Q}^T (\delta A x - \delta b).$$

Dus geldt voor de norm hiervan wegens de driehoeksongelijkheid en de producteigenschap van normen (merk op dat  $\|\tilde{Q}^T\| = 1$ ):

$$1.5 \quad \|x - \tilde{x}\| \leq \|\tilde{M}^{-1}\| \|\delta A^T\| \|r\| + \|\tilde{R}^{-1}\| (\|\delta A\| \|x\| + \|\delta b\|).$$

Voor  $\|\tilde{M}^{-1}\|$  kunnen we de volgende afschatting geven:

$$\|\tilde{M}^{-1}\| \leq \|M^{-1}\| + \|\tilde{M}^{-1}\| \|\tilde{I} - \tilde{M}^{-1}\|.$$

Indien  $\|\tilde{I} - \tilde{M}^{-1}\| < 1$  is geldt dus

$$\|\tilde{M}^{-1}\| \leq \|M^{-1}\| / (1 - \|\tilde{I} - \tilde{M}^{-1}\|).$$

Hieruit volgt:

$$1.6 \quad \|\tilde{M}^{-1}\| \leq \|M^{-1}\| /$$

$$\{1 - (\|A^T\| \|\delta A\| + \|\delta A^T\| \|A\| + \|\delta A^T\| \|\delta A\|) \|M^{-1}\|\},$$

mits de noemer van deze laatste expressie positief is.

Analoog hieraan kunnen we voor  $\|\tilde{R}^{-1}\|$  vinden:

$$1.7 \quad \|\tilde{R}^{-1}\| \leq \|R^{-1}\| / \{1 - \|\tilde{R} - R\| \|R^{-1}\|\},$$

mits ook in deze laatste expressie de noemer positief is. Hieraan is zeker voldaan als de noemer in 1.6 positief is.

Tot slot willen we over de betekenis van 1.5 het volgende opmerken:

Het rechterlid van 1.5 bestaat uit 2 termen waarvan de eerste  $||\tilde{M}^{-1}||$  als factor heeft en de tweede  $||\tilde{R}^{-1}||$ . Daar  $\tilde{M} = \tilde{R}^T \tilde{R}$  geldt ruwweg  $||\tilde{M}^{-1}|| = O(||\tilde{R}^{-1}||^2)$ .

Van een goed gesteld kleinste-kwadratenprobleem mogen we echter verwachten dat het residu  $r$  een kleine norm heeft, zodat  $||r||$  helpt om de eerste term klein te krijgen.

Kwalitatief kunnen we nu zeggen dat  $||x - \tilde{x}||$  niet te groot wordt, als  $||\delta A||$  en  $||\delta b||$  klein genoeg zijn,  $||r||$  klein is en de conditie van het probleem niet al te slecht is, opdat  $||\tilde{R}^{-1}||$  en  $||\tilde{M}^{-1}|| ||r||$  niet te groot worden.

## §2. Voorwaartse foutenanalyse van de "klassieke" methode

Zij  $c = A^T b$ , dan wordt het beschouwde probleem teruggebracht tot het oplossen van het vierkante stelsel:

$$Mx = c.$$

Veronderstel dat we werken met een floating-point arithmetiek met grondtal  $\beta$  en een woordlengte van  $t$  digits.

We beschouwen de fouten in  $\bar{M}$  en  $\bar{c}$ :

$$2.1 \quad \bar{M} = \text{fl}(A^T A) = A^T A + E,$$

$$2.2 \quad \bar{c} = \text{fl}(A^T b) = A^T b + e.$$

Er geldt (vgl. [1], blz. 22):

$$2.3 \quad ||E|| \leq m \beta^{-t} ||A^T|| ||A||,$$

$$2.4 \quad ||e|| \leq m \beta^{-t} ||A^T|| ||b||.$$

Hierin is  $\beta^{-t} = 1.06C\beta^{-t}$  waarin  $C$  afhangt van de arithmetiek in de computer; indien  $\text{fl}(a \times b)$  optimaal is, geldt  $C = \beta/2$ .

Voor een zekere  $\hat{x}$  geldt:

$$2.5 \quad \bar{M} \hat{x} = \bar{c}.$$

Dit is het stelsel dat in de computer opgelost wordt; de numeriek gevonden oplossing is gelijk aan  $\bar{x}$ .

Er bestaat een matrix  $F$  waarvoor geldt:

$$2.6 \quad (\bar{M}+F)\bar{x} = \bar{c}.$$

Indien 2.5 opgelost wordt met de methode van Cholesky kunnen we  $\|F\|$  afschatten met:

$$2.7 \quad \|F\| \leq n^3 \beta^{-t} \|\bar{M}\|.$$

(Vergelijk [1], (5.5.2); doordat  $\bar{M}$  positief definitief is, weten we dat de in (5.5.2) voorkomende  $f(n)$  gelijk is aan 1).

Combinatie van 2.1, 2.2 en 2.6 levert ons:

$$(M+E+F)\bar{x} = c + e.$$

Hieruit volgt:

$$\bar{x} - x = M^{-1}\{e - (E+F)\bar{x}\},$$

of voor de norm hiervan:

$$2.8 \quad \|\bar{x} - x\| \leq \beta^{-t} \|M^{-1}\| \{m \|A^T\| \|b\| + (m \|A^T\| \|A\| + n^3 \|\bar{M}\|) \|\bar{x}\|\}.$$

Indien  $f1(\bar{M}^{-1})$  berekend wordt, kunnen we een afchatting voor  $\|M^{-1}\|$  geven:

$$\begin{aligned} M^{-1} &= f1(\bar{M}^{-1}) + M^{-1}\{I - Mf1(\bar{M}^{-1})\} \\ &= f1(\bar{M}^{-1}) + M^{-1}\{I - \bar{M}f1(\bar{M}^{-1}) + \bar{M}f1(\bar{M}^{-1}) - Mf1(\bar{M}^{-1})\}. \end{aligned}$$

We weten dat ([1], blz. 40):

$$\|\bar{M}^{-1} - f1(\bar{M}^{-1})\| \leq n^3 \beta^{-t} \|f1(\bar{M}^{-1})\|$$

en:

$$\bar{M} - M = E.$$

Hieruit volgt

$$2.9 \quad \|M^{-1}\| \leq \frac{\|f1(\bar{M}^{-1})\|}{1 - \beta^{-t} (n^3 \|\bar{M}\| + m \|A^T\| \|A\|) \|f1(\bar{M}^{-1})\|}.$$

Combinatie van 2.9 en 2.8 levert ons de gewenste afchatting van  $\|\bar{x} - x\|$ .

Laten we nu de samenhang tussen 2.8 en 1.5 bekijken. Stel daartoe

$$m\beta^{-1} \|A^T\| = \|\delta A^T\|.$$

Er geldt dan:

$$2.10 \quad \|\bar{x} - x\| \leq \|M^{-1}\| \|\delta A^T\| (\|A\| \|\bar{x}\| + \|b\|) \\ + n^3 \beta^{-1} \|M^{-1}\| \|\bar{M}\| \|\bar{x}\|.$$

We zien als verschillen:

- 1<sup>o</sup>) in plaats van  $\|\tilde{R}^{-1}\| \|\delta A\| \|x\|$  in 1.5 treedt  $\|M^{-1}\| \|\delta \bar{M}\| \|\bar{x}\|$  in 2.10 op;
- 2<sup>o</sup>) in plaats van  $\|r\|$  in 1.5 treedt  $\|A\| \|\bar{x}\| + \|b\|$  in 2.10 op;
- 3<sup>o</sup>) in plaats van de term  $\|\tilde{R}^{-1}\| \|\delta b\|$  in 1.5 treedt geen expliciete corresponderende term in 2.10 op; een equivalente bijdrage is in  $\|M^{-1}\| \|\delta A^T\| \|b\|$  verwerkt.

Van deze verschilpunten is dat onder 1<sup>o</sup> het belangrijkste.

Door dit verschil kan de bovengrens 2.10 te hoog uitvallen, ook als de gemaakte rekenfouten kleiner zijn dan  $\delta A$  en  $\delta b$  uit §1.

Al door de bepaling van  $A^T A$  kunnen ernstige fouten geïntroduceerd worden. Een remedie hiertegen is de berekening van  $A^T A$  en de Cholesky-ontbinding in dubbele-lengtearithmetiek uit te voeren.

### §3. Voorwaartse foutenanalyse van de QR-methode

Bij gegeven  $A$  bestaan een matrix  $Q$  en een matrix  $R$  waarvoor  $A = QR$  (zie §1). Het kleinste-kwadratenprobleem kunnen we reduceren tot:

$$3.1 \quad Rx = Q^T b = c.$$

De numeriek gevonden matrix  $\bar{R}$  verschilt echter van  $R$ . Wel bestaat een orthogonale matrix  $\hat{Q}$  waarvoor:

$$3.2 \quad \bar{R} = \hat{Q}^T (A+E),$$

waarin



$$3.3 \quad ||E|| \leq 12.5 n^{3/2} \beta^{-t_1} ||A||, \quad (\text{zie [2], [3]}).$$

Eveneens geldt:

$$3.4 \quad \bar{c} = \hat{Q}^T (b+e),$$

waarin

$$3.5 \quad ||e|| \leq 12.5 n \beta^{-t_1} ||b||, \quad (\text{zie [2], [3]}).$$

Zij  $\hat{x}$  de exacte kleinste-kwadratenoplossing van  $[A+E;b+e]$  en  $x$  die van  $[A;b]$  dan geldt volgens 1.4:

$$3.6 \quad x - \hat{x} = -M^{-1} E^T \hat{r} - R^{-1} Q^T (E\hat{x} - e),$$

waarbij:

$$\hat{r} = (A+E)\hat{x} - (b+e).$$

Daar we hier een foutenafschatting geven zonder residuberekening zullen we  $||\hat{r}||$  afschatten met  $||b|| + ||e||$ . Er geldt namelijk voor  $r = Ax - b$ :

$$\begin{aligned} ||r||^2 &= r^T r = x^T A^T Ax - x^T A^T b - b^T Ax + b^T b \\ &= x^T (A^T Ax - A^T b) - x^T A^T b + b^T b \\ &= b^T b - x^T A^T b \\ &= ||b||^2 - x^T A^T Ax \\ &= ||b||^2 - ||Ax||^2 \leq ||b||^2. \end{aligned}$$

Voor de norm van  $x - \hat{x}$  geldt dan:

$$3.7 \quad ||x - \hat{x}|| \leq ||M^{-1}|| ||E^T|| (||b|| + ||e||) + \\ + ||R^{-1}|| (||E|| ||\hat{x}|| + ||e||).$$

De berekende oplossingsvector  $\bar{x}$  wordt verkregen door een terugsubstitutie toe te passen op het stelsel 3.1.

Er geldt voor zekere matrix F:

$$3.8 \quad (\bar{R}+F)\bar{x} = \bar{R}\hat{x},$$

waarin (volgens [1], (5.4.3) met gebruikmaking van  $\|F\|_2 \leq \sqrt{n}\|F\|_1$  (vgl. [1], pag. 10) en  $g = \|\bar{R}\|$  wegens  $\bar{R}^T\bar{R}$  pos. def.):

$$3.9 \quad \|F\| \leq n^{5/2} \|\bar{R}\| \beta^{-t}.$$

Nu geldt:

$$3.10 \quad \|\bar{x}-\hat{x}\| \leq \|\bar{R}^{-1}\| \|F\| \|\bar{x}\|$$

en

$$3.11 \quad \|\hat{x}\| \leq (1 + \|\bar{R}^{-1}\| \|F\|) \|\bar{x}\|.$$

Uiteindelijk krijgen we dus als bovengrens voor  $\|\bar{x}-x\|$ :

$$3.12 \quad \begin{aligned} \|\bar{x}-x\| &\leq \|\bar{R}^{-1}\| \|F\| \|\bar{x}\| \\ &\quad + \|M^{-1}\| \|E^T\| (\|b\| + \|e\|) \\ &\quad + \|\bar{R}^{-1}\| \{ \|E\| (1 + \|\bar{R}^{-1}\| \|F\|) \|\bar{x}\| + \|e\| \}. \end{aligned}$$

Ter completering van 3.12 moeten we nog de volgende afschattingen geven:

$$3.13 \quad \|\bar{R}^{-1}\| \leq \frac{\|\bar{R}^{-1}\|}{1 - \|\bar{R}-R\| \|\bar{R}^{-1}\|},$$

$$3.14 \quad \|\bar{R}^{-1}\| \leq \frac{\|fl(\bar{R}^{-1})\|}{1 - n^3 \beta^{-t} \|fl(\bar{R}^{-1})\|}.$$

Voor de afschatting van  $\|M^{-1}\|$  zie 2.9.

In de praktijk geldt, als  $n^3 \beta^{-t} \|fl(\bar{R}^{-1})\| \ll 1$ , dan is ook  $\|\bar{R}-R\| \|\bar{R}^{-1}\| \ll 1$ , dus, blijkens formules 3.13 en 3.14, zijn dan  $\|\bar{R}^{-1}\|$  en  $\|\bar{R}^{-1}\|$  praktischeven groot als de door berekening verkregen grootheid  $\|fl(\bar{R}^{-1})\|$ .

#### §4. Foutenanalyse met residuberekening

Als we een vector  $\bar{x}$  berekend hebben als benadering voor de oplossingsvector van een kleinste-kwadratenprobleem, dan kunnen we de residuvector  $r = A\bar{x}-b$

berekenen. We vinden:

$$4.1 \quad \bar{r} = \text{fl}(A\bar{x}-b) = A\bar{x} - b + \delta r,$$

met:

$$4.2 \quad \|\delta r\| \leq \{(n+1) \|A\| \|\bar{x}\| + \|b\|\} \beta^{-t_1}.$$

Zij verder  $s = A^T \bar{r}$ , dan geldt:

$$4.3 \quad \bar{s} = \text{fl}(A^T \bar{r}) = A^T \bar{r} + \delta s,$$

met:

$$4.4 \quad \|\delta s\| \leq m \beta^{-t_1} \|A^T\| \|\bar{r}\|.$$

Uit 4.1 volgt:

$$4.5 \quad A^T \bar{r} = A^T A \bar{x} - A^T b + A^T \delta r.$$

Uit 4.5, 4.3 en de relatie  $A^T A x = A^T b$ , volgt:

$$4.6 \quad \bar{x} - x = M^{-1} \{\bar{s} - \delta s - A^T \delta r\},$$

dus geldt voor de norm:

$$4.7 \quad \begin{aligned} \|\bar{x} - x\| \leq & \|M^{-1}\| \{ \|\bar{s}\| + m \beta^{-t_1} \|A^T\| \|\bar{r}\| \} \\ & + \|R^{-1}\| \beta^{-t_1} \{(n+1) \|A\| \|\bar{x}\| + \|b\|\}. \end{aligned}$$

De afschatting voor  $\|M^{-1}\|$  is gegeven in 2.9, die voor  $\|R^{-1}\|$  in 3.13.

Indien  $\bar{r}$  en  $\bar{s}$  in dubbele-lengteprecisie worden berekend, geldt bij benadering:

$$4.8 \quad \|\bar{x} - x\| \leq \|M^{-1}\| \|\bar{s}\|.$$

Literatuur:

1. T.J. Dekker, Numerieke Algebra.  
MC Syllabus 12, Mathematisch Centrum, Amsterdam 1971.
2. Å. Björck, Iterative refinement of linear least squares solutions I,  
BIT 7, (1967) 257-278.
3. G.H. Golub & J.H. Wilkinson,  
Note on the iterative refinement of least squares solution,  
Num. Mat. 9 (1966) 139-148.