



ANDRZEJ KIELBASIŃSKI (Warszawa) i KRYSZYNA ZIĘTAK (Wrocław)

## Analiza numeryczna typowych zadań z unitarnym przekształceniem Householdera

(Praca przyjęta do druku 24.4.1975)

**0. Wstęp.** W pracy tej przedstawiamy dowód numerycznej poprawności unitarnego przekształcenia Householdera. Rozważamy różne arytmetyki numeryczne rzeczywiste lub zespolone oraz typowe lub ogólne warunki definiujące transformację. Pewne szczegóły tej analizy wydają się nam nowe, zasadnicza jednak jej struktura pochodzi z pracy J. H. Wilkinsona [6].

Główną zatem przyczyną naszej publikacji jest chęć spopularyzowania tej ważnej transformacji oraz metod analizy Wilkinsonowskiej.

Numeracja wzorów jest oddzielna w każdym z paragrafów. Cytując wzory podane w innym paragrafie podajemy dwie liczby: numer paragrafu i numer wzoru.

**1. Zespolone i rzeczywiste przekształcenia Householdera.** Niech  $K^n$  oznacza zespoloną lub rzeczywistą kartezyjską przestrzeń wektorów, utożsamianych z macierzami jednokolumnowymi, z normą euklidesową:

$$\|\vec{x}\| = (\vec{x}^H \vec{x})^{1/2}, \quad \vec{x} \in K^n.$$

*Unitarnym przekształceniem Householdera* nazywamy przekształcenie dające się zapisać w postaci

$$(1) \quad P = I - \vec{w} \cdot 2 \cdot \vec{w}^H, \quad \|\vec{w}\| = 1.$$

Celowo zastosowaliśmy tu zapis postaci  $\vec{w} \cdot 2 \cdot \vec{w}^H$ , gdyż jest to poprawnie określony iloczyn macierzy wymiarów  $n \times 1$ ,  $1 \times 1$ ,  $1 \times n$  (przez  $B^H$  rozumiemy macierz sprzężoną z macierzą  $B$ ).

Przekształcenie (1) przyporządkowuje wektorowi  $\vec{x} \in K^n$  jego odbicie zwierciadlane  $\vec{y}$  względem hiperpłaszczyzny prostopadłej do wektora  $\vec{w}$ . Istotnie,

$$\vec{y} = P\vec{x} = \vec{x} - \vec{r} \cdot 2,$$

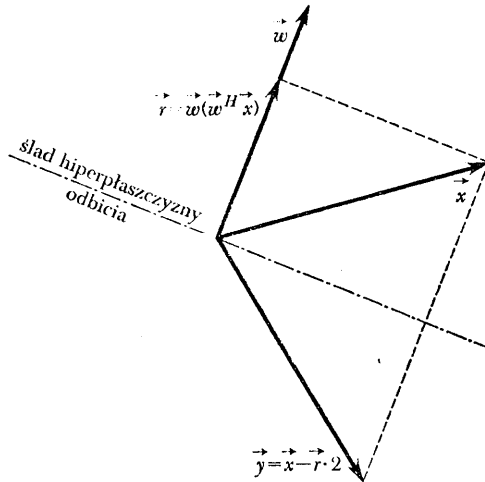
gdzie

$$\vec{r} = \vec{w}(\vec{w}^H \vec{x})$$

jest rzutem prostopadłym wektora  $\vec{x}$  na wektor  $\vec{w}$ .

Unitarne przekształcenia Householdera nazywamy więc niekiedy *przekształceniami odbić zwierciadlanych*. Z tej interpretacji geometrycznej natychmiast wynika unitarność przekształcenia  $P$  (tzn.  $\forall \vec{x} \in K^n$ ,  $\|P\vec{x}\| = \|\vec{x}\|$ ) oraz dalsze, łatwo sprawdzalne własności:

$$PP = I, \quad P = P^{-1} = P^H.$$



Rys. 1 (płaszczyzna rysunku rozpięta na wektorach  $\vec{w}$  i  $\vec{x}$ )

Możemy również sprawdzić, że  $\det(P) = -1$ . Rozważmy bowiem macierz

$$B = [\vec{w}, \vec{a}_2, \vec{a}_3, \dots, \vec{a}_n],$$

gdzie  $\vec{a}_2, \vec{a}_3, \dots, \vec{a}_n$  są dowolnymi liniowo niezależnymi wektorami, rozpinającymi hiperpłaszczyznę prostopadłą do wektora  $\vec{w}$ . Wówczas  $\det(B) \neq 0$ .

Zbadajmy wyznacznik macierzy

$$PB = [P\vec{w}, P\vec{a}_2, \dots, P\vec{a}_n].$$

Ponieważ  $P\vec{w} = -\vec{w}$ ,  $P\vec{a}_i = \vec{a}_i$ , więc

$$\det(PB) = \det[-\vec{w}, \vec{a}_2, \dots, \vec{a}_n] = -\det(B).$$

Wobec tego, że  $\det(PB) = \det(P) \cdot \det(B)$ , otrzymujemy  $\det(P) = -1$ . Przekształcenie Householdera może być oczywiście zapisane również za pomocą dowolnego wektora niezerowego  $\vec{u}$ :  $P = I - \frac{2}{\|\vec{u}\|^2} \vec{u} \vec{u}^H$ , gdy  $\vec{w} = \vec{u} / \|\vec{u}\| \cdot c$ ,  $|c| = 1$ ,  $c \in K$ .

**2. Podstawowe zadanie (H) z przekształceniem Householdera.** Rozważmy następujące zadanie:

Dla danej pary niezerowych wektorów zespolonych  $\vec{a}$  i  $\vec{e}$  należących do  $K^n$  chcemy wyznaczyć przekształcenie unitarne Householdera, które przeprowadza wektor  $\vec{a}$  na wektor  $P\vec{a}$  o kierunku wektora  $\vec{e}$ , tzn.

$$(1) \quad P\vec{a} = \vec{e}k, \quad |k| = \|\vec{a}\| / \|\vec{e}\|.$$

Ponieważ przekształcenie  $P$  jest hermitowskie, więc wyrażenie  $\vec{a}^H P \vec{a} = \vec{a}^H \vec{e} k$  musi być rzeczywiste. Okazuje się (zob. [6], str. 59), że jest to zarazem warunek dostateczny istnienia przekształcenia spełniającego (1.1) i (1). Zatem należy przyjąć

$$(2) \quad k = \begin{cases} \pm \frac{\vec{e}^H \vec{a}}{|\vec{e}^H \vec{a}|} \frac{\|\vec{a}\|}{\|\vec{e}\|}, & \text{gdy } \vec{e}^H \vec{a} \neq 0, \\ \pm \frac{\|\vec{a}\|}{\|\vec{e}\|}, & \text{gdy } \vec{e}^H \vec{a} = 0. \end{cases}$$

Z wzoru (1) wynika

$$\vec{w} = (\vec{a} - \vec{e}k) / \|\vec{a} - \vec{e}k\|.$$

Niech

$$(3) \quad \vec{u} = \vec{a} - \vec{e}k,$$

$$(4) \quad R = \|\vec{u}\|^2/2.$$

Wówczas przekształcenie (1.1) możemy zapisać w postaci

$$(5) \quad P = I - \vec{u} \frac{1}{R} \vec{u}^H.$$

Parametr  $k$  określiliśmy z dokładnością do znaku. Znak ten wybieramy tak, by wektor  $\vec{u}$  miał największą długość. Łatwo sprawdzić, że

$$(6) \quad R = \|\vec{a}\|^2 \mp |\vec{e}^H \vec{a}| \|\vec{a}\| / \|\vec{e}\|.$$

Wobec tego we wzorach (2) wybieramy znak minus. A więc z (2), (3), (6) mamy

$$(7) \quad \vec{u} = \begin{cases} \vec{a} + \vec{e} \frac{\vec{e}^H \vec{a}}{|\vec{e}^H \vec{a}|} \frac{\|\vec{a}\|}{\|\vec{e}\|}, & \text{gdy } \vec{e}^H \vec{a} \neq 0, \\ \vec{a} + \vec{e} \frac{\|\vec{a}\|}{\|\vec{e}\|}, & \text{gdy } \vec{e}^H \vec{a} = 0, \end{cases}$$

$$(8) \quad R = \|\vec{a}\|^2 + |\vec{e}^H \vec{a}| \|\vec{a}\| / \|\vec{e}\|.$$

Taki wybór znaku uzasadnimy później. Na razie ograniczymy się do uwagi, że zapewnia on numeryczną poprawność (zob. [3]) najważniejszych algorytmów korzystających z przekształcenia Householdera.

Z (2), (4), (6) otrzymujemy dla przypadku rzeczywistych wektorów  $\vec{a}$ ,  $\vec{e}$  następujące wzory na rzeczywiste przekształcenie Householdera:

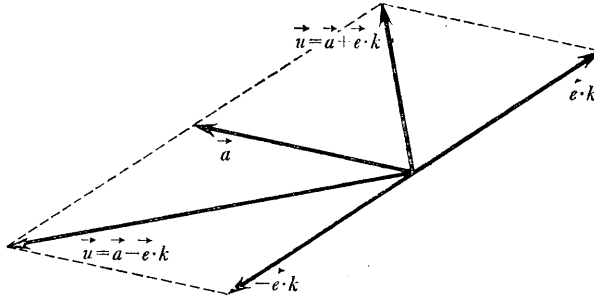
$$(9) \quad \vec{u} = \begin{cases} \vec{a} + \vec{e} \operatorname{sign}(\vec{e}^T \vec{a}) \|\vec{a}\| / \|\vec{e}\|, & \text{gdy } \vec{e}^T \vec{a} \neq 0, \\ \vec{a} + \vec{e} \|\vec{a}\| / \|\vec{e}\|, & \text{gdy } \vec{e}^T \vec{a} = 0, \end{cases}$$

$$(10) \quad R = \|\vec{a}\|^2 + |\vec{e}^T \vec{a}| \|\vec{a}\| / \|\vec{e}\|.$$

Konstrukcja rzeczywistego przekształcenia Householdera w rozważanym zadaniu ma prostą interpretację geometryczną. Mianowicie, wektor  $\vec{u}$  leży na dwusiecznej kąta między wektorami  $\vec{a}$  i  $-\vec{e}$  lub  $\vec{a}$  i  $\vec{e}$  (patrz rys. 2).

Dla powiększenia dokładności przekształcenia, wybieramy dwusieczną mniejszego kąta (spośród kątów między wektorami  $\vec{a}$  i  $\vec{e}$  oraz  $\vec{a}$  i  $-\vec{e}$ ). Wówczas kierunek odpowiedniej dwusiecznej jest „lepiej wyznaczony”.

Podstawowym zadaniem korzystającym z przekształcenia Householdera jest wyznaczenie wektora  $\vec{c} = P\vec{b}$  dla danego wektora  $\vec{b}$  i danych wektorów  $\vec{a}$ ,  $\vec{e}$ , defi-



Rys. 2 (płaszczyzna rysunku rozpięta na wektorach  $\vec{a}$  i  $\vec{e}$ )

niujących przekształcenie  $P$ . Zadanie to nazywać będziemy *zadaniem* (H). W dalszej części pracy wykażemy, że wyznaczenie wektora  $\vec{c}$  może być wykonane w sposób numerycznie poprawny, gdy przekształcenie  $P$  jest określone zależnościami (5), (7), (8). Chcąc wyznaczyć obraz  $\vec{c}$  dowolnego wektora zespolonego  $\vec{b}$ , nie musimy w sposób efektywny konstruować macierzy przekształcenia (5), gdyż

$$(11) \quad \vec{c} = P\vec{b} = \vec{b} - \vec{u} \left( \frac{1}{R} \vec{u}^H \vec{b} \right).$$

Wygodnie jest zatem reprezentować przekształcenie  $P$  współzrędnymi wektora  $\vec{u}$  i liczbą  $R$ . Reprezentacja  $P$  zajmuje wtedy tylko około  $n$  miejsc (liczb rzeczywistych lub zespolonych) w pamięci maszyny (macierz  $P$  zajęłaby około  $n^2/2$  miejsc) zaś koszt otrzymania  $\vec{c}$  ze wzoru (11) jest rzędu tylko  $2n$  mnożeń. (Zamiast  $n^2$  przy posługiwaniu się pełną macierzą  $P$ .)

**3. Zespolona arytmetyka zmiennopozycyjna**  $cf_1$ ,  $cf_{1_2}$ ,  $cf_1$ . Zakładamy, że dysponujemy dwójkową  $t$ -cyfrową arytmetyką zmiennopozycyjną  $fl$  (zob. np. [5], str. 16, [6], str. 110) określoną tak, że obliczone wyniki operacji arytmetycznych i pierwiastkowania mogą być zapisane w postaci:

$$fl(x \diamond y) = (x \diamond y)(1 - \varepsilon), \quad |\varepsilon| \leq 2^{-t},$$

$$fl(\sqrt{x}) = \sqrt{x} (1 - \delta), \quad |\delta| \leq p \cdot 2^{-t},$$

gdzie  $\diamond$  oznacza dowolny z operatorów  $\pm$ ,  $\times$ ,  $/$ , liczba  $p$  zaś należy na ogół do przedziału  $[1, 2]$ . Będziemy stosować relację  $\leq$  (wprowadzoną w [2] lub w [4]), co pozwala na pomijanie w oszacowaniach błędów mniej istotnych składników rzędu  $2^{-2t}$ ,  $2^{-3t}$ , itd. w obecności składników rzędu  $2^{-t}$ .

Niech  $z_k = x_k + i\bar{y}_k$  ( $x_k, y_k$  rzeczywiste,  $i$  — jednostka urojona). Wówczas dodawanie w arytmetyce cfl jest wykonane zgodnie z regułą

$$\begin{aligned} \text{cfl}(z_1 + z_2) &= \text{fl}(x_1 + x_2) + i\bar{\text{fl}}(y_1 + y_2) = \\ &= (x_1 + x_2)(1 - \varepsilon_1) + i\bar{(y_1 + y_2)}(1 - \varepsilon_2) = (z_1 + z_2)(1 - \varepsilon_3), \end{aligned}$$

gdzie  $|\varepsilon_j| \leq 2^{-t}$  ( $\varepsilon_1, \varepsilon_2$  — rzeczywiste,  $\varepsilon_3$  na ogół zespolona).

Podobnie wykonywane są w cfl pozostałe podstawowe operacje. Mianowicie, łatwo pokazać, że

$$(1) \quad \begin{aligned} \text{cfl}(z_1 z_2) &= (z_1 z_2)(1 - \varepsilon), & |\varepsilon| &\leq \frac{1}{1} (1 + \sqrt{2}) 2^{-t}, \\ \text{cfl}(z_1 / z_2) &= (z_1 / z_2)(1 - \varphi), & |\varphi| &\leq \frac{1}{1} (4 + \sqrt{2}) 2^{-t}, \\ \text{cfl}(|z_1|) &= |z_1|(1 - \nu), & |\nu| &\leq \frac{1}{1} (p + 2,25) 2^{-t}, \\ \text{cfl}(rz_1) &= (rz_1)(1 - \delta), & |\delta| &\leq 2^{-t}, \end{aligned}$$

gdzie  $r$  jest liczbą rzeczywistą, a  $|z_1|, z_1/z_2$  obliczamy według następujących wzorów (por. [7]):

$$|z_1| = \begin{cases} 0, & \text{gdy } x_1 = x_2 = 0, \\ |x_1| \sqrt{1 + (y_1/x_1)^2}, & \text{gdy } |x_1| \geq |y_1|, \\ |y_1| \sqrt{1 + (x_1/y_1)^2}, & \text{gdy } |x_1| < |y_1|, \end{cases}$$

$$z = \frac{z_1}{z_2} = \begin{cases} \frac{x_1 + y_1 h}{x_2 + y_2 h} + i \frac{y_1 - x_1 h}{x_2 + y_2 h}, & h = \frac{y_2}{x_2}, \quad \text{gdy } |x_2| \geq |y_2|, \\ \frac{x_1 h + y_1}{x_2 h + y_2} + i \frac{y_1 h - x_1}{x_2 h + y_2}, & h = \frac{x_2}{y_2}, \quad \text{gdy } |x_2| < |y_2|. \end{cases}$$

Algorytmy te są „bezpieczne”. Pozwalają na uniknięcie niedomiaru lub nadmiaru, o ile dane i wyniki są dobrze reprezentowane w cfl. Dla przykładu podamy analizę numeryczną algorytmu dzielenia. Ze względu na „symetrię” algorytmu wystarczy zbadać jeden przypadek. Np.  $|x_2| \geq |y_2|$ . Niech

$$h = y_2/x_2, \quad m = x_2 + y_2 h, \quad u = \frac{x_1 + y_1 h}{m}, \quad v = \frac{y_1 - x_1 h}{m}.$$

Wówczas

$$(2) \quad \bar{z} = \text{cfl}\left(\frac{z_1}{z_2}\right) = \frac{x_1 + y_1 h(1 - \delta_1)}{x_2 + y_2 h(1 - \delta_2)}(1 - \delta_3) + i \frac{y_1 - x_1 h(1 - \delta_4)}{x_2 + y_2 h(1 - \delta_2)}(1 - \delta_5),$$

gdzie

$$|\delta_1|, |\delta_2|, |\delta_4| \leq \frac{1}{1} 2 \cdot 2^{-t}, \quad |\delta_3|, |\delta_5| \leq \frac{1}{1} 3 \cdot 2^{-t}.$$

Ponieważ  $\text{sign}(x_2) = \text{sign}(y_2 h)$ , a ponadto  $|y_2 h/m| \leq \frac{1}{2}$ , więc

$$x_2 + y_2 h(1 - \delta_2) = m \left(1 - \frac{y_2 h}{m} \delta_2\right) = m(1 - \varphi_2), \quad |\varphi_2| \leq \frac{1}{1} 2^{-t}.$$

Zatem z (2) mamy

$$\tilde{z} = \frac{1}{m} \{ [x_1 + y_1 h(1 - \delta_1)](1 - \varphi_3) + i [y_1 - x_1 h(1 - \delta_4)](1 - \varphi_5) \}, \quad |\varphi_i| \leq 4 \cdot 2^{-t}.$$

(3)

Niech

$$w = \frac{1}{m} [(x_1 + y_1 h)(1 - \varphi_3) + i (y_1 - x_1 h)(1 - \varphi_5)].$$

Wówczas z (3) otrzymujemy

$$\frac{|z - w|^2}{|z|^2} = \frac{(x_1 + y_1 h)^2 \varphi_3^2 + (y_1 - x_1 h)^2 \varphi_5^2}{(x_1 + y_1 h)^2 + (y_1 - x_1 h)^2} \leq 16 \cdot 2^{-2t},$$

$$\frac{|\tilde{z} - w|^2}{|z|^2} = \frac{y_1^2 h^2 \delta_1^2 (1 - \varphi_3)^2 + x_1^2 h^2 \delta_4^2 (1 - \varphi_5)^2}{(x_1 + y_1 h)^2 + (y_1 - x_1 h)^2} \leq \frac{h^2}{1 + h^2} 4 \cdot 2^{-2t} \leq 2 \cdot 2^{-2t}.$$

Stąd, wykorzystując nierówność  $|\tilde{z} - z| \leq |\tilde{z} - w| + |w - z|$ , otrzymujemy

$$\frac{|\tilde{z} - z|}{|z|} \leq (4 + \sqrt{2}) 2^{-t},$$

co należało pokazać.

Zbadamy teraz błąd wytworzony przy obliczaniu w arytmetyce cfl iloczynu skalarnego dwóch wektorów  $\vec{w}$  i  $\vec{z}$  o składowych  $w_j, z_j$ . Niech

$$(4) \quad s = \sum_{j=1}^n z_j \bar{w}_j.$$

Łatwo pokazać, że (por. [5], str. 30 i następne)

$$(5) \quad \tilde{s} = \text{cfl} \left( \sum_{j=1}^n z_j \bar{w}_j \right) = \left[ \sum_{j=1}^n z_j \bar{w}_j (1 - \xi_j) \right] (1 - \xi),$$

$$|\xi_j| \leq (n - j + 1 + \sqrt{2}) 2^{-t}, \quad |\xi| \leq 2^{-t}.$$

A więc obliczony wynik  $\tilde{s}$  jest zaokrąglonym iloczynem skalarnym nieco zniekształconych wektorów  $\vec{w}$  lub  $\vec{z}$ , gdyż czynnik  $(1 - \xi_j)$  możemy uważać za zaburzenie współrzędnej  $w_j$  lub  $z_j$ . Dla dowolnych wektorów  $\vec{w}$  i  $\vec{z}$  nie możemy twierdzić, że iloczyn skalarny jest obliczony z małym błędem względnym. Dużą dokładność względną obliczania  $\vec{w}^H \vec{z}$  otrzymujemy w szczególnych przypadkach, np. gdy  $\vec{w} = \vec{z}$ . Mianowicie,

$$(6) \quad \text{cfl}(\vec{z}^H \vec{z}) = \text{fl} \left( \sum_{j=1}^n [(\text{re } z_j)^2 + (\text{im } z_j)^2] \right) = (\vec{z}^H \vec{z})(1 - \varphi),$$

$$|\varphi| \leq (n + 1) 2^{-t}.$$

Zajmiemy się teraz arytmetyką  $\text{cfl}_2$  i  $\text{cfl}$ . Zakładamy, że czytelnik zna zasady arytmetyki  $\text{fl}_2$  (zob. [5], str. 37) i  $\text{fl}$  (zob. [1]). Niech

$$u = \sum_{j=1}^n u_j = \sum_{j=1}^n (a_j + ib_j).$$

Wówczas w arytmetyce  $\text{cfl}_2$  mamy

$$(7) \quad \tilde{u} = \text{cfl}_2 \left( \sum_{j=1}^n (a_j + ib_j) \right) = \text{fl}_2 \left( \sum_{j=1}^n a_j \right) + i \text{fl}_2 \left( \sum_{j=1}^n b_j \right) = \left[ \sum_{j=1}^n z_j (1 - \psi_j) \right] (1 - \psi),$$

$$|\psi_j| \leq \frac{3}{2} (n - j + 1) 2^{-2t}, \quad |\psi| \leq 2^{-t}.$$

Dla arytmetyki  $\text{cfl}$  otrzymujemy

$$(8) \quad \tilde{u} = \text{cfl} \left( \sum_{j=1}^n (a_j + ib_j) \right) = \left[ \sum_{j=1}^n z_j (1 - \psi_j) \right] (1 - \psi_0), \quad |\psi_j| \leq 2^{-t}.$$

Ostatecznie, łatwo pokazać, że dla iloczynu (4) obliczonego w arytmetyce  $\text{cfl}_2$  i  $\text{cfl}$  otrzymujemy zależności (por. (5), (7), (8)):

$$(9) \quad \tilde{s} = \text{cfl}_2 \left( \sum_{j=1}^n z_j \bar{w}_j \right) = \left( \sum_{j=1}^n z_j \bar{w}_j (1 - \varphi_j) \right) (1 - \varphi),$$

$$|\varphi_j| \leq \frac{3}{2} (n - j + 1 + \sqrt{2}) 2^{-2t}, \quad |\varphi| \leq 2^{-t},$$

$$(10) \quad \tilde{s} = \text{cfl} \left( \sum_{j=1}^n z_j \bar{w}_j \right) = \left[ \sum_{j=1}^n z_j \bar{w}_j (1 - \psi_j) \right] (1 - \psi),$$

$$|\psi_j| \leq (2 + \sqrt{2}) 2^{-t}, \quad |\psi| \leq 2^{-t}.$$

Natomiast, jeśli  $\tilde{z} = \bar{w}$ , to mamy (por. (6))

$$(11) \quad \text{cfl}_2(\tilde{z}^H \tilde{z}) = (\tilde{z}^H \tilde{z}) (1 - \sigma), \quad |\sigma| \leq 2^{-t},$$

$$(12) \quad \text{cfl}(\tilde{z}^H \tilde{z}) = (\tilde{z}^H \tilde{z}) (1 - \psi), \quad |\psi| \leq 3 \cdot 2^{-t}.$$

Wobec tego kwadrat normy wektora  $\tilde{z}$ , a więc również jego norma, jest obliczany w arytmetykach  $\text{cfl}_2$ ,  $\text{cfl}$  z dużą dokładnością względną, niezależną od wymiaru  $n$  przestrzeni.

Pojęcie arytmetyki  $\text{fl}_2$ , a więc i  $\text{cfl}_2$ , nie jest ściśle określone. Najistotniejszą cechą arytmetyk objętych tą nazwą jest opisana powyżej zdolność kumulowania iloczynów skalarnych na rejestrze podwójnej precyzji. W niektórych emc jest to jednak kosztowne, więc stosuje się kumulację jedynie w przypadku sum (iloczynów skalarnych) wielkiej liczby składników. Jeśli dostęp do rejestru podwójnej precyzji nie jest kosztowny, to możemy wykorzystywać go również do kumulacji nawet sum kilku składników, do operacji dzielenia z dzielną daną w podwójnej precyzji oraz do pierwiastkowania argumentu danego w podwójnej precyzji (por. [5], str. 38). Będziemy w takim przypadku mówili o *pełnej arytmetyce*  $\text{fl}_2$  ( $\text{cfl}_2$ ).

4. Zespolone zadanie (H) dla  $\vec{e} = \vec{e}_1$ . Niech  $a_j, b_j, u_j$  oznaczają składowe wektorów  $\vec{a}, \vec{b}, \vec{u}$ . Zakładamy, że składowe  $\vec{a}$  i  $\vec{b}$  nie są obciążone żadnym błędem i są równe swoim reprezentacjom w arytmetyce cfl. Niech

$$\vec{e} = \vec{e}_1 = (1, 0, 0, \dots, 0)^T.$$

Jest to przypadek o szczególnie ważnych zastosowaniach. Wzory (2.7), (2.8) można obecnie przekształcić do postaci (zob. [6], str. 308)

$$R = A + S, \quad u_1 = a_1 \left(1 + \frac{A}{S}\right), \quad u_j = a_j \quad (j = 2, 3, \dots, n),$$

gdzie

$$A = \sum_{j=1}^n |a_j|^2, \quad S = \sqrt{|a_1|^2 A}.$$

Wyznaczając  $R$  i  $u_1$  z powyższych wzorów (zamiast z zależności

$$R = \|\vec{a}\|^2 + |a_1| \|\vec{a}\|, \quad u_1 = a_1 \left(1 + \frac{\|\vec{a}\|}{|a_1|}\right)$$

wykonujemy o jedno pierwiastkowanie mniej.

Postaramy się podać teraz realistyczne oszacowania błędów wytworzonych przy numerycznej realizacji zespolonego zadania (H). Przez  $\tilde{F}$  będziemy oznaczać numerycznie obliczoną wartość wyrażenia  $F$  (w pseudo-algolu możemy więc zapisać  $\tilde{F} := F$ ), a  $\hat{\varphi}2^{-t}$  oznacza oszacowanie błędu  $\varphi$  (tzn.  $|\varphi| \leq \hat{\varphi}2^{-t}$ ).

Niech  $\alpha, \beta, \varrho, \mu$  oznaczają błędy względne, z jakimi są wyznaczone w danej arytmetyce wielkości  $\tilde{A}, \tilde{S}, \tilde{R}, \tilde{u}_1$ , tzn.

$$(1) \quad \tilde{A} = A(1 - \alpha), \quad \tilde{S} = S(1 - \sigma), \quad \tilde{R} = R(1 - \varrho), \quad \tilde{u}_1 = u_1(1 - \mu),$$

$$\vec{v} = [\tilde{u}_1, u_2, \dots, u_n].$$

Wektor  $\vec{v}$  jest identyczny z numerycznie wyznaczonym wektorem  $\vec{u}$ . Zbadajmy związki zachodzące pomiędzy błędami  $\alpha, \beta, \varrho, \mu$ . Pozwoli nam to lepiej oszacować normę macierzy błędów

$$(2) \quad B = \tilde{P} - P, \quad \tilde{P} = I - \vec{v} \frac{1}{R} \vec{v}^H$$

niż gdybyśmy próbowali od razu korzystać ze związków pomiędzy oszacowaniami  $\hat{\alpha}, \hat{\sigma}, \hat{\varrho}, \hat{\mu}$ . Łatwo sprawdzić, że

$$(3) \quad \begin{aligned} \sigma &= \frac{1}{2}\alpha + \varepsilon_1 + \Delta_1, & |\varepsilon_1| &\leq (1,5 + p)2^{-t}, \\ \mu &= \frac{\frac{1}{2}\alpha - \varepsilon_1 + \varepsilon_2}{1 + h} + \varepsilon_3 + \Delta_2, & |\varepsilon_2| &\leq 2^{-t}, & |\varepsilon_3| &\leq 2 \cdot 2^{-t}, \\ \varrho &= \frac{\alpha + (\frac{1}{2}\alpha + \varepsilon_1)h}{1 + h} + \varepsilon_4 + \Delta_3, & |\varepsilon_4| &\leq 2^{-t}, \end{aligned}$$



gdzie

$$h = \frac{|a_1|}{\|\bar{a}\|},$$

a błędy  $\Delta_1, \Delta_2, \Delta_3$  mają oszacowania na poziomie  $2^{-2t}$ . Błędy  $\alpha, \varepsilon_1, \varepsilon_2, \varepsilon_4, \sigma, \varrho$  są w istocie wyrażeniami rzeczywistymi, chociaż rozważamy tutaj zespolony przypadek zadania (H).

Zbadamy teraz, jak zaburzenia wektora  $\bar{u}$  i liczby  $R$  wpływają na przekształcenie  $P$ . Z (1) i (2) łatwo otrzymamy następujące wyrażenie dla  $B$ :

$$(4) \quad B = \frac{1}{R} [-\bar{u}\bar{u}^H \varrho_1 + \bar{u}\bar{e}_1^T \bar{u}_1 \bar{\mu} + \bar{e}_1 \bar{u}^H u_1 \mu - \bar{e}_1 \bar{e}_1^T u_1 \bar{u}_1 \mu \bar{\mu}],$$

gdzie

$$\varrho_1 = \varrho + \varrho^2 + \varrho^3 + \dots = \varrho + \Delta_4, \quad |\Delta_4| \leq 2^{-2t},$$

$\bar{u}_1, \bar{\mu}$  są to liczby zespolone sprzężone z liczbami  $u_1, \mu$ . Macierz  $B$  jest więc hermitowska i ma następującą budowę:

$$(5) \quad B = \frac{1}{R} \begin{pmatrix} (\mu + \bar{\mu} - \mu \bar{\mu} - \varrho_1) u_1 \bar{u}_1 & (\mu - \varrho_1) u_1 \bar{u}_2 & (\mu - \varrho_1) u_1 \bar{u}_3 & \dots & (\mu - \varrho_1) u_1 \bar{u}_n \\ (\bar{\mu} - \varrho_1) u_2 \bar{u}_1 & -\varrho_1 u_2 \bar{u}_2 & -\varrho_1 u_2 \bar{u}_3 & \dots & -\varrho_1 u_2 \bar{u}_n \\ (\bar{\mu} - \varrho_1) u_3 \bar{u}_1 & -\varrho_1 u_3 \bar{u}_2 & -\varrho_1 u_3 \bar{u}_3 & \dots & -\varrho_1 u_3 \bar{u}_n \\ \dots & \dots & \dots & \dots & \dots \\ (\bar{\mu} - \varrho_1) u_n \bar{u}_1 & -\varrho_1 u_n \bar{u}_2 & -\varrho_1 u_n \bar{u}_3 & \dots & -\varrho_1 u_n \bar{u}_n \end{pmatrix}.$$

Skorzystamy teraz z tej szczególnej postaci  $B$ . Zauważmy, że pomijając elementy pierwszego wiersza, wszystkie kolumny tej macierzy są równoległe do wektora  $(0, u_2, u_3, \dots, u_n)^T$ . Jeśli więc pomnożymy lewostronnie macierz  $B$  przez macierz odbicia zwierciadlanego  $G$ , przeprowadzającego  $(n-1)$ -wymiarowy wektor  $[u_2, u_3, \dots, u_n]^T$  na wektor o kierunku wektora  $\bar{e}_2$ , to otrzymamy macierz  $GB$ , mającą niezmienny pierwszy wiersz, zaś wiersz drugi postaci

$$\frac{1}{R} (U \bar{u}_1 (\mu - \varrho_1) \quad -U \bar{u}_2 \varrho_1 \quad \dots \quad -U \bar{u}_n \varrho_1),$$

gdzie

$$U = \left( \sum_{j=2}^n |u_j|^2 \right)^{1/2}.$$

Następne wiersze zawierają wyłącznie elementy zerowe. Zatem wszystkie wiersze macierzy  $GB$ , z pominięciem elementów pierwszej kolumny, są bądź równoległe do wektora  $(0, \bar{u}_2, \bar{u}_3, \dots, \bar{u}_n)^T$ , bądź składają się z samych zer. Zatem macierz  $GBG^H$  ma wszystkie elementy równe zeru z wyjątkiem czterech elementów w lewym górnym rogu. Ten niezerowy blok ma postać

$$(6) \quad C = \frac{1}{R} \begin{pmatrix} (\mu + \bar{\mu} - \mu \bar{\mu} - \varrho_1) u_1 \bar{u}_1 & (\mu - \varrho_1) u_1 U \\ (\bar{\mu} - \varrho_1) \bar{u}_1 U & -\varrho_1 U^2 \end{pmatrix}.$$

Oczywiście norma spektralna macierzy  $B$  jest równa normie spektralnej macierzy  $C$ . Aby wyznaczyć tę normę, musimy znaleźć największy moduł wartości własnej macierzy  $C$ . Wielomian charakterystyczny ma postać

$$(7) \quad \lambda^2 + \lambda[2\varrho_1 - (1+h)(\mu + \bar{\mu} - \mu\bar{\mu})] - (1-h^2)(1-\varrho_1)\mu\bar{\mu} = 0.$$

Widzimy więc (zob. (3)), że maksymalny moduł wartości własnej jest równy

$$|\lambda|_{\max} = \frac{1}{2} \left\{ |2\varrho_1 - (1+h)(\mu + \bar{\mu} - \mu\bar{\mu})| + \sqrt{[2\varrho_1 - (1+h)(\mu + \bar{\mu} - \mu\bar{\mu})]^2 + 4(1-h^2)(1-\varrho_1)\mu\bar{\mu}} \right\}.$$

Skorzystamy teraz ze związków (3)

$$(8) \quad |\lambda| = \left| \frac{\frac{1}{2}\alpha + \varepsilon_1 h}{1+h} + \varepsilon_1 - \varepsilon_2 + \varepsilon_4 - (1+h) \frac{\varepsilon_3 + \bar{\varepsilon}_3 + \Delta_5}{2} \right| + \\ + \sqrt{\left( \frac{\frac{1}{2}\alpha + \varepsilon_1 h}{1+h} + \varepsilon_1 - \varepsilon_2 + \varepsilon_4 - (1+h) \frac{\varepsilon_3 + \bar{\varepsilon}_3 + \Delta_5}{2} \right)^2 + \\ + (1-h)^2 \left| \frac{\frac{1}{2}\alpha + \varepsilon_2 - \varepsilon_1}{1+h} + \varepsilon_3 + \Delta_6 \right|^2},$$

gdzie  $\Delta_i$  są rzędu wielkości  $2^{-2i}$ .

Oznaczamy oszacowanie prawej strony powyższej równości przez  $\hat{\gamma}2^{-t}$ . Zatem dla  $\hat{\alpha}$  (mnożnik z oszacowania błędu  $\alpha$ ) i  $\hat{\gamma}$  otrzymujemy dla  $p = 1$  (zob. (3.6), (3.11), (3.12), (3)), wartości przedstawione w poniższej tabelce:

	cf1	cf $\bar{1}$	cf1 $_2$
$\hat{\alpha}$	$n+1$	3	1
$\hat{\gamma}$	$\frac{1}{2}n+9 + \sqrt{(\frac{1}{2}n+9)^2 + (\frac{1}{2}n+6)^2} \leq \\ \leq (1+\sqrt{2})(\frac{1}{2}n+9)$	24	23,5

Z powyższych rozważań mamy oszacowanie postaci

$$(6) \quad \|P - \tilde{P}\|_2 = \|B\|_2 \leq 2^{-t}\hat{\gamma}.$$

Niech  $\vec{g}$  oznacza numerycznie wyznaczony obraz wektora  $\vec{b}$ . Porównajmy go z wektorem  $\vec{c} = P\vec{b}$ :

$$(10) \quad \|\vec{g} - \vec{c}\| \leq \|\vec{g} - \tilde{P}\vec{b}\| + \|\tilde{P}\vec{b} - P\vec{b}\| \leq \|\vec{g} - \tilde{P}\vec{b}\| + 2^{-t}\hat{\gamma}\|\vec{b}\|.$$

Oszacujemy teraz pierwszy składnik. Wektor  $\vec{g}$  wyznaczamy w dwóch krokach (zob. (2.11))

$$\vec{w} := \frac{\vec{v}^H \vec{b}}{\tilde{R}}, \quad \vec{g} := \vec{b} - \vec{v}\vec{w}.$$

Wobec tego z (3.1), (3.5), (3.9), (3.10) otrzymujemy

$$(11) \quad \vec{w} = \frac{\vec{v}^H (I - D_1) \vec{b}}{\tilde{R}}, \quad \vec{g} = (I - D_2)(\vec{b} - (I - D_3)\vec{v}\vec{w}),$$

gdzie

$$D_1 = \text{diag}(\delta_j), \quad D_2 = \text{diag}(\beta_j), \quad D_3 = \text{diag}(\alpha_j),$$

a błędy  $\delta_j, \beta_j, \alpha_j$ , mają oszacowania z mnożnikami  $\hat{\delta}_j, \hat{\beta}_j, \hat{\alpha}_j$  odpowiednio równymi:

	cfl	cfl̄	cfl <sub>2</sub>
$\hat{\delta}_j$	$n-j+3+\sqrt{2}$	$4+\sqrt{2}$	2

$\hat{\alpha}_j = 1 + \sqrt{2}$ ,  $\hat{\beta}_j = 1$  dla każdej z rozważanych arytmetyk. Zatem stąd oraz z (2), (3) mamy następujące oszacowanie:

$$(12) \quad \|\vec{g} - \tilde{P}\vec{b}\|_1 \leq \|D_2\| \|\vec{b}\| + \frac{\|\vec{v}\vec{v}^H\|}{\tilde{R}} (\|D_1\| + \|D_2\| + \|D_3\|) \|\vec{b}\| \leq_1 \\ \leq_1 (2\hat{\delta}_1 + 3\hat{\beta}_1 + 2\hat{\alpha}_1) \|\vec{b}\| 2^{-t} \equiv \hat{\delta} \|\vec{b}\| 2^{-t}.$$

Ostatecznie więc z (10), (12) wynika, że

$$\|\vec{g} - \vec{c}\|_1 \leq (\hat{\gamma} + \hat{\delta}) \|\vec{b}\| 2^{-t} \equiv K_w \|\vec{c}\| 2^{-t},$$

przy czym (przyjeliśmy  $p = 1$ )

	cfl	cfl̄	cfl <sub>2</sub>
$\hat{\delta}$	$2n+11+3\sqrt{2}$	$13+4\sqrt{2}$	$9+2\sqrt{2}$
$K_w$	$3,2n+36$	42,6	35

Widzimy więc, że sądząc z oszacowań błędu, arytmetyka cfl̄ jest prawie równie dobra, jak arytmetyka cfl<sub>2</sub>.

W obu arytmetykach wektor  $\vec{c}$  jest wyznaczany z dużą dokładnością względną, niezależną od wymiaru  $n$ .

**5. Rzeczywiste zadanie H dla  $\vec{e} = \vec{e}_1$ .** Jeśli  $\vec{e} = \vec{e}_1$ , to wzory (2.9), (2.10) przyjmują postać

$$R = A + |a_1|N, \quad u_1 = \begin{cases} a_1 + N, & \text{gdy } a_1 \geq 0, \\ a_1 - N, & \text{gdy } a_1 < 0, \end{cases} \quad u_j = a_j \quad (j = 2, 3, \dots, n),$$

gdzie

$$A = \|\vec{a}\|^2, \quad N = \|\vec{a}\|.$$

J. H. Wilkinson przedstawia w [6], str. 148 i następną, analizę błędu wytworzonego przy numerycznej realizacji w pełnej arytmetyce fl<sub>2</sub> rzeczywistego zadania H dla  $\vec{e} = \vec{e}_1$ . W niniejszej pracy podamy oszacowania potencjalnie nieco lepsze od oszacowania otrzymanego przez Wilkinsona (zob. również [1]). Nadal stosujemy oznaczenia wprowadzone w § 4. Niech

$$\tilde{A} = A(1 - \alpha), \quad \tilde{N} = N(1 - \nu), \quad \tilde{u}_1 = u_1(1 - \mu), \quad \tilde{R} = R(1 - \varrho).$$

Wprowadzimy obecnie zależności pomiędzy błędami  $\alpha$ ,  $\nu$ ,  $\varrho$ ,  $\mu$ . Łatwo sprawdzić, że

$$1 - \nu = \sqrt{1 - \alpha} (1 - \varepsilon_1), \quad |\varepsilon_1| \leq p2^{-t},$$

$$1 - \mu = \left( 1 - \frac{\nu N}{a_1 + \text{sign}(a_1)N} \right) (1 - \varepsilon_2), \quad |\varepsilon_2| \leq 2^{-t},$$

$$1 - \varrho = \left( 1 - \frac{\alpha A + (\nu + \varepsilon_3 - \varepsilon_3 \nu) |a_1| N}{A + |a_1| N} \right) (1 - \varepsilon_4), \quad |\varepsilon_3|, |\varepsilon_4| \leq 2^{-t}.$$

Stąd wynika, że

$$(1) \quad \begin{aligned} \nu &= \frac{1}{2}\alpha + \varepsilon_1 + \Delta_1, \\ \mu &= \frac{\frac{1}{2}\alpha + \varepsilon_1}{1+h} + \varepsilon_2 + \Delta_2, \\ \varrho &= \frac{\alpha + (\frac{1}{2}\alpha + \varepsilon_1 + \varepsilon_3)h}{1+h} + \varepsilon_4 + \Delta_3, \end{aligned}$$

gdzie

$$h = \frac{|a_1|}{\|\vec{a}\|},$$

a błędy  $\Delta_j$  mają oszacowania na poziomie  $2^{-2t}$ . Chcemy podkreślić, że rozważamy tu obliczenie  $R$  w „niepełnej” arytmetyce  $fl_2$ , w następujący sposób

$$(2) \quad \tilde{R} = fl[fl(fl_2(\|\vec{a}\|^2) + \sqrt{fl(fl_2(\|\vec{a}\|^2))} |a_1|)].$$

Pokażemy teraz, jak zniekształcenie wektora  $\vec{u}$  i parametru  $\tilde{R}$  wpływa na przekształcenie  $P$ . Niech tak jak w § 4

$$\|B\| = \|\tilde{P} - P\| \leq 2^{-t} \hat{\gamma}.$$

Podobnie jak w poprzednim punkcie

$$B = \frac{1}{R} (-\vec{u}\vec{u}^T \varrho_1 + \vec{e}_1 \vec{u}^T u_1 \mu + \vec{u} \vec{e}_1^T u_1 \mu - \vec{e}_1 \vec{e}_1^T u_1^2 \mu^2),$$

gdzie

$$\varrho_1 = \varrho + \varrho^2 + \varrho^3 + \dots = \varrho + \Delta_4, \quad |\Delta_4| \leq 2^{-2t}.$$

Błąd  $\alpha$  dla arytmetyki  $fl_2$  i  $\tilde{fl}$  ma takie samo oszacowanie jak w § 4 dla arytmetyki  $cfl_2$  i  $c\tilde{fl}$ , a dla arytmetyki  $fl$  możemy pokazać, że  $|\alpha| \leq n2^{-t}$ . Odpowiednia tabelka dla  $\hat{\gamma}$  wygląda więc następująco:

	$fl$	$\tilde{fl}$	$fl_2$
$\hat{\gamma}$	$\frac{1}{2}n+4 + \sqrt{(\frac{1}{2}n+4)^2 + (\frac{1}{2}n+2)^2} < (1+\sqrt{2})(\frac{1}{2}n+4)$	12	9,7

Uwzględniając analogiczne zmiany w oszacowaniach elementów macierzy  $D_i$  otrzymujemy dla błędu otrzymanego obrazu  $\vec{g}$  wektora  $\vec{b}$  oszacowanie:

$$\|\vec{g} - \vec{c}\| \leq (\hat{\gamma} + \delta) \|\vec{b}\| 2^{-t} \equiv K_w \|\vec{c}\| 2^{-t},$$

przy czym

	f1	f̃1	f1 <sub>2</sub>
$\hat{\delta}$	2n+7	13	9
$K_w$	3,2n+17	25	18,7

Otrzymane wartości wskaźnika kumulacji algorytmu,  $K_w$ , są dla arytmetyk  $\tilde{f}_1$  i  $f_{1_2}$  większe od wartości podanych w [1] (porównaj również [6], str. 154). Jak już wspominaliśmy, jest to związane m.in. ze sposobem obliczania wartości parametru  $R$  (zob. (2)). Ponadto, ze względu na ograniczenie zasięgu stosowania podwójnej precyzji w arytmetyce  $f_{1_2}$  jedynie do kumulowania sum wielu składników, otrzymujemy większe oszacowania błędu wytworzonego przy obliczaniu  $\tilde{P}\tilde{b}$ . Mianowicie, pokazaliśmy, że mnożnik  $\hat{\delta} = 9$ , gdy tymczasem Wilkinson, zakładając pełne wykorzystanie podwójnej precyzji w arytmetyce  $f_{1_2}$ , otrzymuje  $\hat{\delta} = 3,35$ . Z tych samych względów istnieje możliwość obniżania oszacowań otrzymanych dla zespolonego zadania  $H$ , gdybyśmy rozszerzyli zakres stosowania podwójnej precyzji w arytmetyce  $cf_{1_2}$ .

Tak więc pokazaliśmy, że zespolone i rzeczywiste zadanie (H) dla  $\vec{e} = \vec{e}_1$  może być realizowane w sposób numerycznie poprawny w każdej z rozważanych arytmetyk.

**6. Numeryczna realizacja ogólnego zadania (H).** Niech  $P_a$  oznacza przekształcenie Householdera przeprowadzające wektor  $\vec{a}$  na wektor o kierunku wektora  $\vec{e}$ . Wówczas  $P_a$  jest reprezentowane przez wektor  $\vec{u}(a)$  i parametr  $R_a$  (zob. (2.5)–(2.7)).

Wprowadzimy pomocnicze oznaczenia:  $\vec{v}$  — numerycznie wyznaczony wektor  $\vec{u}(a)$ ,  $\vec{f}$  — zaburzony wektor  $\vec{a}$ , taki, by  $(\vec{e}^H \vec{a})_{\text{ob1}} = \vec{e}^H \vec{f}$ ,  $\vec{P}_a$  — numerycznie wyznaczane przekształcenie Householdera reprezentowane przez  $\vec{v}$  i  $\tilde{R}_a$ ,  $\vec{g}$  — numerycznie wyznaczony wektor  $\tilde{P}_a \vec{b}$ .

W przeprowadzonej przez nas analizie numerycznej ogólnego zadania (szczegóło pomijamy) otrzymaliśmy następujące oszacowania

$$\begin{aligned} \|\vec{f} - \vec{a}\| &\leq K_d \|\vec{a}\| 2^{-t}, & \|\vec{v} - \vec{u}(f)\| &\leq \hat{\mu} \|\vec{u}(f)\| 2^{-t}, \\ \|\tilde{P}_a - P_f\| &\leq \hat{\gamma} 2^{-t}, & \|\vec{g} - P_f \vec{b}\| &\leq K_w \|\vec{b}\| 2^{-t} = K_w \|\vec{c}\| 2^{-t}, \end{aligned}$$

gdzie  $\vec{c} = P_a \vec{b}$ ,  $P_f$  oznacza przekształcenie Householdera skonstruowane dla wektora  $\vec{f}$  i  $\vec{e}$ , a  $K_d$ ,  $K_w$ ,  $\hat{\gamma}$ ,  $\hat{\mu}$  są odpowiednio równe:

	f1	f̃1	f1 <sub>2</sub>	cf1	cf̃1	cf1 <sub>2</sub>
$\hat{\mu}$	7n+2,8	4,9	3,5	2,1n+11,6	15,8	9,6
$\hat{\gamma}$	6,8n+13	33,4	22,8	14,5n+56	91,6	50,2
$K_d$	0	0	0	$n + \sqrt{2}$	$2 + \sqrt{2}$	$\langle 2^{-t} \rangle$
$K_w$	8,8n+20	47	32	16,5n+71	110	62

Widzimy więc, że w przypadku zespolonym otrzymany obraz  $\vec{g}$  wektora  $\vec{b}$  dla transformacji  $P_a$  może być interpretowany jako trochę zaburzony obraz wektora  $\vec{b}$  dla transformacji  $P_{\vec{f}}$ , przy czym wektor  $\vec{f} - \vec{a}$  jest małym zaburzeniem wektora  $\vec{a}$ . W tym sensie realizację ogólnego zadania (H) możemy uważać za numerycznie poprawną. W przypadku rzeczywistym  $\vec{f} = \vec{a}$ , w przypadku zespolonym zaś dla arytmetyki  $\text{cfl}_2$  zaburzenia  $\vec{f} - \vec{a}$  są na poziomie  $2^{-2^t} \|\vec{a}\|$ . W trakcie analizy numerycznej zespolonego zadania (H) okazało się, że o dokładności obliczonego wektora  $P_a \vec{b}$  decyduje przede wszystkim dokładność, z jaką oblicza się iloczyn skalarny  $\vec{e}^H \vec{a}$ . O wrażliwości wektora  $\vec{u}(\vec{a})$  na zaburzenia wektora  $\vec{a}$  decyduje bowiem wrażliwość ilorazu  $\vec{e}^H \vec{a} / \|\vec{e}^H \vec{a}\|$  na małe zmiany składowych wektora  $\vec{a}$  (a niekiedy może być ona duża).

Gdyby z jakichś powodów potrzebna była (również dla przypadku zespolonego) wysoka dokładność względna wektora  $\vec{g}$ , należałoby wtedy stosować jedynie arytmetykę  $\text{cfl}_2$ .

#### Bibliografia

- [1] A. Kielbasiński, *Algorytm sumowania z poprawkami i niektóre jego zastosowania*, *Matematyka Stosowana* 1 (1973), str. 23–41.
- [2] — *Analiza numeryczna algorytmu ortogonalizacji Grama–Schmidta*, *ibid.* 2 (1974), str. 15–35.
- [3] — *Podstawowe pojęcia analizy błędów w metodach numerycznych algebry liniowej*, *ibid.* 4 (1975), str. 5–27.
- [4] J. Majchrowska i A. Smoktunowicz, *Analiza numeryczna algorytmu Ortegi–Householdera w dziedzinie zespolonej*, ten tom, str. 55–66.
- [5] J. H. Wilkinson, *Błędy zaokrągleń w procesach algebraicznych*, Warszawa 1967.
- [6] — *Алгебраическая проблема собственных значений*, Москва 1970.
- [7] J. H. Wilkinson and C. Reinsch, *Handbook for automatic computation*, Springer-Verlag, 1971, str. 359–371 (lub *Numer. Math.* 12 (1968), str. 349–368).