



Chebyshev acceleration of iterative refinement

M Arioli, J Scott

June 2011

©2011 Science and Technology Facilities Council

Enquiries about copyright, reproduction and requests for additional copies of this report should be addressed to:

RAL Library
STFC Rutherford Appleton Laboratory
R61
Harwell Oxford
Didcot
OX11 0QX

Tel: +44(0)1235 445384
Fax: +44(0)1235 446403
email: libraryral@stfc.ac.uk

Science and Technology Facilities Council reports are available online at: <http://epubs.stfc.ac.uk>

ISSN 1358- 6254

Neither the Council nor the Laboratory accept any responsibility for loss or damage arising from the use of information contained in any of their reports or in any communication about their tests or investigations.

Chebyshev Acceleration of Iterative Refinement

Mario Arioli^{1,2} and Jennifer Scott^{1,2}

Abstract. We analyse how variants of the Chebyshev algorithm can be used to accelerate the iterative refinement procedure without loss of numerical stability and at a computational cost at each iteration that is only marginally greater than that of iterative refinement. An error analysis of the procedure is presented and numerical tests on selected sparse test problems are used to corroborate the theory and illustrate the potential savings offered by Chebyshev acceleration.

Key words. Chebyshev method, iterative refinement, Gaussian factorization, sparse matrices.

AMS subject classifications. 65F05, 65F50, 65F10, 65G50

1. Introduction. The combination of Gaussian factorization with partial pivoting followed by a few steps of iterative refinement can compute an approximate solution of a linear system of equations that is backward stable, i.e. the residual norm is less than or equal to machine precision times the norm of the data. However, when threshold partial pivoting is used to limit fill-in in the factorization of large sparse systems, the number of iterative refinement steps can be large. Furthermore, when the Gaussian factorization is computed in single precision and then a double precision backward stable approximate solution is recovered using iterative refinement, the number of refinement steps can also be very large and the cost prohibitive. It is particularly important to limit the number of refinement steps on modern multicore architectures where the solve phase of a sparse direct solver can represent a potential bottleneck (see, for example, [10]).

Following the results of [6, 7, 8, 13, 14, 15], in this paper we analyse how variants of the Chebyshev algorithm can be used to accelerate the iterative refinement procedure without loss of numerical stability and at a computational cost for each iteration that is only marginally greater than that of the iterative refinement.

In the rest of this section, we introduce the notation that we will use throughout and then summarize iterative refinement and its properties and potential weakness. In Section 2, we describe two Chebyshev acceleration algorithm variants. We present an error analysis of the main algorithm in Section 3 and in Section 4 we discuss how to automatically choose some of the parameters. In Section 5, we present numerical results for sparse systems that arise from practical applications that corroborate the theoretical results of the previous sections and, in Section 6, we give some final comments. Finally, in the following $\|\cdot\|$ will denote the Euclidean norm for \mathbb{R}^n and the corresponding induced norm for the matrices.

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$, with $\text{rank}(\mathbf{A}) = n$. The linear system

$$\mathbf{Ax} = \mathbf{b} \tag{1.1}$$

has a unique solution $\hat{\mathbf{x}}$. We assume Gaussian elimination is performed using floating-point arithmetic with relative precision ϵ . Thus, the computed factors $\hat{\mathbf{L}}$ and $\hat{\mathbf{U}}$ satisfy the relation [5, 12]

$$\mathbf{A} + \mathbf{F} = \hat{\mathbf{L}}\hat{\mathbf{U}} = \mathbf{M}, \tag{1.2}$$

¹STFC Rutherford Appleton Laboratory, Didcot, Oxon, OX11 0QX, UK.
e-mail: mario.arioli@stfc.ac.uk, jennifer.scott@stfc.ac.uk

²This work was supported by EPSRC Grant EP/E053351/1.

where $\mathbf{F} \in \mathbb{R}^{n \times n}$ and

$$|\mathbf{F}| \leq c(n)\epsilon|\widehat{\mathbf{L}}|\widehat{\mathbf{U}}|. \quad (1.3)$$

Here and elsewhere, $|\mathbf{B}|$ denotes the matrix of entries equal to the absolute values of the corresponding entries in the matrix \mathbf{B} . From (1.2), it follows that

$$\mathbf{x} = \mathbf{M}^{-1}(\mathbf{b} + \mathbf{F}\mathbf{x}) = \mathbf{M}^{-1}(\mathbf{b} - \mathbf{A}\mathbf{x} + \mathbf{M}\mathbf{x}) = \mathbf{M}^{-1}(\mathbf{r}(\mathbf{x}) + \mathbf{M}\mathbf{x}), \quad (1.4)$$

where $\mathbf{r}(\mathbf{x})$ is the residual $\mathbf{b} - \mathbf{A}\mathbf{x}$. Thus, \mathbf{x} is the fixed point of $\mathfrak{F}(\mathbf{x})$, where

$$\mathfrak{F}(\mathbf{x}) = \mathbf{x} + \mathbf{M}^{-1}\mathbf{r}(\mathbf{x}).$$

If β denotes the scaled residual

$$\beta = \frac{\|\mathbf{r}(\mathbf{x})\|}{\|\mathbf{A}\|\|\mathbf{x}\| + \|\mathbf{b}\|}, \quad (1.5)$$

then given a convergence tolerance $\eta > 0$, it is straightforward to write down the basic algorithm for iterative refinement

ALGORITHM 1.1. *Iterative refinement*
 Let $\mathbf{x}^{(0)} = \mathbf{M}^{-1}\mathbf{b}$ and $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$.
 Initialise $k = 0$.
while $\beta^{(k)} > \eta$ **do**
 $\delta\mathbf{x} = \mathbf{M}^{-1}\mathbf{r}^{(k)}$;
 $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \delta\mathbf{x}$;
 $\mathbf{r}^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)}$;
 $\beta^{(k+1)} = \|\mathbf{r}^{(k+1)}\| / (\|\mathbf{A}\|\|\mathbf{x}^{(k+1)}\| + \|\mathbf{b}\|)$;
 $k = k + 1$.
end while

If the spectral radius of $\sigma(\mathbf{M}^{-1}\mathbf{F})$ of $\mathbf{M}^{-1}\mathbf{F}$ satisfies

$$\sigma(\mathbf{M}^{-1}\mathbf{F}) < 1$$

in exact arithmetic, Algorithm 1.1 produces a sequence $\mathbf{x}^{(k)}$ that converges to $\hat{\mathbf{x}}$. Furthermore, from (1.4),

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \mathbf{A}\mathbf{M}^{-1}\mathbf{r}^{(k)} = \mathbf{F}\mathbf{M}^{-1}\mathbf{r}^{(k)}. \quad (1.6)$$

Therefore, if $\sigma(\mathbf{F}\mathbf{M}^{-1}) < 1$ the residuals converge to zero.

REMARK 1. Let $\mathbf{W}\mathbf{J}\mathbf{W}^{-1}$ be the Jordan form of $\mathbf{M}^{-1}\mathbf{F}$. It is immediate to see that $\mathbf{M}\mathbf{W}\mathbf{J}\mathbf{W}^{-1}\mathbf{M}^{-1}$ is the Jordan form of $\mathbf{F}\mathbf{M}^{-1}$. Thus, $\mathbf{M}^{-1}\mathbf{F}$ and $\mathbf{F}\mathbf{M}^{-1}$ have the same spectrum and $\sigma(\mathbf{M}^{-1}\mathbf{F}) = \sigma(\mathbf{F}\mathbf{M}^{-1})$. In the following, we will use interchangeably both $\sigma(\mathbf{M}^{-1}\mathbf{F})$ and $\sigma(\mathbf{F}\mathbf{M}^{-1})$.

REMARK 2. An alternative formulation for (1.1) based on relation (1.4) is

$$\mathbf{x} = \mathbf{M}^{-1}\mathbf{F}\mathbf{x} + \mathbf{M}^{-1}\mathbf{b}. \quad (1.7)$$

This is not of practical use, but it will be useful when developing the theory of Chebyshev acceleration that follows in Section 2.

REMARK 3. Assume that $\sigma(\mathbf{M}^{-1}\mathbf{F}) = 0.5$. To achieve a reduction of three orders of magnitude in the initial residual, the required number of steps of iteration refinement is

$$\text{iter} = \lceil \frac{\log_{10}(10^{-3})}{\log_{10}(0.5)} \rceil,$$

which is approximately 10. The cost of performing this number of iterations may be unacceptably high, for example, if the factors are held out-of-core. In the next section, we propose a variant of Chebyshev acceleration algorithm that may improve the rate of convergence.

2. Chebyshev acceleration. Chebyshev polynomials can be defined by the following 2-term recurrence formula (see [8, page 46]):

$$\begin{cases} T_0(z) = 1, & T_1(z) = z \\ T_{k+1}(z) = 2zT_k(z) - T_{k-1}(z) & k \geq 1. \end{cases}$$

The optimal properties of Chebyshev polynomials given in Theorem 4.2.1 of [8, page 47] can be summarised as follows: let $d > 1$ and set

$$\mathcal{F}_k(z) = \frac{T_k(z)}{T_k(d)},$$

then \mathcal{F}_k has minimum l_∞ norm on the interval $[-1, 1]$ over all polynomials Q_k of degree less than or equal to n and satisfying the condition $Q_k(d) = 1$, and

$$\max_{z \in [-1, 1]} |\mathcal{F}_k(z)| = \frac{1}{T_k(d)}.$$

We now summarize some classical results on Chebyshev acceleration. We refer the reader to [8, Chapters 4 and 12] and [13, 14] for further details. If all the eigenvalues of $\mathbf{M}^{-1}\mathbf{F}$ lie in the interior of an ellipse that is centred at the origin, symmetric with respect to the real axis (the matrix is real so the eigenvalues either are real or are in complex conjugate pairs) and has principal semi-axes a and b , i.e.

$$\left(\frac{x}{a}\right)^2 + \left(\frac{y}{b}\right)^2 = 1, \tag{2.1}$$

then the following theorem holds (see [8, Theorem 12-2.1]).

THEOREM 2.1. Let \mathcal{D} be the region enclosed by (2.1) where $b < a < 1$. If \mathcal{S}_j is the set of all real polynomials $p_j(z)$ of degree at most j such that $p_j(1) = 1$, then the polynomial

$$\wp_j(z) = \frac{T_j(z/c)}{T_j(1/c)}, \quad \text{where } c^2 = a^2 - b^2,$$

is the unique polynomial in the set \mathcal{S}_j such that

$$\max_{z \in \mathcal{D}} |\wp_j(z)| \leq \max_{z \in \mathcal{D}} |p_j(z)|, \quad p_j(z) \in \mathcal{S}_j.$$

Manteuffel [13] showed that this result cannot be extended to the case $a < b < 1$. In this case, c is purely imaginary. However, the $\wp_j(z)$ are still real and the following weaker result can be proved [13]:

$$\lim_{j \rightarrow \infty} \left(\max_{z \in \mathcal{D}} |\wp_j(z)| \right)^{1/j} \leq \lim_{j \rightarrow \infty} \left(\max_{z \in \mathcal{D}} |p_j(z)| \right)^{1/j}, \quad p_j(z) \in \mathcal{S}_j. \quad (2.2)$$

From formula (2.2), we have that the polynomials $\wp_j(z)$ are asymptotically optimal and, furthermore, it has been noted [13, 8] that the asymptotic behaviour is very rapidly reached.

Following along the lines of [8], we now describe the Chebyshev acceleration algorithm. The polynomials $\wp_j(z)$ are defined as follows:

$$\begin{cases} \wp_0 = 1, & \wp_1 = z \\ \wp_{j+1}(z) = \varrho_{j+1} z \wp_j(z) + (1 - \varrho_{j+1}) \wp_{j-1}(z) \\ \varrho_{j+1} = \frac{2}{c} \frac{T_j(1/c)}{T_{j+1}(1/c)} \end{cases}$$

The Chebyshev relations for problem (1.7) are then given by

$$\begin{cases} \mathbf{x}^{(0)} = \mathbf{M}^{-1} \mathbf{b}, & \varrho_1 = 1 \\ \mathbf{x}^{(j+1)} = \varrho_{j+1} (\mathbf{M}^{-1} \mathbf{F} \mathbf{x}^{(j)} + \mathbf{M}^{-1} \mathbf{b}) + (1 - \varrho_{j+1}) \mathbf{x}^{(j-1)}, & j = 0, \dots, \end{cases}$$

Using (1.4), this can be simplified:

$$\begin{cases} \mathbf{x}^{(0)} = \mathbf{M}^{-1} \mathbf{b}, & \varrho_1 = 1 \\ \mathbf{x}^{(j+1)} = \varrho_{j+1} (\mathbf{M}^{-1} \mathbf{r}(\mathbf{x}^{(j)}) + \mathbf{x}^{(j)}) + (1 - \varrho_{j+1}) \mathbf{x}^{(j-1)}, & j = 0, \dots, \end{cases} \quad (2.3)$$

We observe that computing the ϱ_j is straightforward:

$$\varrho_{j+1} = \begin{cases} 1, & \text{if } j = 0 \\ (1 - \frac{1}{2}c^2)^{-1}, & \text{if } j = 1 \\ (1 - \frac{1}{4}c^2\varrho_j)^{-1}, & \text{if } j \geq 2 \end{cases} \quad (2.4)$$

From the maximum modulus principle and the analyticity of \wp_j , $\wp_j(z)$ will take its maximum on the ellipse (2.1). Moreover, we have [8] that

$$\max_{z \in \mathcal{D}} |\wp_j(z)| = \left[\frac{a+b}{1 + \sqrt{1-c^2}} \right]^j. \quad (2.5)$$

Finally, assuming the ϱ_j have been precomputed, simple algebraic manipulations leads to the following algorithm:

ALGORITHM 2.1. *Chebyshev acceleration of iterative refinement.*

Let $\mathbf{x}^{(0)} = \mathbf{M}^{-1}\mathbf{b}$, $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$.

Initialise $k = 0$.

while $\beta^{(k)} > \eta$ **do**

$\mathbf{w}^{(k)} = \mathbf{x}^{(k)} + \mathbf{M}^{-1}\mathbf{r}^{(k)}$;

$\mathbf{x}^{(k+1)} = \varrho_{k+1}\mathbf{w}^{(k)} + (1 - \varrho_{k+1})\mathbf{x}^{(k-1)}$;

$\mathbf{r}^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)}$;

$\beta^{(k+1)} = \|\mathbf{r}^{(k+1)}\| / (\|\mathbf{A}\|\|\mathbf{x}^{(k+1)}\| + \|\mathbf{b}\|)$;

$k = k + 1$.

end while

We have followed the analysis and the evidence given in [15, 7, 6] and have chosen to compute the residuals explicitly. Recursive expressions can easily be computed but they can be less stable [6].

REMARK 4. *For the successful convergence of Algorithm 2.1 it is necessary to forecast the equation of an ellipse that envelops the whole spectrum of $\mathbf{M}^{-1}\mathbf{F}$. If $\sigma(\mathbf{M}^{-1}\mathbf{F})$ lies outside the chosen ellipse and $c^2 \ll 1$ or the ellipse degenerates to a circle ($a = b$), the asymptotic behaviour of Algorithm 2.1 will be the same as that of iterative refinement and thus it will give no acceleration. In the first case, $\varrho_\infty = \lim_{j \rightarrow \infty} \varrho_j = 2/(1 + \sqrt{1 - c^2}) \approx 1$, while in the case $a = b$, $\varrho_j = 1 \forall j$.*

REMARK 5. *In Remark 3, we observed that if $\sigma(\mathbf{M}^{-1}\mathbf{F}) = 0.5$, iterative refinement will require approximately 10 steps to reduce the initial residual by three orders of magnitude. The asymptotic rate of convergence, i.e. the logarithm of the right-hand side of (2.5), describes the number of steps that Algorithm 2.1 needs to reduce the residual by one order of magnitude. The number of steps required to reduce it by p orders of magnitude is*

$$\text{iter} = \lceil \frac{\log_{10}(10^{-p})}{\log_{10}(\frac{a+b}{1+\sqrt{1-c^2}})} \rceil. \quad (2.6)$$

If, in our example, the ellipse

$$\left(\frac{x}{0.5}\right)^2 + \left(\frac{y}{0.05}\right)^2 = 1$$

contains all the eigenvalues, the Chebyshev accelerated algorithm will need approximately 6 steps to obtain a reduction of three orders of magnitude. This illustrates the potential savings offered by Algorithm 2.1.

We can also introduce a simple variant of Algorithm 2.1, based on (2.3):

ALGORITHM 2.2. *Chebyshev acceleration of iterative refinement.*

Let $\mathbf{x}^{(0)} = \mathbf{M}^{-1}\mathbf{b}$, $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$.

Initialise $k = 0$.

while $\beta^{(k)} > \eta$ **do**

$\Delta\mathbf{x}^{(k+1)} = \varrho_{k+1}\mathbf{M}^{-1}\mathbf{r}^{(k)} - (1 - \varrho_{k+1})\Delta\mathbf{x}^{(k)}$;

$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta\mathbf{x}^{(k+1)}$;

$\mathbf{r}^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)}$;

$\beta^{(k+1)} = \|\mathbf{r}^{(k+1)}\| / (\|\mathbf{A}\|\|\mathbf{x}^{(k+1)}\| + \|\mathbf{b}\|)$;

$k = k + 1$.

end while

This variant is slightly more awkward to analyse from a roundoff point of view. However, the numerical results do not differ significantly from those obtained using Algorithm 2.1, which we analyse in the next section.

3. Iterative refinement and Chebyshev error analysis. We assume finite precision arithmetic with relative precision ε is used, i.e. the arithmetic operations $\diamond \in \{+, -, *, /\}$ satisfy

$$fl(g \diamond r) = (1 + \xi)g \diamond r, \quad |\xi| \leq \varepsilon,$$

where with $fl(\cdot)$ denotes the actual results in finite precision. Taking into account formulae (2.4), we assume that the ϱ_j values have been precomputed using extended precision and that they are the correct rounded results to ε accuracy. From the formulae in Algorithm 2.1 and using standard techniques [12], we have

$$\left. \begin{aligned} \bar{\mathbf{r}}^{(k)} &= fl(\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}) = (\mathbf{I} + \mathbf{\Gamma}_3^{(k)})(\mathbf{b} - (\mathbf{M} + \mathbf{G}^{(k)})\mathbf{x}^{(k)}), \\ |\mathbf{G}^{(k)}| &\leq c_G(n)\varepsilon|\mathbf{A}| \ll 1 \quad \text{and} \quad |\mathbf{\Gamma}_3^{(k)}| \leq \varepsilon \mathbf{I}. \end{aligned} \right\} \quad (3.1)$$

Furthermore, from (3.1)

$$\left. \begin{aligned} \bar{\mathbf{r}}^{(k)} &= \mathbf{r}^{(k)} + \mathbf{g}^{(k)}, \\ \mathbf{g}^{(k)} &= \mathbf{\Gamma}_3^{(k)}\mathbf{r}^{(k)} - (\mathbf{I} + \mathbf{\Gamma}_3^{(k)})\mathbf{G}^{(k)}\mathbf{x}^{(k)}, \\ |\mathbf{g}^{(k)}| &\leq 3\varepsilon(|\mathbf{b}| + |\mathbf{A}||\mathbf{x}^{(k)}|). \end{aligned} \right\} \quad (3.2)$$

In Algorithm 2.1, the linear system

$$\mathbf{M}\mathbf{z}^{(k)} = \bar{\mathbf{r}}^{(k)}$$

must be solved. Taking into account the properties of forward and backward substitutions, the computed solution satisfies

$$(\mathbf{M} + \mathbf{E}^{(k)})fl(\mathbf{z}^{(k)}) = \bar{\mathbf{r}}^{(k)}, \quad |\mathbf{E}^{(k)}| \leq c_0(n)\varepsilon|\widehat{\mathbf{L}}||\widehat{\mathbf{U}}|.$$

Setting $\widetilde{\mathbf{M}}_k = \mathbf{M} + \mathbf{E}^{(k)}$, $\bar{\mathbf{x}}^{(k)} = fl(\mathbf{x}^{(k)})$, and $\bar{\mathbf{w}}^{(k)} = fl(\mathbf{w}^{(k)})$, we have

$$\left. \begin{aligned} \bar{\mathbf{w}}^{(k+1)} &= (\mathbf{I} + \mathbf{\Gamma}_1^{(k)})(\widetilde{\mathbf{M}}_k^{-1}\bar{\mathbf{r}}^{(k)} + \bar{\mathbf{x}}^{(k)}), \\ |\mathbf{\Gamma}_1^{(k)}| &\leq \varepsilon \mathbf{I}, \end{aligned} \right\} \quad (3.3)$$

and, finally,

$$\left. \begin{aligned} \bar{\mathbf{x}}^{(k+1)} &= (\mathbf{I} + \mathbf{\Gamma}_4^{(k)}) \left[\varrho_{k+1}(\mathbf{I} + \mathbf{\Gamma}_2^{(k)})\bar{\mathbf{w}}^{(k)} + (1 - \varrho_{k+1})(\mathbf{I} + \mathbf{\Gamma}_3^{(k)})\bar{\mathbf{x}}^{(k-1)} \right], \\ |\mathbf{\Gamma}_i^{(k)}| &\leq \varepsilon \mathbf{I}, \quad i = 2, 3, 4. \end{aligned} \right\} \quad (3.4)$$

From (3.3) and (3.4), we deduce that

$$\bar{\mathbf{x}}^{(k+1)} = \varrho_{k+1}(\mathbf{I} + \widehat{\mathbf{\Gamma}}_1^{(k)})\bar{\mathbf{w}}^{(k)} + (1 - \varrho_{k+1})(\mathbf{I} + \widehat{\mathbf{\Gamma}}_2^{(k)})\bar{\mathbf{x}}^{(k-1)} \quad (3.5)$$

$$= \varrho_{k+1}(\mathbf{I} + \widehat{\mathbf{\Gamma}}_3^{(k)})(\widetilde{\mathbf{M}}_k^{-1}\bar{\mathbf{r}}^{(k)} + \bar{\mathbf{x}}^{(k)}) + (1 - \varrho_{k+1})(\mathbf{I} + \widehat{\mathbf{\Gamma}}_2^{(k)})\bar{\mathbf{x}}^{(k-1)}, \quad (3.6)$$

with $|\widehat{\Gamma}_i^{(k)}| \lesssim 3\varepsilon \mathbf{I}$ for all k and $i = 1, 2, 3$. Although they are uniformly bounded, the $\widehat{\Gamma}_i^{(k)}$ and $\mathbf{E}^{(k)}$ depend non-linearly on $\bar{\mathbf{w}}^{(k)}$ and $\bar{\mathbf{x}}^{(k)}$. Furthermore, from (3.2), (3.5) and (3.6), the exact residual $\mathbf{r}^{(j)} = \mathbf{b} - \mathbf{A}\bar{\mathbf{x}}^{(j)}$ satisfies

$$\begin{aligned} \mathbf{r}^{(k+1)} &= \varrho_{k+1} \left[\mathbf{b} - \mathbf{A}(\mathbf{I} + \widehat{\Gamma}_3^{(k)}) (\widetilde{\mathbf{M}}_k^{-1} \bar{\mathbf{r}}^{(k)} + \bar{\mathbf{x}}^{(k)}) \right] + \\ &\quad (1 - \varrho_{k+1}) \left[\mathbf{b} - \mathbf{A}(\mathbf{I} + \widehat{\Gamma}_2^{(k)}) \bar{\mathbf{x}}^{(k-1)} \right] \\ &= \varrho_{k+1} \left[\mathbf{r}^{(k)} - \mathbf{A}(\mathbf{I} + \widehat{\Gamma}_3^{(k)}) \widetilde{\mathbf{M}}_k^{-1} \bar{\mathbf{r}}^{(k)} - \mathbf{A} \widehat{\Gamma}_3^{(k)} \bar{\mathbf{x}}^{(k)} \right] + \\ &\quad (1 - \varrho_{k+1}) \left[\mathbf{r}^{(k-1)} - \mathbf{A} \widehat{\Gamma}_2^{(k)} \bar{\mathbf{x}}^{(k-1)} \right] \\ &= \varrho_{k+1} \left[\mathbf{I} - \mathbf{A}(\mathbf{I} + \widehat{\Gamma}_3^{(k)}) \widetilde{\mathbf{M}}_k^{-1} \right] \mathbf{r}^{(k)} + (1 - \varrho_{k+1}) \mathbf{r}^{(k-1)} \\ &\quad - \varrho_{k+1} (\mathbf{g}^{(k)} + \mathbf{A} \widehat{\Gamma}_3^{(k)} \bar{\mathbf{x}}^{(k)}) - (1 - \varrho_{k+1}) \mathbf{A} \widehat{\Gamma}_2^{(k)} \bar{\mathbf{x}}^{(k-1)}. \end{aligned}$$

Therefore, we have the following recursive expression

$$\mathbf{r}^{(k+1)} = \varrho_{k+1} \mathbf{H}_k \mathbf{r}^{(k)} + (1 - \varrho_{k+1}) \mathbf{r}^{(k-1)} + \mathbf{f}^{(k+1)}, \quad (3.7)$$

where

$$\begin{aligned} \mathbf{H}_k &= \mathbf{I} - \mathbf{A}(\mathbf{I} + \widehat{\Gamma}_3^{(k)}) \widetilde{\mathbf{M}}_k^{-1}, \\ \mathbf{f}^{(k+1)} &= -\varrho_{k+1} (\mathbf{g}^{(k)} + \mathbf{A} \widehat{\Gamma}_3^{(k)} \bar{\mathbf{x}}^{(k)}) - (1 - \varrho_{k+1}) \mathbf{A} \widehat{\Gamma}_2^{(k)} \bar{\mathbf{x}}^{(k-1)}, \end{aligned}$$

and, from the bounds in (3.2), (3.3), and (3.4), it follows that

$$|\mathbf{f}^{(k+1)}| < 3\varepsilon |1 - \varrho_{k+1}| |\mathbf{A}| |\bar{\mathbf{x}}^{(k-1)}| + 6\varepsilon \varrho_{k+1} (|\mathbf{A}| |\bar{\mathbf{x}}^{(k)}| + |\mathbf{b}|).$$

3.1. Analysis of \mathbf{H}_k . We assume that numerical exceptions (overflows or underflows) do not occur during the execution of Algorithm 2.1. This is a necessary condition for continuous dependence of the errors on the data. Moreover, we assume that

$$\|\mathbf{E}_k \mathbf{M}_k^{-1}\| < 1. \quad (3.8)$$

With this assumption, $\widetilde{\mathbf{M}}_k$ is nonsingular and

$$\widetilde{\mathbf{M}}_k^{-1} = \mathbf{M}^{-1} (\mathbf{I} + \mathbf{E}_k \mathbf{M}_k^{-1})^{-1} = \mathbf{M}^{-1} (\mathbf{I} - \mathbf{E}_k \mathbf{M}_k^{-1} (\mathbf{I} + \mathbf{E}_k \mathbf{M}_k^{-1})^{-1}). \quad (3.9)$$

From (1.2), (1.3), (3.8), and (3.9), it follows that

$$\begin{aligned} \mathbf{H}_k &= \mathbf{I} - \mathbf{A}(\mathbf{I} + \widehat{\Gamma}_3^{(k)}) \mathbf{M}^{-1} (\mathbf{I} - \mathbf{E}_k \widetilde{\mathbf{M}}_k^{-1}) \\ &= \mathbf{I} - (\mathbf{M} - \mathbf{F}) \mathbf{M}^{-1} (\mathbf{I} + \mathbf{M} \widehat{\Gamma}_3^{(k)} \mathbf{M}^{-1}) (\mathbf{I} - \mathbf{E}_k \widetilde{\mathbf{M}}_k^{-1}) \\ &= \mathbf{F} \mathbf{M}^{-1} - \mathbf{M} \widehat{\Gamma}_3^{(k)} \mathbf{M}^{-1} + \mathbf{E}_k \widetilde{\mathbf{M}}_k^{-1} + \mathcal{E}, \end{aligned}$$

where $\mathcal{E} = \mathcal{O}(\varepsilon^2)$. Thus, if for each k the matrices $\mathbf{F} \mathbf{M}^{-1}$, $\mathbf{M} \widehat{\Gamma}_3^{(k)} \mathbf{M}^{-1}$, and $\mathbf{E}_k \widetilde{\mathbf{M}}_k^{-1}$ have Euclidean norm strictly less than 0.25, $\|\mathbf{H}_k\| < 1$ and $\mathbf{I} - \mathbf{H}_k$ and $\mathbf{I} + \mathbf{H}_k$ are invertible.

3.2. Error bounds. Taking into account the previous results, and if we assume \mathbf{H}_k and $\mathbf{f}^{(k)}$ depend continuously on the data, equation (3.6) has a fixed point on a large enough compact convex set of \mathbb{R}^n (Generalized Brouwer Fixed Point Theorem [3, Theorem 3.2]). Assuming in finite precision arithmetic Algorithm 2.1 computes a sequence $\bar{\mathbf{x}}^{(k)}$ that converges to the point $\bar{\mathbf{x}}^\infty$, then the residuals $\mathbf{r}^{(k)}$ converge to \mathbf{r}^∞ . From (3.7) and taking into account that ϱ_k converges to a finite $\varrho_\infty = 2/(1 + \sqrt{1 - c^2})$, we have

$$\mathbf{r}^\infty = \varrho_\infty \mathbf{H}_\infty \mathbf{r}^\infty + (1 - \varrho_\infty) \mathbf{r}^\infty + \mathbf{f}^\infty, \quad (3.10)$$

where

$$|\mathbf{f}^\infty| \leq 3\varepsilon |1 - \varrho_\infty| |\mathbf{A}| |\bar{\mathbf{x}}^\infty| + 6\varepsilon \varrho_\infty (|\mathbf{A}| |\bar{\mathbf{x}}^\infty| + |\mathbf{b}|).$$

From (3.10) it follows that

$$\varrho_\infty (\mathbf{I} - \mathbf{H}_\infty) \mathbf{r}^\infty = -\mathbf{f}^\infty$$

and thus

$$\mathbf{r}^\infty = -\frac{1}{\varrho_\infty} (\mathbf{I} - \mathbf{H}_\infty)^{-1} \mathbf{f}^\infty.$$

If $\sigma(\mathbf{H}_\infty) \ll 1$, we have

$$\begin{aligned} |\mathbf{r}^\infty| &\leq \frac{1}{\varrho_\infty} 3\varepsilon |1 - \varrho_\infty| |\mathbf{A}| |\bar{\mathbf{x}}^\infty| + 6\varepsilon (|\mathbf{A}| |\bar{\mathbf{x}}^\infty| + |\mathbf{b}|) + \mathcal{O}(\varepsilon^2) \\ &\leq 9\varepsilon (|\mathbf{A}| |\bar{\mathbf{x}}^\infty| + |\mathbf{b}|) + \mathcal{O}(\varepsilon^2). \end{aligned}$$

Thus if the sequence computed by Algorithm 2.1 converges, there exists k^* such that $\forall k > k^*$

$$|\mathbf{r}^{(k)}| \leq 9\varepsilon (|\mathbf{A}| |\bar{\mathbf{x}}^{(k)}| + |\mathbf{b}|) + \mathcal{O}(\varepsilon^2).$$

However, if $\sigma(\mathbf{H}_\infty) < 1$ the norm of \mathbf{r}^∞ can be bounded by

$$\|\mathbf{r}^\infty\| \leq \frac{9\varepsilon}{1 - \|\mathbf{H}_\infty\|} (|\mathbf{A}| \|\bar{\mathbf{x}}^\infty\| + \|\mathbf{b}\|) + \mathcal{O}(\varepsilon^2).$$

Hence there exists k^* such that $\forall k > k^*$

$$\|\mathbf{r}^{(k)}\| \leq \frac{9\varepsilon}{1 - \|\mathbf{H}_k\|} (|\mathbf{A}| \|\bar{\mathbf{x}}^{(k)}\| + \|\mathbf{b}\|) + \mathcal{O}(\varepsilon^2). \quad (3.11)$$

REMARK 6. *If Algorithm 2.1 converges and $\mathbf{H}_\infty \ll 1$, it converges to the solution of a linear system that is a perturbation of the original system (1.1). Therefore, if we choose $\eta = 9\varepsilon (|\mathbf{A}| |\bar{\mathbf{x}}^{(k)}| + |\mathbf{b}|)$, the computation will terminate with a vector $\bar{\mathbf{x}}$ that is the solution of*

$$\begin{aligned} (\mathbf{A} + \delta\mathbf{A})\bar{\mathbf{x}} &= \mathbf{b} + \delta\mathbf{b} \\ |\delta\mathbf{A}| &\leq 9\varepsilon |\mathbf{A}|, \quad |\delta\mathbf{b}| \leq 9\varepsilon |\mathbf{b}|. \end{aligned}$$

REMARK 7. If mixed precision is used in Algorithm 2.1 (the factorization is computed using arithmetic of relative precision ε_1 and all the other operations are performed using arithmetic of relative precision $\varepsilon_2 = \varepsilon_1^2$) then, with

$\eta = \frac{9\varepsilon_2}{1 - \|\mathbf{H}_k\|} \left(\|\mathbf{A}\|\bar{\mathbf{x}}^{(k^*)} + \|\mathbf{b}\| \right)$ and provided the condition numbers of \mathbf{A} and \mathbf{M} are less than ε_1^{-1} , the computation will terminate with a vector $\bar{\mathbf{x}}$ that satisfies

$$\|\bar{\mathbf{x}} - \mathbf{x}\| \leq \frac{9\varepsilon_1 \|\mathbf{x}\|}{1 - \|\mathbf{H}_k\|} + \mathcal{O}(\varepsilon_2).$$

3.3. Best achievable accuracy. The inequality (3.11) can give the false impression that we can always achieve convergence after a finite number of steps to a residual lying within the ball of radius ε . When $\sigma(\mathbf{H}_\infty)$ is very close to 1, since $\sigma(\mathbf{H}_k) < \|\mathbf{H}_k\|$ the ratios

$$\frac{\varepsilon}{1 - \|\mathbf{H}_k\|}$$

will increase. Thus the ‘best’ achievable accuracy $\omega(\varepsilon)$ will be

$$\omega(\varepsilon) = \frac{\varepsilon}{1 - \|\mathbf{H}_\infty\|}.$$

In practice, for $\sigma(\mathbf{H}_k)$ less than about 0.9, $\omega(\varepsilon) \approx \varepsilon$. However, for $\sigma(\mathbf{H}_k)$ greater than 0.99, we start to see that $\omega(\varepsilon)$ can be much larger than ε .

4. How to choose the ellipse. In Remark 5 in Section 2, we discussed how Chebyshev acceleration can significantly reduce the number of iterations required for convergence. If the spectral radius $\sigma(\mathbf{M}^{-1}\mathbf{F}) = 1 - \sigma(\mathbf{M}^{-1}\mathbf{A})$ lies between (0, 1), we can scale the ellipse so that

$$\left. \begin{aligned} a &= 1 - \sigma(\mathbf{M}^{-1}\mathbf{A}) \\ b &= t * a \end{aligned} \right\}.$$

Here t must be chosen such that the spectrum is contained within the ellipse. Recall that the number of steps *iter* to reduce the residual by p orders of magnitude is given by equation (2.6). In Figures 4.1 and 4.2, we present the graphs of *iter* for $p = 1$ and 8, respectively. Results are plotted for $t = 0, 0.1, 0.01$ and 1 (iterative refinement) and for $\sigma(\mathbf{M}^{-1}\mathbf{F})$ between 0 and 1. From the graphs, we can forecast that if $\sigma > 0.4$ and t is chosen to be 0.01, Algorithms 2.1 and 2.2 will require significantly fewer steps than iterative refinement. The reduction will potentially be very important for values of $\sigma(\mathbf{M}^{-1}\mathbf{A})$ close to 1, and the best achievable accuracy $\omega(\varepsilon)$ will be rapidly obtained.

We note that after the first step of Algorithm 2.1 (which is equal to the first step of iterative refinement) we can estimate the value of σ to be

$$\sigma(\mathbf{M}^{-1}\mathbf{F}) \approx \rho_1 = \frac{\|\bar{\mathbf{r}}^{(1)}\|}{\|\bar{\mathbf{r}}^{(0)}\|}. \quad (4.1)$$

More generally, the ratio between the computed residuals at the k th and $(k - 1)$ th steps

$$\rho_k = \frac{\|\bar{\mathbf{r}}^{(k)}\|}{\|\bar{\mathbf{r}}^{(k-1)}\|} \quad (4.2)$$

may be used to estimate σ .

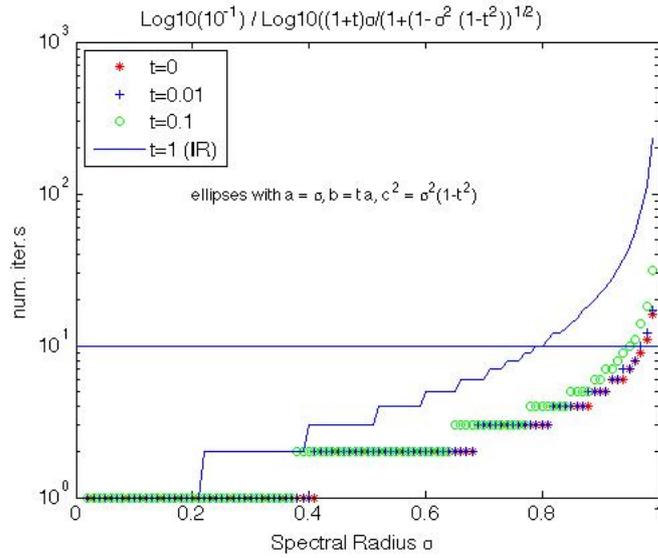


FIG. 4.1. Asymptotic rate of convergence for reducing the initial residual of 10^{-1}

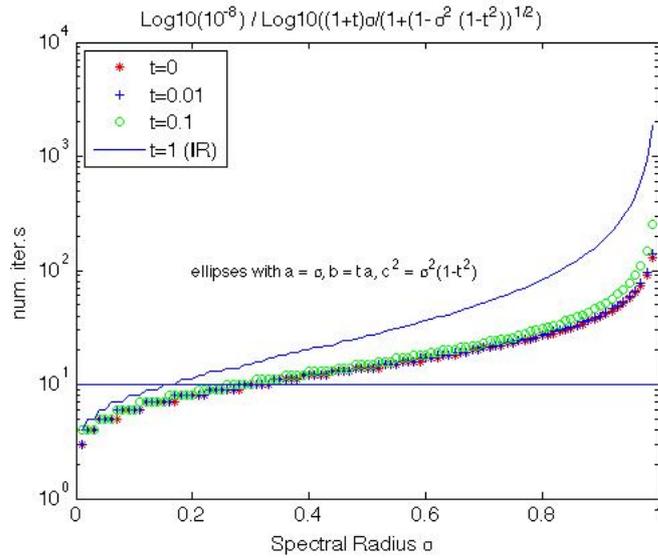


FIG. 4.2. Asymptotic rate of convergence for reducing the initial residual of 10^{-8}

5. Tests on sparse linear systems. In our experiments on sparse systems, we factorize the matrix \mathbf{A} using the single precision version of the new sparse multifrontal solver HSL_MA97 [11], store the computed factors in double precision and then perform refinement using double precision arithmetic. We ran this mixed precision approach on a large number of real symmetric problems taken from the University of Florida Sparse Matrix Collection [4]. For each example, the right-hand side vector \mathbf{b} was generated by setting each component x_i of \mathbf{x} to be a random number in the range $(-1, 1)$. In many cases, only one or two steps of iterative refinement were required to

achieve a scaled residual β (see (1.5)) of less than $5 * 10^{-15}$ (see also the results in [9]). For some ill-conditioned problems, iterative refinement converged to the requested accuracy but required ten or more iterations. These problems are reported on in Table 5.1. Our expectation for these problems is that, with appropriately chosen ellipse parameters a and b , Chebyshev acceleration will be able to reduce the number of iterations.

In our tests, we set $b = 0.01 * a$ and experimented with a range of values of a . For some problems (including HB/nos2 and HB/nos7) setting a equal to the estimate (4.1) gives very good results. In some cases, ρ_k given by (4.2) rapidly converges and setting $a = \rho_k$ for a small value of $k > 1$ minimises the number of iterations. For example, for test example GHS_indef/bratu3d, using ρ_1 requires 18 iterations whereas using ρ_2 reduces the number of iterations to 14 (iterative refinement needs 22 steps). However, we also observed that for some examples, ρ_1 is very small and $a = \rho_1$ gives no improvement on iterative refinement, whereas using a later value of ρ_k can result in savings. This is illustrated by problem GHS_indef/cont-300. In this instance, $\rho_1 = 0.08$ and using $a = \rho_1$ requires 142 iterations (only one less than iterative refinement), while using $\rho_9 = 0.87$ reduces the iteration count to 43.

In Table 5.1, we present results for $a = \rho_1$ and for the a that in our tests resulted in the smallest number of iterations (where applicable, we indicate which ρ_k was used). These results confirm our expectations that Chebyshev acceleration can significantly

TABLE 5.1

Comparison of the number of steps (iter) required by iterative refinement (IR) and Chebyshev accelerated iterative refinement. Chebyshev refinement is run with $a = \rho_1$ and the a that minimizes the number of iterations, denoted by a_{best} (if $a_{best} = \rho_k$, the value of k is given in parentheses).

Problem	IR		Chebyshev IR		
	iter	$a = \rho_1$		$a = a_{best}$	
		iter	ρ_1	iter	a_{best}
HB/nos2	20	12	0.43	12	0.43 (1)
HB/nos7	18	10	0.53	9	0.54
HB/bcsstm27	27	14	0.56	13	0.60
GHS_indef/bratu3d	22	18	0.34	14	0.25 (2)
Cylshell/s3rmt3m1	19	15	0.40	12	0.46 (2)
Cylshell/s3rmq4m1	11	8	0.25	7	0.26
GHS_indef/ncvxbqp1	19	16	0.28	11	0.39 (2)
GHS_indef/cont-300	143	142	0.08	43	0.87 (9)
Oberwolfach/gyro	21	19	0.28	14	0.42 (3)
GHS_indef/sparsine	30	17	0.51	17	0.51 (1)

reduce the number of calls to the solve phase of the direct solver and also that the savings achieved can be very dependent on choosing appropriate ellipse parameters. In many cases, it is worthwhile to perform 2 or 3 steps of iterative refinement to obtain a suitable value for a and then to use Chebyshev acceleration.

For problems with ρ_k almost equal to 1 for k sufficiently large, iterative refinement converges very slowly. For these examples, we are interested in seeing whether

Chebyshev acceleration is able to significantly improve the rate of convergence. We experimented with setting $a = 0.99, 0.9999, 0.999999$ (again, with $b = 0.01*a$) and ran iterative refinement and Chebyshev accelerated iterative refinement for 200 steps. Our findings are reported in Table 5.2. Here we give the initial scaled residual, the scaled residual for iterative refinement ($\beta_{IR}^{(200)}$) and for Chebyshev refinement ($\beta_C^{(200)}(a)$). As a approaches 1, $\beta_C^{(k)}(a)$ reduces more rapidly than $\beta_{IR}^{(k)}$. In particular, after 200 steps, the scaled residual for Chebyshev refinement with $a = 0.999999$ is two orders of magnitude smaller than for iterative refinement. Furthermore, with $a \geq 0.9999$, only 15 steps are required to reduce $\beta_C^{(k)}(a)$ below $\beta_{IR}^{(200)}$.

TABLE 5.2

Comparison of the scaled residuals for iterative refinement and Chebyshev accelerated iterative refinement after 200 steps. These are denoted by $\beta_{IR}^{(200)}$ and $\beta_C^{(200)}(a)$, respectively; $\beta^{(0)}$ is the initial residual.

Problem	$\beta^{(0)}$	$\beta_{IR}^{(200)}$	$\beta_C^{(200)}(a)$		
			$a = 0.99$	0.9999	0.999999
Boeing/crystk03	$4.38 * 10^{-8}$	$2.18 * 10^{-10}$	$3.09 * 10^{-11}$	$3.79 * 10^{-12}$	$2.28 * 10^{-12}$
Oberwolfach/t2dal	$5.97 * 10^{-8}$	$3.09 * 10^{-10}$	$4.39 * 10^{-11}$	$5.33 * 10^{-12}$	$3.18 * 10^{-12}$
Oberwolfach/t3dh_a	$5.31 * 10^{-8}$	$2.74 * 10^{-10}$	$3.89 * 10^{-11}$	$4.68 * 10^{-12}$	$2.75 * 10^{-12}$

So far, our results have shown that, with an appropriate choice of ellipse, Chebyshev acceleration offers advantages over iterative refinement. However, we note that if we choose an ellipse that is too large, the performance of Chebyshev acceleration may be significantly worse than that of iterative refinement. For example, iterative refinement requires 20 steps for problem HB/nos2; with $a = \rho_1 = 0.43$, Chebyshev acceleration reduced this to 12 iterations but other values of a require more iterations. This is illustrated in Table 5.3. As expected, for small a , the performance is as for iterative refinement while for sufficiently large a (in this case, $a > 0.7$), the performance is worse than for iterative refinement. Note that, although $a = \rho_1 = 0.43$ minimises the number of iterations, the precise choice of a is not critical: for a in the approximate range 0.4 to 0.5, Chebyshev acceleration offers worthwhile savings.

TABLE 5.3

The number of iterations required for convergence of Chebyshev accelerated iterative refinement for problem HB/nos2 using a range of values of a . $a = \rho_1$ is in bold.

ρ	0.1	0.2	0.4	0.43	0.5	0.6	0.7	0.8	0.9
<i>iter</i>	20	19	14	12	13	16	19	25	37

6. Conclusions. We have analysed Chebyshev accelerated iterative refinement from the point of view of roundoff and have presented numerical results for sparse linear systems arising from practical applications that support the theory. As our experiments illustrate, using a estimate of the spectral radius obtained by performing a small number of steps of iterative refinement gives good convergence. Moreover, if the ellipse is chosen to be too small (so that it does not contain the complete spectrum), Chebyshev accelerated iterative refinement performs no worse than iterative refinement. We also point out that in all our numerical tests on real symmetric matrices, ellipses with $b < a < 1$ were optimal. This suggests that for such problems

$\sigma(\mathbf{M}^{-1}\mathbf{F})$ generally corresponds either to a real eigenvalue or to one with its real part much larger than its imaginary part. We could not find a theoretical justification of this phenomenon and it is possible to build small artificial examples where this is not the case. Finally, we remark when $\sigma(\mathbf{M}^{-1}\mathbf{F}) > 1$ and iterative refinement does not converge, the use of Flexible GMRES could be a better and more efficient option [1, 2].

REFERENCES

- [1] M. ARIOLI, I. S. DUFF, S. GRATTON, AND S. PRALET, *A note on GMRES preconditioned by a perturbed LDL^T decomposition with static pivoting*, SIAM J. Sci. Comput., 29 (2007), pp. 2024–2044.
- [2] M. ARIOLI AND I. S. DUFF, *Using FGMRES to obtain backward stability in mixed-precision*, Electronic Transactions on Numerical Analysis, 33 (2009), pp. 31–44.
- [3] R. F. BROWN, *A Topological Introduction to Nonlinear Analysis*, Birkhauser, Boston, USA, 1993.
- [4] T. A. DAVIS, *The University of Florida Sparse Matrix Collection*, Technical Report, University of Florida (2007). See <http://www.cise.ufl.edu/~davis/techreports/matrices.pdf>
- [5] G. H. GOLUB AND C. F. V. LOAN, *Matrix Computations*, Johns Hopkins University Press, Baltimore, MD, USA, second ed., 1989.
- [6] M. H. GUTKNECHT AND S. RÖLLIN, *The Chebyshev iteration revisited*, Parallel Computing, 28 (2002), pp. 263–283.
- [7] M. H. GUTKNECHT AND Z. STRAKOŠ, *Accuracy of two thrice-term and three two-term recurrences for Krylov space solvers*, SIAM J. Matrix Anal. Appl., 22 (1) (2000), pp. 213–229.
- [8] L. A. HAGEMAN AND D. M. YOUNG, *Applied Iterative Methods*, Academic Press, New York and London, 1981.
- [9] J. D. HOGG AND J. A. SCOTT, *A fast and robust mixed precision solver for the solution of sparse symmetric linear systems*, ACM Transactions on Mathematical Software, 37 (2010), pp. 17–24.
- [10] ———, *A note on the solve phase of a multicore solver*, Technical Report RAL-TR-2010-007, Rutherford Appleton Laboratory (2010).
- [11] ———, *HSL_MA97: a multifrontal code for sparse symmetric systems*, Technical Report, Rutherford Appleton Laboratory, in preparation (2011).
- [12] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms, Second Edition*, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2002.
- [13] T. A. MANTEUFFEL, *The Tchebychev iteration for nonsymmetric linear systems*, Numerische Mathematik, 28 (1977), pp. 307–327.
- [14] ———, *Adaptive procedure for estimating parameters for the nonsymmetric Tchebychev iteration.*, Numerische Mathematik, 31 (1978), pp. 183–208.
- [15] H. RUTISHAUSER, *Theory of gradient methods*, in: Refined Iterative Methods for Computation of the Solution and the Eigenvalues of Self-Adjoint Boundary Value Problems, Mitt. Inst. Angew. Math. ETH Zurich, Nr. 8, Birkhauser, Basel, 1959, pp. 2449.