

Maintaining convergence properties of BiCGstab methods in finite precision arithmetic

Gerard L.G. Sleijpen and Henk A. van der Vorst

Mathematical Institute, University of Utrecht, P.O. Box 80 010, NL-3508 TA Utrecht, The Netherlands

E-mail: sleijpen@math.ruu.nl; vorst@math.ruu.nl

Received 26 September 1994; revised 16 March 1995

Communicated by C. Brezinski

It is well-known that Bi-CG can be adapted so that hybrid methods with computational complexity almost similar to Bi-CG can be constructed, in which it is attempted to further improve the convergence behavior. In this paper we will study the class of BiCGstab methods.

In many applications, the speed of convergence of these methods appears to be determined mainly by the incorporated Bi-CG process, and the problem is that the Bi-CG iteration coefficients have to be determined from the BiCGstab process. We will focus our attention to the accuracy of these Bi-CG coefficients, and how rounding errors may affect the speed of convergence of the BiCGstab methods. We will propose a strategy for a more stable determination of the Bi-CG iteration coefficients and by experiments we will show that this indeed may lead to faster convergence.

Keywords: Non-symmetric linear systems, iterative solvers, Bi-CG, Bi-CGSTAB, BiCGstab(ℓ).

AMS subject classification: 65F10.

1. Introduction

The BiCGstab methods can be viewed as Bi-CG combined with repeated low degree GMRES processes, like GMRES(1) in Bi-CGSTAB. Therefore, we start with a brief overview of Bi-CG.

Bi-CG [7,12] is an iterative solution method for linear systems

$$Ax = b \tag{1}$$

in which the $n \times n$ matrix A is nonsingular. In typical applications n will be large and A will be sparse. For ease of presentation, we assume A and b to be real.

Starting with an initial guess x_0 for the solution x and a “shadow” residual \tilde{r}_0 (most often one takes $\tilde{r}_0 = r_0$), Bi-CG produces sequences of approximations x_k , residuals r_k , and search directions u_k by

$$u_k = r_k - \beta_k u_{k-1}, \quad x_{k+1} = x_k + \alpha_k u_k, \quad r_{k+1} = r_k - \alpha_k A u_k, \tag{2}$$

where the Bi-CG coefficients α_k and β_k are such that r_k and Au_k are orthogonal to the shadow Krylov subspace $\mathcal{K}_k(A^T; \tilde{r}_0)$.

In principle we are free to select any basis for the shadow Krylov subspace that suits our purposes. We will represent the basis vectors of this subspace in polynomial form.

If (ψ_k) is a sequence of polynomials of degree k with a non-trivial leading coefficient θ_k then the vectors $\psi_0(A^T)\tilde{r}_0, \dots, \psi_{k-1}(A^T)\tilde{r}_0$ form a basis of $\mathcal{K}_k(A^T; \tilde{r}_0)$ and we have (see [24] or [21]):

$$\beta_k = \frac{\theta_{k-1}}{\theta_k} \frac{\rho_k}{\sigma_{k-1}} \quad \text{and} \quad \alpha_k = \frac{\rho_k}{\sigma_k}, \quad (3)$$

where

$$\begin{aligned} \rho_k &:= (r_k, \psi_k(A^T)\tilde{r}_0), \\ \sigma_k &:= (Au_k, \psi_k(A^T)\tilde{r}_0). \end{aligned}$$

In finite precision arithmetic computation the values of the iteration coefficients depend quite critically on the choice of the basis vectors for the shadow Krylov subspace. For example, if we make the straightforward choice $\psi_k(t) = t^k$, then the basis vectors tend to be more and more in the direction of the dominating eigenvector of A^T . Depending on how well the dominant eigenvalue of A^T is separated from the others, this would imply that eventually the new vectors in $\mathcal{K}_k(A; r_0)$ are effectively made only orthogonal with respect to this dominating eigenvector. This then would lead to a new vector r_k , that is almost orthogonal to the k th basis vector for the shadow Krylov subspace, and hence we may expect large relative errors in the new iteration coefficients.

Of course, even if all computational steps are done as accurately as possible (in finite precision), eventually the computed Bi-CG coefficients will differ in all digits from the exact ones. As is well known for the Bi-CG process itself (cf. [1,10,16]), this can be attributed to a global loss of bi-orthogonality. Since Bi-CG seems to work rather well as long as some local bi-orthogonality is maintained (that means that the local bi-orthogonalization is done accurately enough, see also [9]), we expect to recover the convergence behavior of the incorporated Bi-CG process (in the BiCGstab methods) if we compute the iteration coefficients as accurately as possible. Therefore, we want to avoid all additional perturbations that might be introduced by an unfortunate choice of the polynomial process that is carried out on top of the Bi-CG process.

In section 2 we will study the choice of the set of ψ_k 's, and we will identify polynomials that lead to sufficiently stable computation of the Bi-CG iteration coefficients.

The polynomials ψ_k can also be used for a different purpose in the Bi-CG process. Sonneveld [24] was the first to suggest to rewrite the inner products, not only to avoid the operations with A^T , e.g.,

$$\rho_k = (r_k, \psi_k(A^T)\tilde{r}_0) = (\psi_k(A)r_k, \tilde{r}_0) = (\mathbf{r}_k, \tilde{\mathbf{r}}_0), \quad (4)$$

but also to allow the construction of recursions for the vectors $\mathbf{r}_k := \psi_k(A)\mathbf{r}_0$. In this way the polynomials ψ_k can be used for a further reduction of the residual in some norm. In fact, Sonneveld suggested the specific choice $\psi_k = \phi_k$, where $r_k = \phi_k(A)r_0$ (i.e., ϕ_k is the Bi-CG iteration polynomial), and this led to the well-known CGS method [24]. More recently, other hybrid Bi-CG methods have emerged as well. In all these approaches the Bi-CG iteration vectors are not computed explicitly.

In the BiCGstab methods [11,23,25] ψ_k is chosen as the product of low degree minimum residual (like GMRES) polynomials. We will study these choices in section 3 in view of our new insights on the choice of the ψ_k . It will turn out that the quest for a stable computation of the iteration coefficients is not always in concordance with optimal residual reducing properties.

Conventions

- (1) Throughout this paper $\|\cdot\|$ will denote the Euclidean norm.
- (2) For methods to be discussed, the residual r_k at the k th step is in $\mathcal{K}_{k+1}(A; r_0)$ and can be written as a k th degree polynomial ψ_k in A acting on r_0 : $r_k = \psi_k(A)r_0$ and $\psi_k(0) = 1$. In connection with this we will call the polynomial for method M , the M -polynomial associated with A and r_0 (see, e.g., [26]), or the $M(k)$ -polynomial if we want to specify the degree k . In particular the *OR-polynomial* ψ_k^{OR} corresponds to the situation where the residual is orthogonal with respect to the Krylov subspace $\mathcal{K}_k(A; r_0)$ (FOM [18] and GENCG [4] define implicitly OR-polynomials). The *MR-polynomial* ψ_k^{MR} defines the residual that is minimal in the Krylov subspace $\mathcal{K}_{k+1}(A; r_0)$ (as in GMRES [19]).
- (3) We will often use phrases like “reduces the residual” or “small residual”. This will always mean that these residual vectors are reduced (or are small) with respect to the Euclidean norm.

2. The Bi-CG iteration coefficients

To understand why the Bi-CG coefficients can be inaccurate, we concentrate first on the ρ_k (see (4)). Then, as we will show, the effects of inaccurate computation of σ_k can be understood in a similar way. The computed ρ_k will be inaccurate if r_k is nearly orthogonal to $\psi_k(A^T)\tilde{r}_0$. As is well known, this will happen if the incorporated Lanczos process nearly breaks down (i.e. $(\phi_k(A)r_0, \psi_k(A^T)\tilde{r}_0) \approx 0$ for any polynomial ψ_k of exact degree k), and this may be attributed to an unlucky choice of \tilde{r}_0 . This kind of breakdown may be circumvented by so-called look-ahead techniques (see, e.g., [9,17]). However, apart from this a bad choice of ψ_k may lead to a small ρ_k as well. Here, we only consider how this choice of ψ_k causes further instability in the computation of the iteration coefficients. We assume that the Lanczos process itself does not (nearly) break down.

The relative error ϵ_k , due to rounding errors, in ρ_k can be bounded sharply by (see, e.g., [6])

$$\begin{aligned} |\epsilon_k| &\leq 1.01n\bar{\xi} \frac{(|r_k|, |\psi_k(A^T)\tilde{r}_0|)}{|(r_k, \psi_k(A^T)\tilde{r}_0)|} \\ &\leq 1.01n\bar{\xi} \frac{\|r_k\| \|\psi_k(A^T)\tilde{r}_0\|}{|(r_k, \psi_k(A^T)\tilde{r}_0)|} \leq \frac{1.01n\bar{\xi}}{\hat{\rho}_k}, \end{aligned} \quad (5)$$

where

$$\hat{\rho}_k := \frac{|(r_k, \psi_k(A^T)\tilde{r}_0)|}{\|r_k\| \|\psi_k(A^T)\tilde{r}_0\|}, \quad (6)$$

$\bar{\xi}$ is the relative machine precision and n the dimension of the problem. For a small relative error we want to have $\hat{\rho}_k$ (the scaled ρ_k) as large as possible.

Because of the orthogonality of r_k with respect to $\mathcal{K}_k(A^T; \tilde{r}_0)$ it follows that

$$(r_k, \psi_k(A^T)\tilde{r}_0) = \gamma_k(r_k, (A^T)^k\tilde{r}_0), \quad (7)$$

where γ_k is the leading coefficient of ψ_k : $\psi_k(t) = \sum_{j \leq k} \gamma_j t^j$. Hence,

$$\hat{\rho}_k = \frac{|\gamma_k|}{\|\psi_k(A^T)\tilde{r}_0\|} \frac{|(r_k, (A^T)^k\tilde{r}_0)|}{\|r_k\|}. \quad (8)$$

The second quotient in this expression for $\hat{\rho}_k$ does not depend on ψ_k , and therefore, the expression for $\hat{\rho}_k$ is maximal over all polynomials ψ_k with fixed leading coefficient γ_k , when ψ_k is an appropriate multiple of the OR-polynomial ψ_k^{OR} associated with A^T and \tilde{r}_0 , that is

$$\psi_k^{\text{OR}}(A^T)\tilde{r}_0 \perp \mathcal{K}_k(A^T; \tilde{r}_0) \quad \text{and} \quad \psi_k^{\text{OR}}(0) = 1.$$

The appropriate multiple is to take care that the polynomial has leading coefficient γ_k , but since the expression for $\hat{\rho}_k$ is invariant under scaling for ψ_k , we conclude that the OR-polynomial is the polynomial that makes $\hat{\rho}_k$ maximal.

For $\sigma_k = (Au_k, \psi_k(A^T)\tilde{r}_0)$ we can follow the same line of reasoning. Since by construction in Bi-CG the vectors Au_k are orthogonal to lower dimensional shadow Krylov subspaces, it follows that the relative error in σ_k , due to rounding errors can be bounded by an expression which is also minimal for ψ_k^{OR} . Note that the Bi-CG iteration coefficients are formed from ratios of ρ_k 's and σ_k 's, so that errors in each of these add to the inaccuracy in the α_k 's and β_k 's.

Since Bi-CG is designed to avoid all the work for the construction of an orthogonal basis, it would be expensive to construct the ψ_k^{OR} as the basis generating polynomials for the shadow Krylov subspace. In Bi-CG a compromise is made by taking the $\psi_k = \phi_k$ which creates a bi-orthogonal basis. At least for (near) symmetric matrices this is (almost) optimal.

3. Small residuals and accurate Bi-CG coefficients

Now the question arises whether we can select polynomials ψ_k , with $\psi_k(0) = 1$, that satisfy the following requirements:

1. ψ_k leads to sufficiently stable Bi-CG coefficients,
2. ψ_k can be used to further reduce the Bi-CG residual, that is $\mathbf{r}_k = \psi_k(A)r_k$ is (much) smaller in norm than r_k ,
3. the ψ_k can be (implicitly) formed by short recurrences.

3.1. Some choices for the polynomial ψ_k

A choice that comes close, in many relevant cases, to fulfilling all these requirements is the one suggested by Sonneveld [24]: $\psi_k = \phi_k$, which leads to CGS. However, there are two disadvantages associated with this choice. The first is that there is no reason why ϕ_k should lead to a further reduction (and often it does not), the second is that all irregularities in the convergence behavior of Bi-CG are magnified in CGS (although the negative effects of this on the accuracy of the approximated solution can be largely reduced [15,22]).

An obvious alternative is to select ψ_k as the product of k first degree MR (or GMRES) polynomials, which leads to Bi-CGSTAB [25], or as n/ℓ factors of ℓ degree MR-polynomials [21,23].

An obvious problem, when replacing the Bi-CG polynomial by other polynomials which are chosen as to reduce the residual vector, is that these polynomials do not necessarily lead implicitly to an optimal basis for the shadow Krylov subspace.

We will first concentrate on suitable (inexpensive) polynomial methods that help to further reduce the Bi-CG residual. For hybrid Bi-CG methods, where $\mathbf{r}_k = \psi_k^H(A)r_k$, we have that ρ_k is computed as $\rho_k = (\mathbf{r}_k, \tilde{r}_0)$ (cf. (4)), and in finite precision arithmetic there will be a relative error ϵ_k^H in the evaluation of this inner product that can be bounded as

$$|\epsilon_k^H| \leq 1.01n\bar{\xi} \frac{(|\mathbf{r}_k|, |\tilde{r}_0|)}{|(\mathbf{r}_k, \tilde{r}_0)|} \leq 1.01n\bar{\xi} \frac{\|\mathbf{r}_k\| \|\tilde{r}_0\|}{|(\mathbf{r}_k, \tilde{r}_0)|} \leq \frac{1.01n\bar{\xi}}{\hat{\rho}_k^H}, \quad (9)$$

where

$$\hat{\rho}_k^H := \frac{|(\mathbf{r}_k, \tilde{r}_0)|}{\|\mathbf{r}_k\| \|\tilde{r}_0\|}. \quad (10)$$

For similar reasons as in section 2, we have

$$\hat{\rho}_k^H = \frac{|\gamma_k|}{\|\psi_k^H(A)r_k\|} \frac{|(A^k r_k, \tilde{r}_0)|}{\|\tilde{r}_0\|}, \quad (11)$$

where γ_k is the leading coefficient of ψ_k^H . Again, the OR-polynomial ψ_k^{OR} (associated with A and r_k), which minimizes $\|\psi_k^H(A)r_k\|/|\gamma_k|$, would lead to a maximal value for $\hat{\rho}_k^H$. Of course, there is no guarantee or reason why this polynomial should also

maximize the expression for $\hat{\rho}_k$ in the Bi-CG representation for the inner product, but given the fact that we want the polynomial to act directly on r_k this is the best we can do.

In the BiCGstab methods the ψ_k^H is chosen as a product of MR-polynomials, for obvious reasons. The MR-polynomial (associated with A and r_k) minimizes $\|r_k\| = \|\psi_k^H(A)r_k\|$ over all polynomials ψ_k^H for which $\psi_k^H(0) = 1$, and hence seems to be more appropriate for creating small residuals r_k . But it would be too expensive to perform k steps of FOM or GMRES in order to compute $r_k = \psi_k^H(A)r_k$ (using r_k as initial residual): we not only strive for a large reduction but also for inexpensive steps. The product of MR(1)-polynomials (i.e. MR-polynomials of degree 1) as in Bi-CGSTAB is a compromise between the wish for small residuals and inexpensive steps. In view of our discussion above, however, we might also consider OR(1)-polynomials in order to achieve more stability in some cases, giving up some of the reduction.

With OR(1)-polynomials $(1 - \omega_k \lambda)$ for which $(I - \omega_k A)\hat{r} \perp \hat{r}$, where $\hat{r} := \psi_{k-1}^H(A)r_k$, we compromise between (locally) more accurate coefficients and inexpensive steps. Although these polynomials occasionally cure stagnation of Bi-CGSTAB (see also [5]) they also may amplify residuals, which again leads to inaccurate approximations (as explained by (7) in [23]) or even to overflow.

We will now show that, if the angle between \hat{r} and $A\hat{r}$ is more than 45° (i.e. $|\langle \hat{r}, A\hat{r} \rangle| \leq \frac{1}{2}\sqrt{2}\|\hat{r}\|\|A\hat{r}\|$) then the OR(1)-polynomial locally amplifies the residual, while the MR(1)-polynomial leads to a smaller value for $\hat{\rho}_k^H$ (cf. (11)).

This property is easily explained with figure 1, where \hat{s} is a scalar multiple of $A\hat{r}$, scaled such that $\|\hat{s}\| = \|\hat{r}\|$. The residual r^{OR} is obtained by applying the OR(1)-polynomial to \hat{r} , r^{MR} results from the MR(1)-polynomial. Clearly, in the situation as sketched in the figure, $\|r^{MR}\| < \|\hat{r}\| < \|r^{OR}\|$, while scaling the polynomials such that the leading coefficient is identical (in the figure $\|\hat{r}\|/\|A\hat{r}\|$) changes the order: $\|r^{OR}\|/|1/\hat{\omega}| < \|\hat{s}\| < \|r^{MR}\|/|\hat{\omega}|$.

We have to be careful with such amplifications when they are extremely large or when they occur in a consecutive number of iteration steps (as will be the case in a stagnation phase of the process). In such cases, any of the two choices may slow down the convergence: the OR(1)-polynomials because of amplifying $\|r_k\|$; the MR(1)-polynomials because of shrinking $\hat{\rho}_k^H$ and thus affecting the accuracy of the Bi-CG coefficients. Apparently, Bi-CGSTAB may also be expected to converge poorly if in a number of consecutive steps the angle between \hat{r} and $A\hat{r}$ is more than 45° . Especially when, in Bi-CGSTAB, GMRES(1) stagnates in a consecutive number of steps, it is important to have accurate Bi-CG coefficients, because, in such a case any convergence is to be expected only from the Bi-CG part of Bi-CGSTAB. Unfortunately, this is precisely the situation where the MR(1)-polynomials spoil the accuracy of the coefficients, while the OR(1)-polynomials spoil the accuracy of the approximations or lead to overflow. Therefore, we have to find a cure by other modifications.

We will suggest other choices for ω_k (as in (28), see also [5]) which occasionally cure these problems. However, they often cannot completely prevent poor

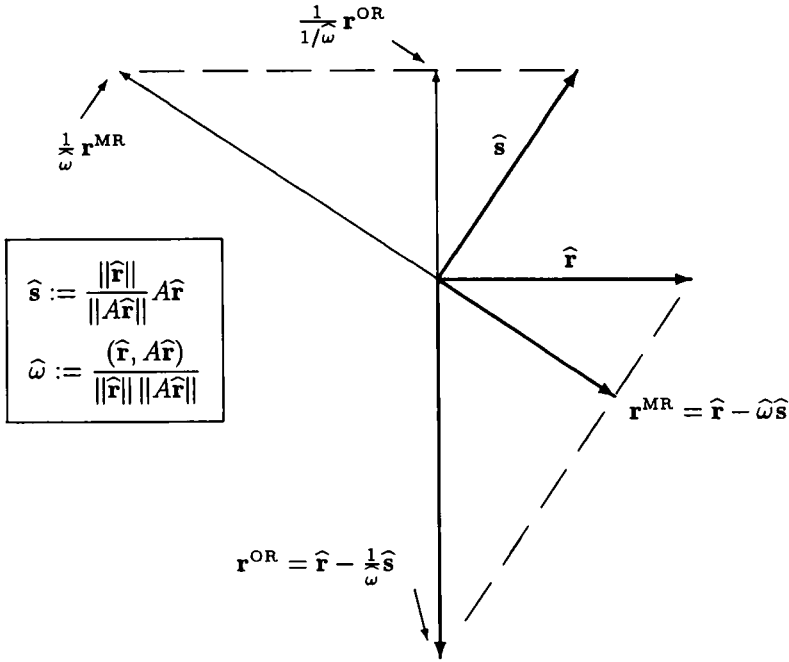


Figure 1. Amplification effects of GMRES(1) and FOM(1).

convergence. One reason is that, with these first degree factors, we are (implicitly) building a power basis for the shadow Krylov subspace (if the ω_k are close to each other), and we have seen in section 1 that this is highly undesirable.

We may expect better convergence results by performing a composition of ℓ steps, using minimizing polynomials of degree ℓ , with $\ell > 1$: that is, by taking ψ_k as a product of MR(ℓ)-polynomials as in BiCGstab(ℓ) (see [11,21,23]), provided that the MR(ℓ)-polynomials lead to significant reductions. We will consider this approach in much more detail.

In BiCGstab(ℓ) the polynomial ψ_k is constructed as a product of polynomials of degree ℓ : for $k = m\ell$, $\psi_k = p_{m-1} \cdots p_0$ where p_j is of degree ℓ . To investigate what properties these polynomial factors p_j should have, we consider $k = m\ell$, concentrate on p_m , and we define

$$\hat{\mathbf{r}} := \psi_k(A)\phi_{k+\ell}(A)r_0 = \psi_k(A)r_{k+\ell}. \quad (12)$$

In BiCGstab(ℓ) the vector $\hat{\mathbf{r}}$ is computed explicitly in the Bi-CG part of the algorithm. The new residual $\mathbf{r}_{k+\ell}$ will be $\mathbf{r}_{k+\ell} := p_m(A)\hat{\mathbf{r}}$ (for implementational details, see [21,23]).

3.2. Minimizing polynomials of low degree

For the derivations and results in this section we assume exact arithmetic. However, we expect that they also have some validity in finite precision arithmetic, and as we will see in section 4, the experiments do largely confirm our expectations.

As before, we wish to maximize $\hat{\rho}_k^H$ for polynomials p_m of degree ℓ for which $p_m(0) = 1$. As in section 2, this means that we have to restrict ourselves to polynomials with a fixed leading coefficient γ_ℓ .

In order to make our formulas more readable, we define

$$|\mathbf{r}| := \frac{\|\mathbf{r}\|}{|\gamma_\ell|}, \quad \text{where} \quad \mathbf{r} := p(A)\hat{\mathbf{r}} \quad \text{and} \quad p(t) = \sum_{j=0}^{\ell} \gamma_j t^j. \quad (13)$$

In order to have $\hat{\rho}_{k+\ell}^H$,

$$\hat{\rho}_{k+\ell}^H = \frac{|(\mathbf{r}, \tilde{\mathbf{r}}_0)|}{\|\mathbf{r}\| \|\tilde{\mathbf{r}}_0\|}, \quad (14)$$

(approximately) maximal, we should have that

$$|p(A)\hat{\mathbf{r}}| \quad \text{is (approximately) minimal.}$$

As is well known \mathbf{r}^{MR} solves

$$\min\{\|\mathbf{r}\| \mid \mathbf{r} = p(A)\hat{\mathbf{r}}, \deg(p) \leq \ell, p(0) = 1\}$$

if and only if

$$\mathbf{r}^{\text{MR}} \perp A\hat{\mathbf{r}}, A^2\hat{\mathbf{r}}, \dots, A^\ell\hat{\mathbf{r}}. \quad (16)$$

In a similar way, one can show that \mathbf{r}^{OR} solves

$$\min\{|\mathbf{r}| \mid \mathbf{r} = p(A)\hat{\mathbf{r}}, \deg(p) \leq \ell, p(0) = 1\}$$

(cf. (15)) if and only if

$$\mathbf{r}^{\text{OR}} \perp \hat{\mathbf{r}}, A\hat{\mathbf{r}}, \dots, A^{\ell-1}\hat{\mathbf{r}}. \quad (17)$$

The residual \mathbf{r}^{MR} is the ℓ th residual of a minimal residual method (as GMRES), and \mathbf{r}^{OR} is the ℓ th residual of an orthogonal residual method (as FOM), each for the problem $Ax = b$ with the same initial residual $\hat{\mathbf{r}}$.

The following theorem compares how good \mathbf{r}^{MR} is in maximizing $\hat{\rho}_k^H$ and how well \mathbf{r}^{OR} helps to reduce $\hat{\mathbf{r}}$. A method like ORTHODIR produces explicitly an orthonormal basis for which theorem 3.1 can be applied.

We will call a sequence of vectors $\hat{\mathbf{r}}_1, \dots, \hat{\mathbf{r}}_\ell$ a *Krylov basis* for $\mathcal{K}_\ell(A; \hat{\mathbf{r}}_1)$ if $\hat{\mathbf{r}}_1, \dots, \hat{\mathbf{r}}_j$ and $\hat{\mathbf{r}}_1, \dots, A^{j-1}\hat{\mathbf{r}}_1$ span the same space for each $j = 1, \dots, \ell$.

Theorem 3.1

Let $\hat{\mathbf{r}}_1, \dots, \hat{\mathbf{r}}_{\ell-1}$ be an orthonormal Krylov basis for $\mathcal{K}_{\ell-1}(A; A\hat{\mathbf{r}})$. Let the vectors $\tilde{\mathbf{r}}_0$ and $\tilde{\mathbf{r}}_\ell$ be obtained by orthogonalizing $\hat{\mathbf{r}}$ and $A^\ell\hat{\mathbf{r}}$ with respect to $\mathcal{K}_{\ell-1}(A; A\hat{\mathbf{r}})$:

$$\tilde{\mathbf{r}}_0 := \hat{\mathbf{r}} - \sum_{j=1}^{\ell-1} (\hat{\mathbf{r}}, \hat{\mathbf{r}}_j) \hat{\mathbf{r}}_j \quad \text{and} \quad \tilde{\mathbf{r}}_\ell := A^\ell\hat{\mathbf{r}} - \sum_{j=1}^{\ell-1} (A^\ell\hat{\mathbf{r}}, \hat{\mathbf{r}}_j) \hat{\mathbf{r}}_j.$$

Let ϱ be defined as

$$\varrho := \frac{(\tilde{\mathbf{r}}_\ell, \tilde{\mathbf{r}}_0)}{\|\tilde{\mathbf{r}}_\ell\| \|\tilde{\mathbf{r}}_0\|}. \quad (18)$$

Then

$$\mathbf{r}^{\text{MR}} = \tilde{\mathbf{r}}_0 - \varrho \frac{\|\tilde{\mathbf{r}}_0\|}{\|\tilde{\mathbf{r}}_\ell\|} \tilde{\mathbf{r}}_\ell \quad \text{and} \quad \mathbf{r}^{\text{OR}} = \tilde{\mathbf{r}}_0 - \frac{1}{\varrho} \frac{\|\tilde{\mathbf{r}}_0\|}{\|\tilde{\mathbf{r}}_\ell\|} \tilde{\mathbf{r}}_\ell, \quad (19)$$

$$\|\mathbf{r}^{\text{MR}}\| = \sqrt{1 - \varrho^2} \|\tilde{\mathbf{r}}_0\| \quad \text{and} \quad \|\mathbf{r}^{\text{OR}}\| = \frac{\sqrt{1 - \varrho^2}}{|\varrho|} \|\tilde{\mathbf{r}}_0\|, \quad (20)$$

$$|\mathbf{r}^{\text{MR}}| = \frac{\sqrt{1 - \varrho^2}}{|\varrho|} \|\tilde{\mathbf{r}}_\ell\| \quad \text{and} \quad |\mathbf{r}^{\text{OR}}| = \sqrt{1 - \varrho^2} \|\tilde{\mathbf{r}}_\ell\|. \quad (21)$$

Proof

One may easily verify that $s := \tilde{\mathbf{r}}_0 - \gamma \tilde{\mathbf{r}}_\ell = p(A)\hat{\mathbf{r}}$ for some polynomial p of degree ℓ with leading coefficient γ and $p(0) = 1$. Moreover, since $\tilde{\mathbf{r}}_0$ and $\tilde{\mathbf{r}}_\ell$ are orthogonal to $A\hat{\mathbf{r}}, \dots, A^{\ell-1}\hat{\mathbf{r}}$, the vector s is also orthogonal to these vectors as well.

Define $\nu := \|\tilde{\mathbf{r}}_0\|/\|\tilde{\mathbf{r}}_\ell\|$. With $\gamma = \varrho\nu$, s is also orthogonal to $\tilde{\mathbf{r}}_\ell$, and consequently, orthogonal to $A^\ell\hat{\mathbf{r}}$. By (16), $s = \mathbf{r}^{\text{MR}}$ for this $\gamma = \varrho\nu$. Similarly, the choice $\gamma = \nu/\varrho$ makes s orthogonal to $\tilde{\mathbf{r}}_0$, which implies that $s = \mathbf{r}^{\text{OR}}$ (cf. (17)), and this completes the proof of (19).

The expressions in (20) for the norms of the optimal residuals follow from Pythagoras' theorem: $\tilde{\mathbf{r}}_\ell \perp \mathbf{r}^{\text{MR}}$ and $\tilde{\mathbf{r}}_0 \perp \mathbf{r}^{\text{OR}}$. Combining these expressions with the values for the leading coefficient γ ($\gamma = \varrho\nu$, $\gamma = \nu/\varrho$, respectively) gives (21). \square

The vector $\tilde{\mathbf{r}}_0$ in the theorem is precisely the residual in the $(\ell - 1)$ th step of a minimal residual method. Therefore, if ζ is the residual reduction in step ℓ by this MR method (i.e. $\zeta := \|\mathbf{r}_\ell^{\text{MR}}\|/\|\mathbf{r}_{\ell-1}^{\text{MR}}\|$) then $\varrho = \sqrt{1 - \zeta^2}$ and we have the following corollary.

Corollary 3.2

$$\frac{\|\mathbf{r}_\ell^{\text{MR}}\|}{\|\mathbf{r}_\ell^{\text{OR}}\|} = |\varrho| \quad \text{and} \quad \frac{|\mathbf{r}_\ell^{\text{MR}}|}{|\mathbf{r}_\ell^{\text{OR}}|} = \frac{1}{|\varrho|}, \quad (22.a)$$

where

$$\sqrt{1 - \varrho^2} = \frac{\|\mathbf{r}_\ell^{\text{MR}}\|}{\|\mathbf{r}_{\ell-1}^{\text{MR}}\|}. \quad (22.b)$$

Property (20) can also be found in [3,26] (there, the authors focussed on GMRES while our formulation follows the ORTHODIR approach).

3.3. Global effects

The concept of Lanczos breakdown is well known. We will see that it can be translated in terms of angles between the Krylov subspaces and their “shadow subspaces”. It is less obvious what *near-breakdown of the Lanczos process* could

mean. We will say that no near-breakdown takes place if the angle between the Krylov subspace $\mathcal{K}_k := \mathcal{K}_k(A; r_0)$ and the shadow Krylov subspace $\tilde{\mathcal{K}}_k := \mathcal{K}_k(A^T; \tilde{r}_0)$ is uniformly (in k) sufficiently smaller than $\pi/2$:

$$\inf_k \cos \angle(\mathcal{K}_k, \tilde{\mathcal{K}}_k) = \inf_k \left(\inf_{v \in \mathcal{K}_k} \sup_{\tilde{v} \in \tilde{\mathcal{K}}_k} \frac{|(v, \tilde{v})|}{\|v\| \|\tilde{v}\|} \right) \geq \kappa > 0. \quad (23)$$

In particular, we then have that

$$\inf_k \sup \left\{ \frac{|(r_k, \tilde{v})|}{\|r_k\| \|\tilde{v}\|} \mid \tilde{v} \in \tilde{\mathcal{K}}_{k+1} \right\} > 0, \quad (24)$$

and, since $r_k \perp \tilde{\mathcal{K}}_k$, we see that $(r_k, \tilde{r}_k) \neq 0$ (no-breakdown of the Lanczos process).

For hybrid Bi-CG methods as CGS, Bi-CGSTAB and BiCGstab(ℓ) we translate property (24) as

$$\kappa_k := \sup \left\{ \frac{|(\psi(A)r_k, \tilde{r}_0)|}{\|\psi(A)r_k\| \|\tilde{r}_0\|} \mid \psi \text{ pol. of degree } k \right\}, \quad \kappa_k \geq \kappa > 0, \text{ for all } k. \quad (25)$$

We may expect to obtain the most accurate ρ_k by using the k th OR-polynomial ϕ_k^{OR} of k steps of an OR method with initial residual r_k . Then, ρ_k may expect to be endowed with a relative error of size $n\bar{\xi}/\kappa_k$. In practise, if we use another polynomial ψ_k of degree k , we may expect an error in ρ_k that is larger than $n\bar{\xi}/\kappa_k$ by a multiplicative factor

$$\frac{|\psi_k(A)r_k|}{|\phi_k^{\text{OR}}(A)r_k|}. \quad (26)$$

If this factor is large, say $\geq \kappa/(n\bar{\xi})$, the scalar ρ_k and hence the Bi-CG coefficients can not be expected to have any correct digit. It is difficult to analyze this factor as a function of k , since the initial residual r_k changes for each k . Clearly, since $|\phi_{k+\ell}^{\text{OR}}(A)r_{k+\ell}| \leq |p_m^{\text{OR}}(A)\psi_k(A)r_{k+\ell}| = |p_m^{\text{OR}}(A)\hat{r}|$ and $\psi_{k+\ell} = p_m\psi_k = p_m \cdots p_0$,

$$\frac{|\psi_{k+\ell}(A)r_{k+\ell}|}{|\phi_{k+\ell}^{\text{OR}}(A)r_{k+\ell}|} = \frac{1}{\varrho_m} \frac{|p_m^{\text{OR}}(A)\hat{r}|}{|\phi_{k+\ell}^{\text{OR}}(A)r_{k+\ell}|} \geq \frac{1}{\varrho_m}, \quad (27)$$

where

$$\varrho_m := \frac{|p_m^{\text{OR}}(A)\hat{r}|}{|p_m(A)\hat{r}|}.$$

In our discussion, we assume that the factor in (26) has at least the order of magnitude of the product $\prod_{j=0}^m 1/\varrho_j$ (with $k = m\ell$) of the factors $1/\varrho_j$ (in (27)) per sweep. The conclusion that we should avoid that $\prod 1/\varrho_j \geq \kappa/(n\bar{\xi})$ seems to be supported by results from numerical experiments (cf. section 4).

3.4. Discussion

The OR(ℓ)-polynomial is the best choice for obtaining more accurate Bi-CG coefficients, but the MR(ℓ)-polynomial will do almost as well if there is a significant

error reduction at the ℓ th step of GMRES (with initial residual $\tilde{\mathbf{r}}$) (cf. (21) and (22)). In that case the effect of $\text{OR}(\ell)$ is practically equivalent with the effect of $\text{MR}(\ell)$ (cf. (20) and (22)), and $\text{MR}(\ell)$ -polynomials may be slightly preferable.

However, if GMRES (with initial residual $\tilde{\mathbf{r}}$) does not reduce the residual well at the ℓ th step then $|\varrho| \ll 1$ (cf. (22)), and, hence, the $\text{MR}(\ell)$ -polynomial may be expected to lead to an inaccurate $\rho_{k+\ell}$, while the $\text{OR}(\ell)$ -polynomial will enlarge the residual significantly. Likewise, if in a consecutive number of sweeps with $\text{BiCGstab}(\ell)$ the factor $|\varrho|$ is less than $\sqrt{2}/2$, say, then the choice of $\text{MR}(\ell)$ -polynomials may lead to inaccurate Bi-CG coefficients as well, because of an accumulation of rounding errors in ρ_k (and σ_k). On the other hand, the $\text{OR}(\ell)$ -polynomials may lead to unacceptable large residuals, or even to overflow, after a consecutive number of amplifications of the residual.

We propose to make a compromise by choosing some intermediate between the $\text{OR}(\ell)$ - and $\text{MR}(\ell)$ -polynomial in case of a poor reduction by the latter one. This choice is inspired by equation (19):

$$\mathbf{r}_{k+\ell} = \tilde{\mathbf{r}}_0 - \hat{\gamma}_k \frac{\|\tilde{\mathbf{r}}_0\|}{\|\tilde{\mathbf{r}}_\ell\|} \tilde{\mathbf{r}}_\ell, \quad \text{where } \hat{\gamma}_k := \frac{\varrho}{|\varrho|} \max(|\varrho|, \Omega), \quad (28)$$

with $\tilde{\mathbf{r}}_0$, $\tilde{\mathbf{r}}_\ell$, and ϱ as in theorem 3.1 and $\Omega \in [0, \infty)$. In section 3.4.1 we will present some strategies for the choice of Ω and we will comment on computational details. However, although this approach often helps to cure our instability problems, it is not a panacea which can always circumvent an occasionally poor convergence of $\text{BiCGstab}(\ell)$. In such a situation, increasing the value of ℓ may be an alternative.

We may expect better convergence for $\text{BiCGstab}(\ell)$ if we increase the value for ℓ , and we will argue why. It helps for our discussion to compare $\text{BiCGstab}(\ell)$ with BiCGSTAB ($=\text{BiCGstab}(1)$), that is, we compare 1 sweep of $\text{BiCGstab}(\ell)$ with ℓ sweeps of Bi-CGSTAB . In ℓ sweeps of Bi-CGSTAB , ℓ times a $\text{MR}(1)$ -polynomial is applied (each time for a different starting vector: $\psi_{k+j}(A)r_{k+j+1}$, $j = 0, \dots, \ell - 1$). By selecting higher degree $\text{MR}(\ell)$ -polynomials, as in $\text{BiCGstab}(\ell)$, we hope to profit for two different reasons:

(1) One sweep of $\text{GMRES}(\ell)$ may be expected to result in a better residual reduction than ℓ steps of $\text{GMRES}(1)$,

Remark

In the case that $\text{GMRES}(1)$ reduces well, we do not have to fear a loss of speed of convergence due to inaccurate Bi-CG coefficients. However, since GMRES may accelerate (superlinear convergence behavior), we may expect to obtain a smaller residual \mathbf{r}_k with $\text{BiCGstab}(\ell)$ than with Bi-CGSTAB .

(2) ℓ steps of $\text{GMRES}(1)$ contribute ℓ times to a decrease of $\hat{\rho}_k$ (hence contributing ℓ times to increasingly larger rounding errors in ρ_k), while one sweep of $\text{GMRES}(\ell)$ contributes only once; the decreasing effect in each single step of ℓ

steps of GMRES(1) may be expected to be comparable or worse than the effect of only one sweep with GMRES(ℓ).

Remark

If the ℓ th step of GMRES does not reduce well (i.e. we have a small $|\varrho|$) then the MR(ℓ)-polynomial “amplifies” the inaccuracy on $\rho_{k+\ell}$ by $1/|\varrho|$ (in comparison with the OR(ℓ)-polynomial). In such a situation we may not expect to obtain a significant reduction by any of the steps of GMRES(1). It is even worse: in ℓ steps of Bi-CGSTAB, we should expect an “amplification” by the GMRES(1) steps in the inaccuracy of $\rho_{k+\ell}$ by a factor like $(1/|\varrho|)^\ell$ or more. More specifically, in a stagnation phase of the iterative method, it is more likely to have rather accurate Bi-CG coefficients, when we use BiCGstab(ℓ), than with Bi-CGSTAB.

In our discussion, the expected superlinear convergence behavior of GMRES plays a rather important role. However, as is well known, GMRES may as well converge slower or even stagnate for quite a while in any phase of the iteration process. In order to profit from a possible good reduction of the MR($\ell - 1$)-polynomial, in cases where the MR(ℓ)-polynomial gives only a poor additional reduction, we may use the modification as suggested in (28).

3.4.1. Computational details

The computation of the leading coefficient and ϱ , as suggested in theorem 3.1 and formula (28), can be done with relatively inexpensive operations involving ℓ -vectors only.

If $\hat{\mathbf{r}}_0 = \hat{\mathbf{r}}, \hat{\mathbf{r}}_1, \dots, \hat{\mathbf{r}}_\ell$ is a Krylov basis of $\mathcal{K}_{\ell+1}(A; \hat{\mathbf{r}})$ and $\mathbf{R} := [\hat{\mathbf{r}}_0 \dots \hat{\mathbf{r}}_\ell]$ then $\hat{\mathbf{r}}_0 = \mathbf{R}\tilde{\eta}_0$ and $\hat{\mathbf{r}}_\ell = \mathbf{R}\tilde{\eta}_\ell$ for some $\tilde{\eta}_0, \tilde{\eta}_\ell \in \mathbb{R}^{\ell+1}$. Consider the $(\ell + 1) \times (\ell + 1)$ -matrix $V := \mathbf{R}^T \mathbf{R}$, the inner product $\langle \tilde{\eta}, \tilde{\mu} \rangle := \tilde{\mu}^T V \tilde{\eta}$ and norm $|\tilde{\eta}| := \sqrt{\langle \tilde{\eta}, \tilde{\eta} \rangle}$. Then (cf. (28))

$$\mathbf{r}_{k+\ell} = \mathbf{R} \left(\tilde{\eta}_0 - \hat{\gamma}_k \frac{|\tilde{\eta}_0|}{|\tilde{\eta}_\ell|} \tilde{\eta}_\ell \right), \quad (29)$$

with

$$\hat{\gamma}_k := \frac{\varrho}{|\varrho|} \max(|\varrho|, \Omega), \quad \varrho = \frac{\langle \tilde{\eta}_\ell, \tilde{\eta}_0 \rangle}{|\tilde{\eta}_\ell| |\tilde{\eta}_0|},$$

and Ω some scalar in $[0, \infty)$.

With $\Omega = 0$ we have the MR(ℓ)-polynomial which gives optimal reduction of $\hat{\mathbf{r}}$ with respect to $\|\cdot\|$. For very small ϱ or for a consecutive number of non-large ϱ (say, $|\varrho| < 1/2$), $\Omega = 0$ may result in inaccurate Bi-CG coefficients, since, as compared with the OR-polynomial, the MR-polynomial amplifies $\hat{\mathbf{r}}$ with respect to $\|\cdot\|$ by $1/|\varrho|$ (i.e. $|\mathbf{r}^{\text{MR}}| = |\mathbf{r}^{\text{OR}}|/|\varrho|$): if in m sweeps of BiCGstab(ℓ) each ϱ is $|\varrho| \leq \beta < 1$ we may expect an amplification of the relative rounding error by at most $(1/\beta)^m$ (comparing local effects of OR- and MR-polynomials and assuming that local differences accumulate).

With $\Omega = 1/|\varrho|$ we have the OR(ℓ)-polynomial. Although this polynomial gives optimal reduction of $\hat{\mathbf{r}}$ with respect to $|\cdot|$, it may lead to large residuals.

With $\Omega > 0$ we may avoid large amplifications of the inaccuracy in the Bi-CG coefficients. The choice $\Omega = 1$ distributes negative effects of small ϱ equally amongst residual reduction (amplifying $\|\mathbf{r}^{\text{MR}}\|$ by at most $\sqrt{2}$: $\|\mathbf{r}_{k+\ell}\| \leq \sqrt{2}\|\mathbf{r}^{\text{MR}}\|$) and rounding error amplification (amplifying $|\mathbf{r}^{\text{OR}}|$ by at most $\sqrt{2}$: $|\mathbf{r}_{k+\ell}| \leq \sqrt{2}|\mathbf{r}^{\text{OR}}|$). For $\ell = 1$, this choice amplifies both $\|\hat{\mathbf{r}}\|$ and $|A\hat{\mathbf{r}}|$ whenever $|\varrho| < \frac{1}{2}$ (that is, $\|\mathbf{r}_{k+1}\| > \|\hat{\mathbf{r}}\|$ and $|\mathbf{r}_{k+1}| > |A\hat{\mathbf{r}}|$), and reduces both quantities for other values of ϱ . A significant amplification of the rounding error in a few steps need not worry us as long as the coefficients are still rather accurate (in 16 digit arithmetic a cumulative amplification by, say, 10^8 may still be very acceptable). Therefore, the best choice of Ω will depend on the length of the stagnation phase: using (28), we expect a cumulative rounding error amplification by at most $(1 + \Omega^{-2})^{L/(2\ell)}$, where $2L$ is number of matrix-vector multiplications in the stagnation phase (again assuming that the local differences accumulate). Therefore, for larger ℓ a smaller Ω may be acceptable. Moreover, in the upper bound, we did not take into account the reducing effect of the OR-polynomial itself: the amplification factor $1/|\varrho|$ of the MR-polynomial may be harmless if $|\mathbf{r}^{\text{OR}}|$ is small. For instance, for $\ell = 1$ and $|\varrho| < \sqrt{2}/2$, we have that $|\mathbf{r}^{\text{MR}}| < |A\hat{\mathbf{r}}|$.

3.5. Minimizing polynomials and σ_k

Following the arguments for ρ_k , using the concept of near-breakdown of the LU-decomposition, we can see that, for accurate σ_k , the $|A\psi_k(A)u_k|$ should be as small as possible. In the m th sweep, where $\psi_k = p_{m-1} \cdots p_0$ is given ($k = m\ell$) and p_m is constructed, $|Ap_m(A)\hat{\mathbf{u}}|$, with $\hat{\mathbf{u}} := \psi_k(A)u_{k+\ell}$, should be as small as possible. Unfortunately, the construction of ψ_k is linked to the residual r_k . Note that we do not have this problem in the Bi-CG process. For Bi-CG, the polynomial ψ_k that will give the most accurate coefficients ρ_k and σ_k is linked to the initial shadow residual \tilde{r}_0 (see section 2): $|\psi_k(A^T)\tilde{r}_0|$ should be minimal.

The following observation links the scaled σ_k (cf. (30)) to the scaled ρ_k and the residual r_k and gives some theoretical support to the strategy of concentrating on non-small scaled ρ_k only.

Since $Au_k = (1/\alpha_k)(r_k - r_{k+1})$ and $(\psi_k(A)r_{k+1}, \tilde{r}_0) = 0$ we have,

$$\frac{1}{\hat{\sigma}_k^{\text{H}}} := \frac{\|A\psi_k(A)u_k\| \|\tilde{r}_0\|}{|(A\psi_k(A)u_k, \tilde{r}_0)|} = \frac{\|\psi_k(A)r_k - \psi_k(A)r_{k+1}\| \|\tilde{r}_0\|}{|(\psi_k(A)r_k, \tilde{r}_0)|}. \quad (30)$$

Hence, with $\mathbf{r}_k := \psi_k(A)r_k$, we can bound the scaled σ_k by the scaled ρ_k and the growth of the residual in step k :

$$\frac{1}{\hat{\sigma}_k^{\text{H}}} \leq \frac{\|\mathbf{r}_k\| \|\tilde{r}_0\|}{|(\mathbf{r}_k, \tilde{r}_0)|} \left(1 + \frac{\|\psi_k(A)r_{k+1}\|}{\|\psi_k(A)r_k\|} \right) = \frac{1}{\hat{\rho}_k^{\text{H}}} \left(1 + \frac{\|\psi_k(A)r_{k+1}\|}{\|\psi_k(A)r_k\|} \right). \quad (31)$$

If our strategy to choose the polynomial ψ_k prevents $\hat{\rho}_k^{\text{H}}$ to become too small and the residuals cannot grow much in one step then our strategy works for σ_k as well.

One can show that (see [14]),

$$\frac{\|r_{k+1}\|}{\|r_k\|} \leq \frac{\theta_+^{(k)}}{\theta_-^{(k)}} \frac{1}{\lambda_k}, \quad \text{with } \lambda_k := \frac{|(r_{k+1}, \tilde{r}_{k+1})|}{\|r_{k+1}\| \|\tilde{r}_{k+1}\|}, \quad (32)$$

and $\theta_-^{(k)}$ and $\theta_+^{(k)}$ are maximal and minimal, respectively, such that,

$$\theta_-^{(k)} \leq \theta_j^{(k)} \leq \theta_+^{(k)} \quad \text{for all } j = 0, \dots, k, \quad (33)$$

where the $\theta_j^{(k)}$ are the singular values of the Lanczos matrix of A associated with $\{r_0, \dots, r_k\}$ and $\{\tilde{r}_0, \dots, \tilde{r}_k\}$. The LU-decomposition in Bi-CG breaks down in step k if and only if the smallest singular value $\theta_-^{(k)}$ is zero. By scaling by $\theta_+^{(k)}$, we obtain $\theta_-^{(k)}/\theta_+^{(k)}$, that can be viewed as a quantification of the *near-breakdown of the LU-decomposition*. (For a relation to the Babuška-Brezzi condition, well-known in mixed finite element theory, cf. [2, section 4]). Apparently, (32) tells us that the growth of the Bi-CG residuals at step k can be bounded in terms of the “distance” λ_k to Lanczos breakdown and the distance $\theta_-^{(k)}/\theta_+^{(k)}$ to LU-decomposition breakdown.

If Bi-CG incorporated in the BiCGstab process does not suffer from near-breakdown of the Lanczos process nor from the LU-decomposition, then we expect

$$\|\psi_k(A)r_{k+1}\|/\|\psi_k(A)r_k\|$$

(cf. (31)) to be of moderate size.

4. Numerical experiments

In the previous sections, we have focussed on the scaled ρ_k (i.e. $\hat{\rho}_k^H$ as defined in (10)) and the scaled σ_k . Although these scaled values (where we have applied Cauchy-Schwartz) are actually smaller than the values that we display in the figures they do not differ significantly in our numerical examples. The figures also show that the values of $\hat{\rho}_k^H$ and $\hat{\sigma}_k^H$ are rather close to each other, as they should in view of our arguments in section 3.5.

All figures show

- the \log_{10} of the norm of the true residual $b - Ax_k$ (curve 1, full line in figure) for $k = m\ell$,
- the \log_{10} of $|(\mathbf{r}_k, \tilde{\mathbf{r}}_0)|/(|\mathbf{r}_k|, |\tilde{\mathbf{r}}_0|)$ (curve 2, dash-dotted line in figure), also for intermediate k ,
- the \log_{10} of $|(\mathbf{A}\mathbf{u}_k, \tilde{\mathbf{r}}_0)|/(|\mathbf{A}\mathbf{u}_k|, |\tilde{\mathbf{r}}_0|)$ (curve 3, dotted line in figure),
- and the $\log_{1/0.7}$ of $|\hat{\gamma}_k|$ (curve 4, the \bullet 's in figure), where $\hat{\gamma}_k$ is the scaled leading coefficient of the polynomial p_m that we actually used (cf. (28)) and 0.7 is the value for Ω if we modify the methods by (28) (i.e., if $\Omega \neq 0$).

We count the iteration phases by numbers of matrix-vector products (MV), since this makes it possible to compare effectively different BiCGstab variants.

The numerical results clearly indicate that limiting the relative size of the leading coefficient of the polynomial may help to cure the effects of stagnation of $\text{BiCGstab}(\ell)$, and also it may help to increase the value of ℓ (our standard choice was $\ell = 1$ in these examples). In some situations one modification may help, while in other situations the other modification may help or a combination of both.

4.1. Example 1

First we consider an advection dominated second order partial differential equation, with Dirichlet boundary conditions, on the unit cube (this equation was taken from [13]):

$$-u_{xx} - u_{yy} - u_{zz} + 1000u_x = f.$$

The function f is defined by the solution

$$u(x, y, z) = \exp(xyz) \sin(\pi x) \sin(\pi y) \sin(\pi z).$$

This equation was discretized using $10 \times 10 \times 10$ finite volumes and central differences for u_x , resulting in a seven-diagonal linear system of order 1000. No preconditioning has been used in this example (or other examples in order to make the differences between approaches more visible).

As we see, Bi-CGSTAB more or less stagnates (see figure 2), which creates the kind of situation that we were particularly interested in, and that we want to cure. After about 30 MVs, the scalars ρ_k and σ_k , and consequently the Bi-CG coefficients, do not have any correct digit. None of the scaled leading coefficients ϱ is extremely small (according to figure 2, $|\varrho| \geq (0.7)^{10} = 0.0282$; the \bullet 's in the figure show $\log_{1/0.7} |\varrho|$). Apparently, the accumulation of the amplifications of $|\hat{\mathbf{r}}|$ by the MR(1)-polynomials leads to this situation. The graph of $|(\mathbf{r}_k, \tilde{\mathbf{r}}_0)| / (|\mathbf{r}_k|, |\tilde{\mathbf{r}}_0|)$ shows nicely the predicted exponential decrease for $k \leq 10$ (that is for less than 20 MVs; $\prod_{k \leq 10} |\varrho| \geq (0.0282)^{10} = 3.18 \times 10^{-16}$). For larger k , $k \geq 11$, these values are smaller than the machine precision.

Although the modification with $\Omega = 0.7$ improves the relative size of the ρ_k and σ_k (only after $k = 50$ we lose all significant digits, cf. figure 3), we now even have divergence: the amplification of $\|\hat{\mathbf{r}}\|$ per step is not compensated by better convergence behavior of the incorporated Bi-CG process. Note that, for $k \leq 10$ (that is, $\#\text{MV} \leq 20$), $|(\mathbf{r}_k, \tilde{\mathbf{r}}_0)| / (|\mathbf{r}_k|, |\tilde{\mathbf{r}}_0|)$ decreases proportionally to $(1/\sqrt{1 + \Omega^{-2}})^k = 0.573^k$.

Using $\text{BiCGstab}(2)$ improves the situation: both versions of $\text{BiCGstab}(\ell)$ converge nicely. According to figure 4, $|\varrho| \geq (0.7)^3 = 0.343$. During the first 60 MVs, we have to deal with these ϱ 's 15 times. This leads us to expect a decrease of $|(\mathbf{r}_k, \tilde{\mathbf{r}}_0)| / (|\mathbf{r}_k|, |\tilde{\mathbf{r}}_0|)$ by at most $\geq (0.343)^{15} = 1.07 \times 10^{-7}$. This is confirmed by the numerical results. Here, $\text{BiCGstab}(2)$ seems to be able to retain locally five correct digits of the Bi-CG coefficients, which is, apparently, enough to keep the incorporated Bi-CG process converging as it should.

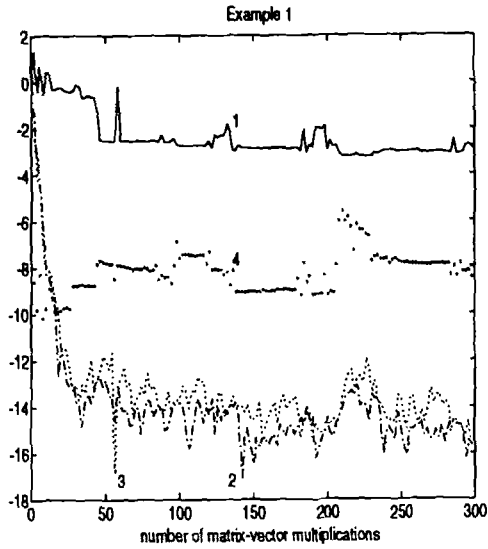


Figure 2. Standard Bi-CGSTAB.

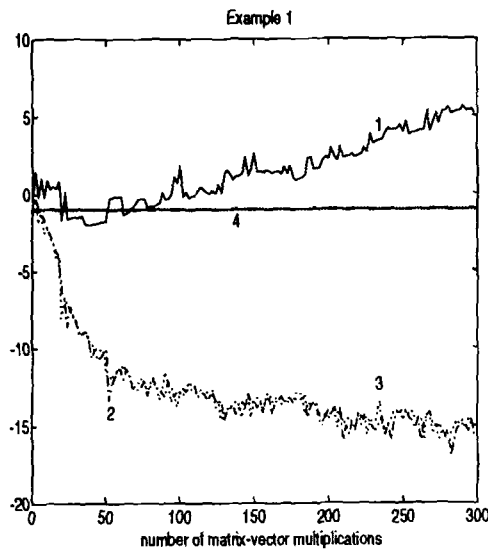


Figure 3. Modified $\Omega = 0.7$.

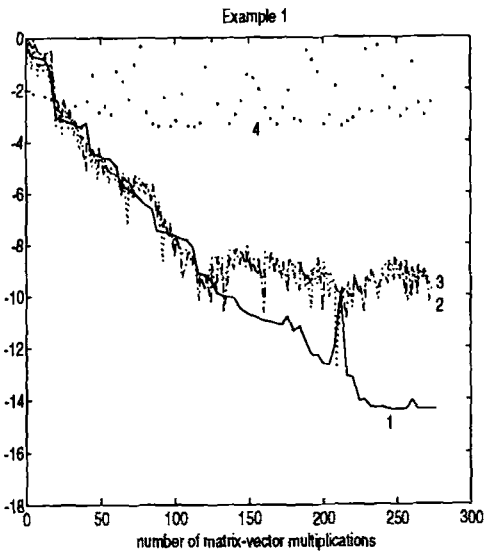


Figure 4. Standard BiCGstab(2).

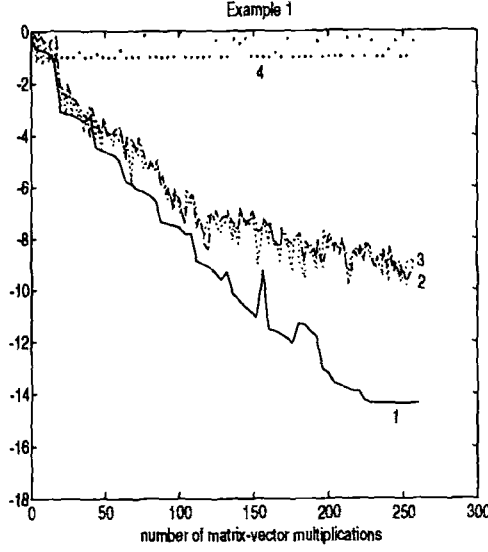


Figure 5. Modified $\Omega = 0.7$.

The modification with $\Omega = 0.7$ further improves the situation (but only slightly).

4.2. Example 2

For this comparison we have chosen the matrix that arises from a 63×63 finite volume discretization of

$$-u_{xx} - u_{yy} + \gamma(xu_x + yu_y) + \beta u = f$$

(34)

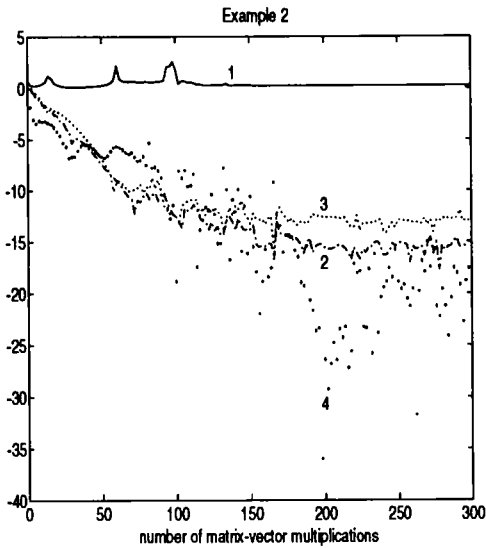


Figure 6. Standard Bi-CGSTAB.

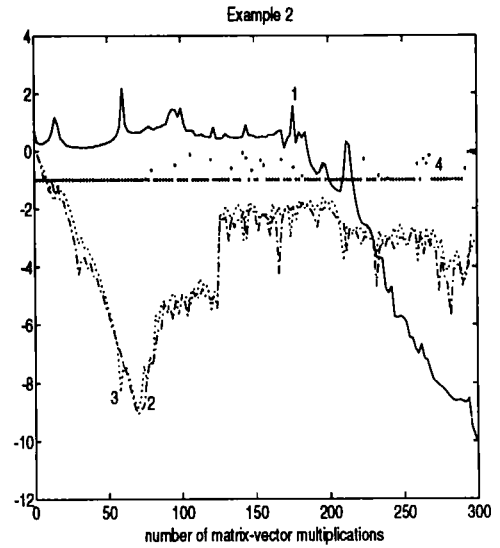
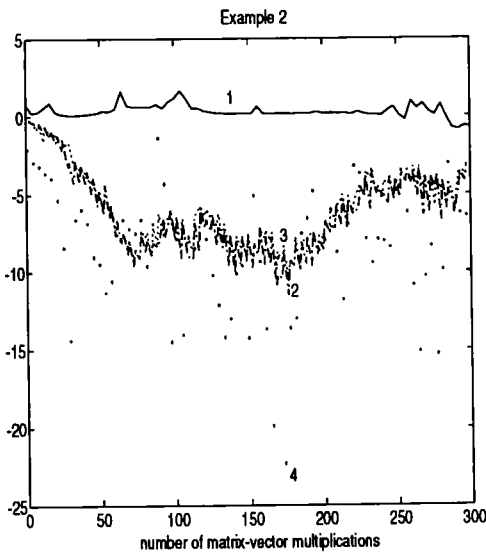
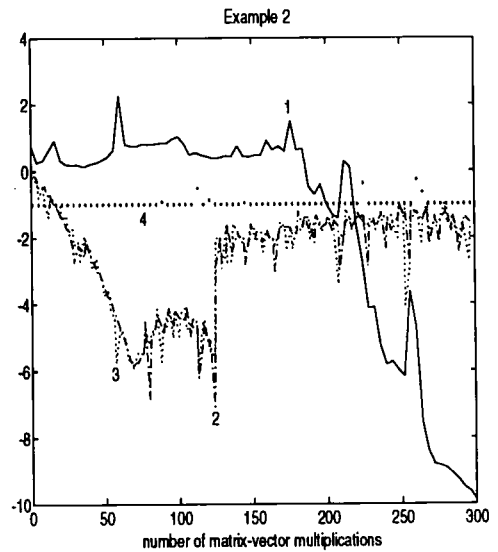
Figure 7. Modified $\Omega = 0.7$.

Figure 8. Standard BiCGstab(2).

Figure 9. Modified $\Omega = 0.7$.

on the unit cube with Dirichlet boundary conditions, with $\gamma = 100$ and $\beta = -200$. This example has been suggested in [8]. We have chosen the right hand side b such that the solution x of the equation $Ax = b$ is the vector $(1, 1, \dots, 1)^T$. The zero vector was used as an initial guess.

Here Bi-CGSTAB stagnates also (see figure 6). After about 150 MVs, the Bi-CG coefficients do not have any correct digit left, due to the cumulative effect of

non-small $|\varrho|$'s ($|\varrho| \approx (0.7)^5 = 0.168$). The modification with $\Omega = 0.7$ (cf. figure 7) improves the relative size of the ρ_k and σ_k significantly, enough to survive the phase where $|\varrho|$'s are ≈ 0.17 . As soon as the MR(1)-polynomial is more effective in reducing the norm of the residual (for $\#MV \geq 80$, we find $|\varrho| \geq 0.7$), the relative size of the ρ_k and σ_k grows.

With BiCGstab(2) the situation does not seem to improve (see figure 8). In the initial phase, the values of $|\varrho|$ for the MR(2)-polynomials are much smaller than the values of $|\varrho|$ for the MR(1)-polynomials. Although, up to $\#MV \approx 80$ the decrease of the scaled ρ_k and σ_k is a little bit better than with Bi-CGSTAB, this is apparently not enough to help survive the phase where the $|\varrho|$ is too small. However, (not shown in the figure) BiCGstab(2) did eventually converge (that is, $\|\mathbf{r}_k\| \leq 10^{-14} \|\mathbf{r}_0\|$) in 420 MV's, needing 120 MV more than the modified Bi-CGSTAB. 120 MVs is the length of the phase where the scaled ρ_k and σ_k are very small ($\approx 10^{-10}$).

The convergence behavior of modified BiCGstab(2) with $\Omega = 0.7$ (see figure 9) is comparable to the one of modified Bi-CGSTAB (in figure 7).

4.3. Example 3

For this example we have selected a problem similar to the one in example 2. Here, the matrix arises from a 66×66 finite volume discretization of (34), now with $\gamma = 1000$ and $\beta = 10$ as in [20].

The example shows that a combination of our strategies for Bi-CGSTAB (increasing ℓ and limiting the size of the scaled leading coefficient $\hat{\gamma}_k$, cf. (29)) can cure the accuracy problems, whereas each of the strategies independently may fail.

As figure 11 shows, limiting the leading coefficient improves the accuracy of the Bi-CG coefficients in the initial phase of the process but can not prevent a loss of all digits (for $\#MV > 200$). Now the amplifying effect on $\|\mathbf{r}_k\|$ of this choice of the polynomial can clearly be seen (from $\#MV = 200$ to $\#MV = 400$ the residual grows with $(\sqrt{1 + \Omega^2})^{100} = 4.6 \times 10^8$).

5. Conclusions

In order to maintain the convergence properties of the Bi-CG component in hybrid Bi-CG methods, it is necessary to select polynomial methods for the hybrid part that permit to compute the Bi-CG coefficients as accurately as possible.

The (dynamical) combination of two strategies for the improvement of the local accuracy seems to be very attractive. These strategies are:

- A. Take products of degree ℓ polynomials with $\ell > 1$ as in BiCGstab(ℓ) rather than of degree 1 polynomials as in Bi-CGSTAB.
- B. Try to limit the size of the leading coefficient of these polynomials, by switching occasionally between FOM and GMRES processes.

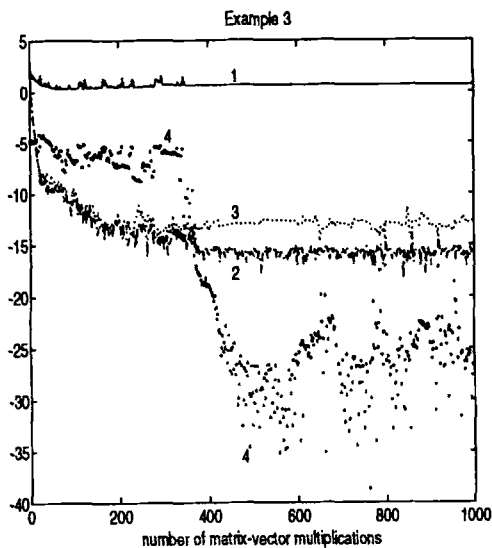


Figure 10. Standard Bi-CGSTAB.

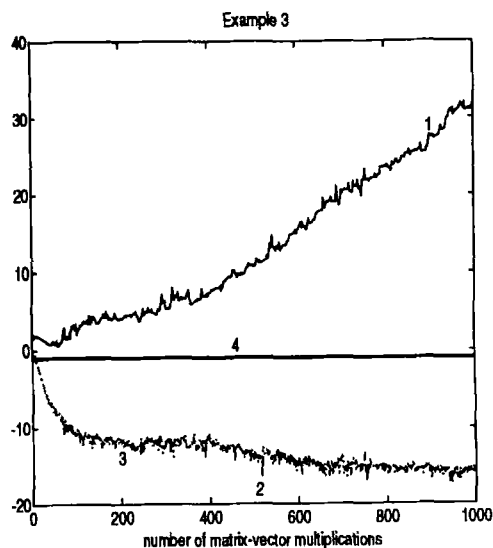
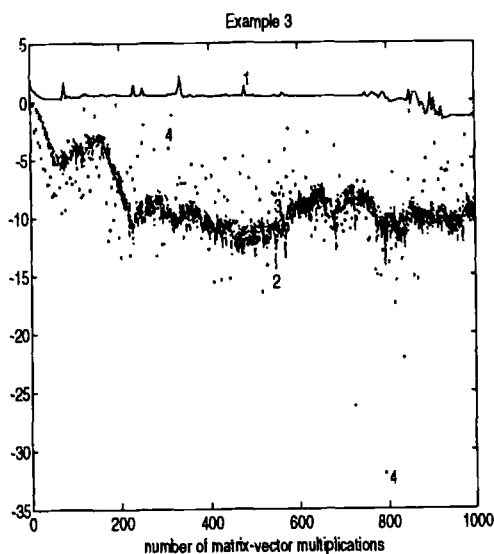
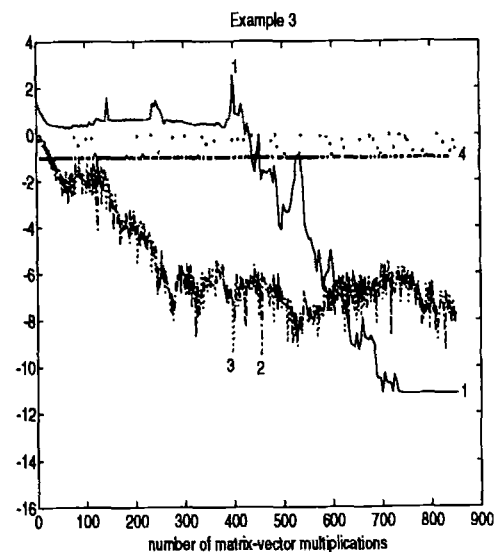
Figure 11. Modified $\Omega = 0.7$.

Figure 12. Standard BiCGstab(2).

Figure 13. Modified $\Omega = 0.7$.

This approach often leads to improved convergence and may help to overcome phases of stagnation. Our strategies are rather inexpensive, relative to the work per matrix-vector product.

Acknowledgement

We thank the referees for their suggestions to improve our presentation.

References

- [1] Z. Bai, Error analysis of the Lanczos algorithm for the nonsymmetric eigenvalue problem, *Math. Comp.* 62 (1994) 209–226.
- [2] R.E. Bank and T.F. Chan, An analysis of the composite step biconjugate gradient method, *Numer. Math.* 66 (1993) 295–319.
- [3] P.N. Brown, A theoretical comparison of the Arnoldi and GMRES algorithms, *SIAM J. Sci. Stat. Comp.* 12 (1991) 58–78.
- [4] S.C. Eisenstat, H.C. Elman and M.H. Schultz, Variational iterative methods for nonsymmetric systems of linear equations, *SIAM J. Numer. Anal.* 20 (1983) 345–357.
- [5] T.F. Chan, E. Gallopoulos, V. Simoncini, T. Szeto and C.H. Tong, A quasi-minimal residual variant of the Bi-CGSTAB algorithm for nonsymmetric systems, *SIAM J. Sci. Comp.* 15 (1994) 338–347.
- [6] G.H. Golub and C.F. Van Loan, *Matrix Computations*, 2nd ed. (The Johns Hopkins University Press, Baltimore and London, 1989).
- [7] R. Fletcher, Conjugate gradient methods for indefinite systems, in: *Proc. Dundee Biennial Conf. on Numerical Analysis*, ed. G. Watson (Springer, New York, 1975).
- [8] R.W. Freund, A transpose-free quasi-minimal residual algorithm for non-Hermitian linear systems, *SIAM J. Sci. Comp.* 14 (1993) 470–482.
- [9] R.W. Freund, M.H. Gutknecht and N. Nachtigal, An implementation of the look-ahead Lanczos algorithm for non-Hermitian matrices, *SIAM J. Sci. Comp.* 14 (1993) 137–158.
- [10] A. Greenbaum, Behavior of slightly perturbed Lanczos and conjugate gradient recurrences, *Lin. Alg. Appl.* 113 (1989) 7–63.
- [11] M.H. Gutknecht, Variants of BiCGStab for matrices with complex spectrum, *SIAM J. Sci. Comp.* 14 (1993) 1020–1033.
- [12] C. Lanczos, Solution of systems of linear equations by minimized iteration, *J. Res. Nat. Bur. Stand.* 49 (1952) 33–53.
- [13] U. Meier Yang, Preconditioned conjugate gradient-like methods for nonsymmetric linear systems, Preprint, Center for Research and Development, University of Illinois at Urbana-Champaign (1992).
- [14] J. Modersitzki and G.L.G. Sleijpen, An error analysis of CSBCG, in preparation.
- [15] A. Neumaier, Oral presentation at the Oberwolfach meeting “Numerical Linear Algebra”, Oberwolfach (April 1994).
- [16] C. Paige, Accuracy and effectiveness of the Lanczos algorithm for the symmetric eigenproblem, *Lin. Alg. Appl.* 34 (1980) 235–258.
- [17] B.N. Parlett, D.R. Taylor and Z.A. Liu, A look-ahead Lanczos algorithm for unsymmetric matrices, *Math. Comp.* 44 (1985) 105–124.
- [18] Y. Saad, Krylov subspace methods for solving large unsymmetric linear systems, *Math. Comp.* 37 (1981) 105–126.
- [19] Y. Saad and M.H. Schultz, GMRES: A generalized minimum residual algorithm for solving nonsymmetric linear systems, *SIAM J. Sci. Stat. Comp.* 7 (1986) 856–869.
- [20] Y. Saad, A flexible inner-outer preconditioned GMRES algorithm, *SIAM J. Sci. Comp.* 14 (1993) 461–469.
- [21] G.L.G. Sleijpen and D.R. Fokkema, BiCGstab(ℓ) for linear equations involving matrices with complex spectrum, *ETNA* 1 (1993) 11–32.
- [22] G.L.G. Sleijpen and H.A. van der Vorst, Reliable updated residuals in hybrid Bi-CG methods, Preprint 886, Dept. Math., University Utrecht (1994), to appear in *Computing*.
- [23] G.L.G. Sleijpen, H.A. van der Vorst and D.R. Fokkema, BiCGstab(ℓ) and other hybrid Bi-CG methods, *Numer. Algor.* 7 (1994) 75–109.
- [24] P. Sonneveld, CGS, a fast Lanczos-type solver for nonsymmetric linear systems, *SIAM J. Sci. Stat. Comp.* 10 (1989) 36–52.

- [25] H.A. van der Vorst, Bi-CGSTAB: A fast and smoothly converging variant of Bi-CG for the solution of nonsymmetric linear systems, *SIAM J. Sci. Stat. Comp.* 13 (1992) 631–644.
- [26] H.A. van der Vorst and C. Vuik, The superlinear convergence behaviour of GMRES, *J. Comp. Appl. Math.* 48 (1993) 327–341.
- [27] D.M. Young and K.C. Jea, Generalized conjugate-gradient acceleration of nonsymmetrizable iterative methods, *Lin. Alg. Appl.* 34 (1980) 159–194.