

# How to make Simpler GMRES and GCR more stable

Pavel Jiránek<sup>†</sup>, Miroslav Rozložník<sup>‡</sup> and Martin H. Gutknecht

Research Report No. 2008-10  
May 2008

Seminar für Angewandte Mathematik  
Eidgenössische Technische Hochschule  
CH-8092 Zürich  
Switzerland

---

<sup>†</sup>Faculty of Mechatronics and Interdisciplinary Engineering Studies, Technical University of Liberec, Hájkova 6, CZ-461 17 Liberec, Czech Republic ([pavel.jiranek@tul.cz](mailto:pavel.jiranek@tul.cz)). The work of this author was supported by the MSMT CR under the project 1M0554 “Advanced Remedial Technologies”.

<sup>‡</sup>Institute of Computer Science, Academy of Sciences of the Czech Republic, Pod vodárenskou věží 2, CZ-182 07 Prague 8, Czech Republic ([miro@cs.cas.cz](mailto:miro@cs.cas.cz)). The work of this author was supported by the project 1ET400300415 within the National Program of Research “Information Society” and by the Institutional Research Plan AV0Z10300504 “Computer Science for the Information Society: Models, Algorithms, Applications”.

# How to make Simpler GMRES and GCR more stable

Pavel Jiránek<sup>†</sup>, Miroslav Rozložník<sup>‡</sup> and Martin H. Gutknecht

Seminar für Angewandte Mathematik  
Eidgenössische Technische Hochschule  
CH-8092 Zürich  
Switzerland

Research Report No. 2008-10

May 2008

## Abstract

In this paper we analyze the numerical behavior of several minimum residual methods, which are mathematically equivalent to the GMRES method. Two main approaches are compared: the one that computes the approximate solution (similar to GMRES) in terms of a Krylov space basis from an upper triangular linear system for the coordinates, and the one where the approximate solutions are updated with a simple recursion formula. We show that a different choice of the basis can significantly influence the numerical behavior of the resulting implementation. While Simpler GMRES and ORTHODIR are less stable due to the ill-conditioning of the basis used, the residual basis is well-conditioned as long as we have a reasonable residual norm decrease. These results lead to a new implementation, which is conditionally backward stable, and, in a sense they explain the experimentally observed fact that the GCR (ORTHOMIN) method delivers very accurate approximate solutions when it converges fast enough without stagnation.

**Keywords:** Large-scale nonsymmetric linear systems, Krylov subspace methods, minimum residual methods, numerical stability, rounding errors.

---

<sup>†</sup>Faculty of Mechatronics and Interdisciplinary Engineering Studies, Technical University of Liberec, Hálkova 6, CZ-461 17 Liberec, Czech Republic ([pavel.jiranek@tul.cz](mailto:pavel.jiranek@tul.cz)). The work of this author was supported by the MSMT CR under the project 1M0554 “Advanced Remedial Technologies”.

<sup>‡</sup>Institute of Computer Science, Academy of Sciences of the Czech Republic, Pod vodárenskou věží 2, CZ-182 07 Prague 8, Czech Republic ([miro@cs.cas.cz](mailto:miro@cs.cas.cz)). The work of this author was supported by the project 1ET400300415 within the National Program of Research “Information Society” and by the Institutional Research Plan AV0Z10300504 “Computer Science for the Information Society: Models, Algorithms, Applications”.

# 1 Introduction

In this paper we consider certain methods for solving a system of linear algebraic equations

$$Ax = b, \quad A \in \mathbb{R}^{N \times N}, \quad b \in \mathbb{R}^N, \quad (1)$$

where  $A$  is a large and sparse nonsingular matrix that is, in general, non-symmetric. For solving such systems, Krylov subspace methods are very popular. They build a sequence of iterates  $x_n$  ( $n = 0, 1, 2, \dots$ ) such that  $x_n \in x_0 + \mathcal{K}_n(A, r_0)$ , where  $\mathcal{K}_n(A, r_0) \equiv \text{span}\{r_0, Ar_0, \dots, A^{n-1}r_0\}$  is the  $n$ th Krylov subspace generated by the matrix  $A$  from the residual  $r_0 \equiv b - Ax_0$  that corresponds to the initial guess  $x_0$ . Many approaches for defining such approximations  $x_n$  have been proposed, see, e.g., the books by Greenbaum [11], Meurant [19], and Saad [23]. In particular, due to their smooth convergence behavior, minimum residual methods satisfying

$$\|r_n\| = \min_{\tilde{x} \in x_0 + \mathcal{K}_n(A, r_0)} \|b - A\tilde{x}\|, \quad r_n \equiv b - Ax_n, \quad (2)$$

are widely used, e.g., the GMRES algorithm of Saad and Schultz [24].

The classical implementation of GMRES makes use of a nested sequence of orthonormal bases of the Krylov subspaces  $\mathcal{K}_n(A, r_0)$ . These bases are generated by an Arnoldi process [2]. With the notation  $\rho_0 \equiv \|r_0\|$ ,  $q_1 \equiv \rho_0^{-1}r_0$ ,  $Q_n \equiv [q_1, \dots, q_n]$ , where the columns of  $Q_n$  form this orthonormal basis of  $\mathcal{K}_n(A, r_0)$ , and with an  $(n+1) \times n$  upper Hessenberg matrix  $H_{n+1,n}$ , its result can be cast in matrix form as

$$[q_1, AQ_n] = Q_{n+1}[e_1, H_{n+1,n}].$$

This can be viewed as the QR factorization of the matrix  $[q_1, AQ_n]$ . Ultimately, an approximate solution  $x_n$  satisfying the minimum residual property (2) is constructed in the form  $x_n = x_0 + Q_n y_n$ , but  $x_n$  is not needed at every step. From the relation

$$\|r_n\| = \|r_0 - AQ_n y_n\| = \|\rho_0 e_1 - H_{n+1,n} y_n\|$$

it follows that  $y_n$  is the solution of the  $(n+1) \times n$  least squares problem  $H_{n+1,n} y_n \approx \rho_0 e_1$ , and that  $\|r_n\|$  equals the norm of its residual  $\rho_0 e_1 - H_{n+1,n} y_n \in \mathbb{R}^{n+1}$ . This problem can be solved via the recursive QR factorization of  $H_{n+1,n}$ , updated by applying  $n$  Givens rotations and determining

a new one in the  $n$ th step. Once the norm of the residual is small enough — which can be seen without explicitly solving the least squares problem — the triangular system with the computed R-factor is solved, and the approximate solution  $x_n$  is computed. In [6, 12, 20] it was shown that this “classical” version of the GMRES method is backward stable provided that the Arnoldi process is implemented using the modified Gram-Schmidt algorithm or Householder reflections.

In this paper we deal with a different approach proposed by Walker and Zhou [29], who called it the Simpler GMRES method. To derive it, we recall that the minimum residual property (2) is equivalent to the orthogonality condition

$$r_n \perp AK_n(A, r_0),$$

where  $\perp$  is the orthogonality relation induced by the standard Euclidean inner product  $\langle \cdot, \cdot \rangle$ . Instead of building an orthonormal basis of  $\mathcal{K}_n(A, r_0)$  we look for an orthonormal basis  $V_n \equiv [v_1, \dots, v_n]$  of  $AK_n(A, r_0)$ . As proposed by Walker and Zhou, we could construct it again by an Arnoldi process. This leads to the QR factorization

$$A[q_1, V_{n-1}] = V_n U_n, \quad (3)$$

where  $U_n$  is an  $n \times n$  upper triangular matrix. We propose a generalization that consists in allowing to replace this Arnoldi process. Instead of using the image  $Av_{n-1}$  of the last constructed orthonormal basis vectors to extend the basis we consider any nested sequence of matrices  $Z_{n-1} \equiv [z_1, \dots, z_{n-1}]$  such that the columns of  $[q_1, Z_{n-1}]$  form a basis of  $\mathcal{K}_n(A, r_0)$ , and we make use of  $Az_{n-1}$  to extend the basis. We may assume that the columns  $z_k$  of  $Z_{n-1}$  have unit length (and we will do so in the error analysis), but they need not be mutually orthogonal. The orthonormal basis  $V_n$  of  $AK_n(A, r_0)$  is thus obtained from the QR factorization of the image of  $[q_1, Z_{n-1}]$ :

$$A[q_1, Z_{n-1}] = V_n U_n. \quad (4)$$

Since  $r_n \in r_0 + AK_n(A, r_0) = r_0 + \mathcal{R}(V_n)$  and  $r_n \perp \mathcal{R}(V_n)$ , we can obtain the residual from  $r_n = (I - V_n V_n^T)r_0$ . Note that  $r_n$  is just the orthogonal projection of  $r_0$  onto the orthogonal complement of  $\mathcal{R}(V_n)$ . To compute it we apply the modified Gram-Schmidt method, which leads to the recursion

$$r_n = r_{n-1} - \alpha_n v_n, \quad \alpha_n \equiv \langle r_{n-1}, v_n \rangle. \quad (5)$$

This recursion can be cast into a matrix relation too. Let  $R_{n+1} \equiv [r_0, \dots, r_n]$ , let  $D_n \equiv \text{diag}(\alpha_1, \dots, \alpha_n)$ , and let  $L_{n+1,n} \in \mathbb{R}^{(n+1) \times n}$  be the bidiagonal matrix with ones on the main diagonal and minus ones on the first subdiagonal; then (5) can be written as

$$R_{n+1}L_{n+1,n} = V_n D_n. \quad (6)$$

Since the columns of  $[q_1, Z_{n-1}]$  are a basis of  $\mathcal{K}_n(A, r_0)$ , we can represent  $x_n$  in the form

$$x_n = x_0 + [q_1, Z_{n-1}]t_n, \quad (7)$$

so that  $r_n = r_0 - A[q_1, Z_{n-1}]t_n = r_0 - V_n U_n t_n$ . Due to the minimum residual property, we have  $r_n \perp \mathcal{R}(V_n)$ , and thus simply

$$U_n t_n = V_n^T r_0 = [\alpha_1, \dots, \alpha_n]^T. \quad (8)$$

Hence, once the residual norm is small enough, we can solve this triangular system and compute  $x_n = x_0 + [q_1, Z_{n-1}]t_n$ . We call this general approach the *simpler approach*. It includes, as a special case, Simpler GMRES, where  $Z_{n-1} \equiv V_{n-1}$ . We will also be interested in the case of the residual basis  $[q_1, Z_{n-1}] = [\frac{r_0}{\|r_0\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$ , which we will call SGMRES/RB, where ‘‘RB’’ refers to ‘‘residual basis’’. (Recently this method has also been derived and implemented by Yvan Notay.)

Recursion (5) reveals the connection between the simpler approach and yet another minimum residual approach. Let us set  $p_n \equiv A^{-1}v_n$ ,  $P_n \equiv [p_1, \dots, p_n]$ . Then, left-multiplying (5) by  $A^{-1}$  yields

$$x_n = x_{n-1} + \alpha_n p_n, \quad \alpha_n = \langle r_{n-1}, A p_n \rangle, \quad (9)$$

or, in matrix form,

$$X_{n+1}L_{n+1,n} = -P_n D_n$$

with  $X_{n+1} \equiv [x_0, \dots, x_n]$ . This shows that  $p_n \in \mathcal{K}_n(A, r_0)$  is a direction vector: it has the direction in which one moves from  $x_{n-1}$  to  $x_n$ . The step length  $\alpha_n$  can be determined from one of the formulas on the right-hand side of (5) or (9). Recall that it follows from the condition  $\langle r_{n-1}, v_n \rangle = 0$ , which enforces the minimization of  $\|r_n\|$  on the line  $\alpha \mapsto r_{n-1} - \alpha v_n$ . So, instead of computing the coordinates  $t_n$  of  $x_n - x_0$  with respect to the columns of  $[q_1, Z_{n-1}]$  first, we can directly update  $x_n$  from (9). However, this requires that we construct the direction vector  $p_n$  (or a scalar multiple of it). Now, note that left-multiplying (4) by  $A^{-1}$  yields

$$[q_1, Z_{n-1}] = P_n U_n. \quad (10)$$

If  $U_n$  is known from (4), a recursion for  $p_n$  can be extracted from this formula. Note that it has the same recurrence coefficients (stored in the columns of  $U_n$ ) that are used in the Gram-Schmidt process in (4); so the two recursions can be run in the same loop. The obvious disadvantages of this approach is that we have to store both all the direction vectors  $p_n$  and all the original orthonormal basis vectors  $v_n = Ap_n$ . Moreover, any roundoff errors in  $U_n$  may have a strong effect on  $P_n$ . However, as we will see, this is the price we have to pay if we want to apply the simple and convenient 2-term update formulas (5) and (9) and spend only one matrix-vector (MV) product per step, namely  $Az_{n-1}$  in (4) (or  $Av_{n-1}$  in (3) if  $Z_{n-1} \equiv V_{n-1}$ ). The case  $Z_{n-1} \equiv V_{n-1}$  of this method was proposed in [22] under the name  $A^T A$ -variant of GMRES. We will use here the terminology *update approach* for this case and, more exactly, refined ORTHODIR for the particular case with  $Z_{n-1} \equiv V_{n-1}$ , since, as we will see, it is a refined version of the residual norm minimizing ORTHODIR algorithm [10, 31]. Likewise the case with  $Z_{n-1} = [\frac{r_1}{\|r_1\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$ , which can be viewed as a refined version of the ORTHOMIN algorithm [28, 31] (or the GCR method of Elman [9, 8]) and is identical to the GMRESR method [27] of van der Vorst and Vuik with the choice  $u_n^{(0)} = r_n$ , will be referred to as refined ORTHOMIN (see our comments below).

The refined ORTHODIR and ORTHOMIN algorithms with residual norm minimization started from the fact that the direction vectors  $p_n$  of the minimum residual method characterized by (2) are  $A^T A$ -orthonormal to each other: since  $V_n = AP_n$ , we have  $P_n^T A^T AP_n = V_n^T V_n = I$ . Because directions are only determined up to a scalar multiple, we might give up the normalization of  $V_n$  and choose instead  $P_n^T A^T AP_n = V_n^T V_n$  to be a nonsingular diagonal matrix. So, in analogy to (4), we can directly compute the columns of  $P_n = [p_1, \dots, p_n]$  and  $U_n$  from (10), and complement this by the explicit successive evaluation of  $V_n = AP_n$  (which, at the same time, serves for extending the Krylov subspace). Again, we can view (10) as either an Arnoldi process for an  $A^T A$ -orthogonal basis if we choose  $Z_{n-1} \equiv AP_{n-1}$ , or as a Gram-Schmidt implementation of a QR decomposition of  $[q_1, Z_{n-1}]$  with respect to the  $A^T A$ -inner product if  $Z_{n-1}$  originates elsewhere. The case where  $Z_{n-1} \equiv AP_{n-1}$ ,  $q_1 \equiv r_0$ , and  $U_n$  is unit triangular corresponds to the original ORTHODIR algorithm [10, 31]; the case where  $Z_{n-1} \equiv [r_1, \dots, r_{n-1}]$ ,  $q_1 \equiv r_0$ , and  $U_n$  is unit triangular yields a version of the ORTHOMIN algorithm as proposed by Young and Jea [31], which was called GCR by Elman [9]. Despite the popularity of the name GCR we will mostly use the older

name ORTHOMIN here, which also underlines the analogy to ORTHODIR. Details can also be found in [3] (choosing  $B = A^T A$  and  $C = I$  there). The cases with short-term recurrences have been treated in detail in [17] and [4].

However, what we have concealed in these descriptions is that we need a second matrix-vector product, namely  $Av_{n-1}$  in ORTHODIR and  $Ar_n$  in ORTHOMIN, to compute the coefficients of the orthogonal projection (i.e., of the Gram-Schmidt algorithm). Due to the  $A^T A$ -orthogonality, in ORTHODIR the relevant projection of  $Ap_{n-1}$  is  $p_n = (I - P_{n-1}(AP_{n-1})^T A)Ap_{n-1}$ , which with  $V_{n-1} = AP_{n-1}$  may be written as  $p_n = (I - P_{n-1}V_{n-1}^T A)v_{n-1}$ . The new vector  $v_n$  can be computed either from  $v_n = (I - V_{n-1}V_{n-1}^T)Av_{n-1}$  or directly as  $v_n = Ap_n$ , but the latter requires an extra MV. An analogue consideration holds for ORTHOMIN. So, in the latter form, these algorithms are not competitive. Some remarks on their stability were drawn in [11]; we will not cover these implementations here.

The well-known remedy suggested by Vinsome [28] and Eisenstadt, Elman, and Schultz [8] consists in computing and storing both  $P_n$  and  $V_n$ . This is achieved by computing  $V_n$  with either the Arnoldi process (3) or with another QR decomposition of  $A[r_0, r_1, \dots, r_{n-1}]$  analogous to (4). But this means that up to the scaling of the bases  $P_n$ ,  $V_n$ , and  $Z_n$  we return to the refined ORTHODIR and refined ORTHOMIN algorithms discussed above. The remaining difference between Vinsome's ORTHOMIN and our refined ORTHOMIN is that we normalize the residuals before orthogonalizing them, and that we use normalized direction vectors. The analog is true for the difference between the usual implementation of ORTHODIR and our refined ORTHODIR. The importance of normalizing the residuals before the orthogonalization will be seen later.

The paper is organized as follows. In Section 2 we analyze first the maximum attainable accuracy of the simpler approach based on (3) or (4) for  $v_n$  and (7), (8) for  $x_n$ . Then we turn to the update approach based on (3) or (4) for  $v_n$ , (10) for  $p_n$ , and (9), (5) for  $x_n$  and  $r_n$ . To keep the text readable, we assume rounding errors only in selected, most relevant parts of the computation. The bounds presented in Theorems 2.1 and 2.2 show that the conditioning of the matrix  $[q_1, Z_{n-1}]$  plays an important role in the numerical stability of these schemes. Both theorems give bounds on the maximum attainable accuracy measured by the normwise backward error. While for the simpler approach this quantity does not depend on the conditioning of  $A$ , the bound for the update approach is proportional to  $\kappa(A)$  (as we will show in our constructed numerical example, the bound is

attainable). However, the dependence on  $\kappa(A)$  is usually an overestimate; in practice, both the simpler and update approaches behave almost equally for the same choice of the basis. This is especially true for the relative errors of the computed approximate solutions, where we give essentially the same upper bound. The situation is completely analogous to results for the GMRES method [24] and the MINRES method [21] given by Sleijpen, van der Vorst and Modersitzki in [26].

In Section 3 we derive particular results for two choices of the basis  $[q_1, Z_{n-1}]$ . First for  $[q_1, Z_{n-1}] = [q_1, V_{n-1}]$  leading to Simpler GMRES by Walker and Zhou [29] and to refined ORTHODIR. Then for  $[q_1, Z_{n-1}] = [\frac{r_0}{\|r_0\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$ , which leads to SGMRES/RB and refined ORTHOMIN, respectively. It appears that the two choices lead to truly different behavior in the condition number of  $U_n$ , which governs the stability of the considered schemes. Since all these methods converge in a finite number of iterations, we fix the iteration index  $n$  such that  $r_0 \notin AK_{n-1}(A, r_0)$ , that is, the exact solution has not yet been reached. Based on this we give conditions on the linear independence of the basis  $[q_1, Z_{n-1}]$ . It is known that  $[r_0, \dots, r_{n-1}]$  can be rank deficient when the GMRES method stagnates (the breakdown occurs in ORTHOMIN and hence also in SGMRES/RB), while this does not happen for  $[q_1, V_{n-1}]$  (Simpler GMRES and ORTHODIR are breakdown-free). On the other hand, we show that while the choice  $[q_1, Z_{n-1}] = [q_1, V_{n-1}]$  leads to inherently less numerically stable schemes, the second selection  $[q_1, Z_{n-1}] = [\frac{r_0}{\|r_0\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$  gives rise to conditionally stable implementations provided we have some reasonable residual decrease. In particular, we show that the SGMRES/RB implementation is conditionally backward stable. Our theoretical results are illustrated by selected numerical experiments. In Section 4 we conclude and give directions for future work.

Throughout the paper, we denote by  $\|\cdot\|$  the Euclidean vector norm and the induced matrix norm, and by  $\|\cdot\|_F$  the Frobenius norm. Moreover, for  $B \in \mathbb{R}^{N \times n}$  ( $N \geq n$ ) of rank  $n$ ,  $\sigma_1(B) \geq \sigma_n(B) > 0$  are the extremal singular values of  $B$ , and  $\kappa(B) = \sigma_1(B)/\sigma_n(B)$  is the spectral condition number. By  $I$  we denote the unit matrix of a suitable dimension, by  $e_k$  ( $k = 1, 2, \dots$ ) its  $k$ th column, and we let  $e \equiv [1, \dots, 1]^T$ . We assume the standard model of finite precision arithmetic with the unit roundoff  $u$  (see Higham [15] for details). In our bounds, instead of distinguishing between several constants (which are in fact polynomials in  $N$  and  $n$  that can differ from place to place), we use a generic constant  $c$ .



## 2 Maximum attainable accuracy of simpler and update approaches

In this section we analyze the numerical stability of the simpler and update approaches formulated in the previous section. In order to make our analysis readable, we assume that only the computations performed in (4), (8) and (10) are affected by rounding errors and that the computed Q-factor in the QR factorization (4) is close to an orthonormal matrix and has been computed in a backward stable way. Hence we assume that the computed (orthogonal) factor  $V_n$  and the upper triangular factor  $U_n$  in the QR factorization (4) satisfy

$$A[q_1, Z_{n-1}] = V_n U_n + F_n, \quad \|F_n\| \leq cu \|A\| \| [q_1, Z_{n-1}] \|, \quad (11)$$

and  $\|V_n - \hat{V}_n\| \leq cu$ , where  $\hat{V}_n$  is the nearest orthonormal matrix satisfying  $\hat{V}_n^T \hat{V}_n = I$ . For simplicity, we will not distinguish between  $V_n$  and  $\hat{V}_n$  and assume that  $V_n$  is exactly orthonormal. For details we refer to [5, 15]. From [30, 15] we have for the computed solution  $\hat{t}_n$  of (8) that

$$(U_n + \Delta U_n) \hat{t}_n = D_n e, \quad |\Delta U_n| \leq cu |U_n|, \quad (12)$$

where the absolute value and inequalities are understood component-wise. The approximation  $\hat{x}_n$  to  $x$  is then computed as

$$\hat{x}_n = x_0 + [q_1, Z_{n-1}] \hat{t}_n. \quad (13)$$

The crucial quantity for the analysis of the maximum attainable accuracy is the gap between the true residual  $b - A\hat{x}_n$  of the computed approximation and the updated residual  $r_n$  obtained from the update formula (5) describing the projection of the previous residual; see [11, 14]. In fact, once the true residual becomes negligible compared to the true one (and in the algorithms considered here it ultimately will), the gap equals the true residual divided by  $\|A\| \|\hat{x}_n\|$ , which therefore can be thought of as the backward error of the ultimate approximate solution  $\hat{x}_n$  (after suitable normalization). Here is our basic result on this gap for the simpler approach.

**Theorem 2.1.** *In the simpler approach, the gap between the true residual  $b - A\hat{x}_n$  and the updated residual  $r_n$  satisfies*

$$\frac{\|b - A\hat{x}_n - r_n\|}{\|A\| \|\hat{x}_n\|} \leq cu \kappa([q_1, Z_{n-1}]) \left( 1 + \frac{\|x_0\|}{\|\hat{x}_n\|} \right).$$

*Proof.* From (13) we have  $b - A\hat{x}_n = r_0 - A[q_1, Z_{n-1}]\hat{t}_n = r_0 - (V_n U_n + F_n)(U_n + \Delta U_n)^{-1} D_n e$ , and (5) gives  $r_n = r_0 - V_n D_n e$ . Using the identity  $I - U_n(U_n + \Delta U_n)^{-1} = \Delta U_n(U_n + \Delta U_n)^{-1}$  and the relation  $[q_1, Z_{n-1}](U_n + \Delta U_n)^{-1} D_n e = [q_1, Z_{n-1}]\hat{t}_n = \hat{x}_n - x_0$  we can express the gap between  $b - A\hat{x}_n$  and  $r_n$  as

$$\begin{aligned} b - A\hat{x}_n - r_n &= (V_n - (V_n U_n + F_n)(U_n + \Delta U_n)^{-1}) D_n e \\ &= (V_n \Delta U_n + F_n)(U_n + \Delta U_n)^{-1} D_n e \\ &= (V_n \Delta U_n + F_n)[q_1, Z_{n-1}]^\dagger [q_1, Z_{n-1}](U_n + \Delta U_n)^{-1} D_n e \\ &= (V_n \Delta U_n + F_n)[q_1, Z_{n-1}]^\dagger (\hat{x}_n - x_0). \end{aligned} \quad (14)$$

Taking the norm, considering (11), and noting that the terms involving  $V_n \Delta U_n$  and  $F_n$  can be subsumed into the generic constant  $c$ , we get

$$\|b - A\hat{x}_n - r_n\| \leq cu \|A\| \| [q_1, Z_{n-1}] \| \| [q_1, Z_{n-1}]^\dagger \| (\|\hat{x}_n\| + \|x_0\|). \quad (15)$$

Division by  $\|A\| \|\hat{x}_n\|$  concludes the proof.  $\square$

In the following we analyze the maximum attainable accuracy of the update approach. In accordance with (11) we assume that in finite precision arithmetic the computed direction vectors satisfy

$$[q_1, Z_{n-1}] = P_n U_n + G_n, \quad \|G_n\| \leq cu \|P_n\| \|U_n\|. \quad (16)$$

Note that the norm of the matrix  $G_n$  cannot be bounded by  $cu \|A\| \| [q_1, Z_{n-1}] \|$  as it is in the case of the QR factorization (11). As in (9) we compute then the approximate solution  $\hat{x}_n$  as

$$\hat{x}_n = \hat{x}_{n-1} + \alpha_n p_n. \quad (17)$$

**Theorem 2.2.** *In the update approach, the gap between the true residual  $b - A\hat{x}_n$  and the updated residual  $r_n$  satisfies*

$$\frac{\|b - A\hat{x}_n - r_n\|}{\|A\| \|\hat{x}_n\|} \leq cu \kappa(A) \kappa([q_1, Z_{n-1}]) \left( 1 + \frac{\|x_0\|}{\|\hat{x}_n\|} \right),$$

*provided that  $\eta_n \equiv 1 - cu \kappa(A) \kappa([q_1, Z_{n-1}]) > 0$ .*

*Proof.* Since  $\hat{x}_n = x_0 + P_n D_n e = x_0 + ([q_1, Z_{n-1}] - G_n) U_n^{-1} D_n e$  and  $r_n = r_0 - V_n D_n e$ , we have that

$$\begin{aligned} b - A\hat{x}_n - r_n &= (V_n - A[q_1, Z_{n-1}]U_n^{-1})D_n e + AG_n U_n^{-1} D_n e \\ &= (-F_n + AG_n)U_n^{-1} D_n e \end{aligned} \quad (18)$$

due to (4). From (4) and (16), we get  $P_n = A^{-1}V_n + (A^{-1}F_n - G_n)U_n^{-1}$ . Taking a norm we obtain  $\|P_n\| \leq \|A^{-1}\| + cu\kappa(A)\|U_n^{-1}\| + cu\|P_n\|\kappa(U_n)$ . The norm of the residual matrix  $G_n$  in (16) can hence be estimated as

$$\|G_n\| \leq cu\kappa(A)\|[q_1, Z_{n-1}]\|. \quad (19)$$

Owing to (17), we have the identity  $U_n^{-1} D_n e = U_n^{-1} P_n^\dagger P_n D_n e = U_n^{-1} P_n^\dagger (\hat{x}_n - x_0)$ , and  $\|U_n^{-1} P_n^\dagger\| \leq \eta_n^{-1} \|[q_1, Z_{n-1}]^\dagger\|$  following from (16). Thus we obtain

$$\|U_n^{-1} D_n e\| \leq \eta_n^{-1} \|[q_1, Z_{n-1}]^\dagger\| (\|\hat{x}_n\| + \|x_0\|), \quad (20)$$

which together with (18), (19), and (11) proves the statement of the theorem.  $\square$

The bound on the ultimate backward error given in Theorem 2.2 is worse than the one of Theorem 2.1. We see that for the simpler approach the normwise backward error is on the order of the roundoff unit, whereas for the update approach we have an upper bound proportional to the condition number of  $A$ . In terms of the residual norms, this leads to the bounds involving  $cu\kappa(A)\kappa([q_1, Z_{n-1}])$  and  $cu\kappa^2(A)\kappa([q_1, Z_{n-1}])$  terms for the simpler and update approach, respectively.

From Theorems 2.1 and 2.2, we can also estimate the ultimate level of the relative 2-norm of the error of both the simpler and update approach. However, as shown below, it appears that *the update approach leads to an approximate solution on essentially the same accuracy level in the error as the simpler approach*. A similar phenomenon was also observed by Sleijpen, van der Vorst and Modersitzki [26] in the symmetric case for GMRES and MINRES.

**Corollary 2.1.** *The gap between the computed approximate solutions  $\hat{x}_n$  and exact approximations  $x_n$  in both the simpler ( $x_n = x_0 + [q_1, Z_{n-1}]t_n$ ) and update ( $x_n = x_{n-1} + \alpha_n A^{-1}v_n$ ) approaches can be bounded by*

$$\frac{\|x_n - \hat{x}_n\|}{\|x\|} \leq cu\kappa(A)\kappa([q_1, Z_{n-1}]) \frac{\|\hat{x}_n\| + \|x_0\|}{\|x\|}, \quad (21)$$

*provided that  $\eta_n \equiv 1 - cu\kappa(A)\kappa([q_1, Z_{n-1}]) > 0$ .*

*Proof.* For the simpler approach, the result follows directly from Theorem 2.1. For the update approach, using (18) we have

$$x_n - \hat{x}_n = x - \hat{x}_n - A^{-1}r_n = (-A^{-1}F_n + G_n)U_n^{-1}D_n e$$

and the statement now follows from (11), (19) and (20).  $\square$

The bound (21) from Corollary 2.1 depends on the quantity  $(\|\hat{x}_n\| + \|x_0\|)/\|x\|$  (or more precisely on  $\|\hat{x}_n - x_0\|/\|x\|$ ), which is, however, strongly influenced by the conditioning of the upper triangular matrix  $U_n$ . As shown in Section 3, the matrix  $U_n$  can be ill-conditioned for a particular case  $[q_1, Z_{n-1}] = [q_1, V_{n-1}]$ , thus leading to an inherently less numerically stable scheme, whereas (under some assumptions) the scheme with  $[q_1, Z_{n-1}] = [\frac{x_0}{\|r_0\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$  gives rise to a well-conditioned triangular matrix  $U_n$ . In the following lemma we give bounds on  $\|\hat{x}_n - x_0\|$  in terms of the singular values of the matrix  $U_n$ .

**Lemma 2.1.** *In the simpler approach, we have*

$$\|\hat{x}_n - x_0\| \leq \|[q_1, Z_{n-1}]\| \|\hat{t}_n\| \leq \|[q_1, Z_{n-1}]\| \|(U_n + \Delta U_n)^{-1} D_n e\|,$$

and in the update approach,

$$\|\hat{x}_n - x_0\| \leq \|P_n D_n e\| \leq (1 + \text{cuk}(A)) \|[q_1, Z_{n-1}]\| \|U_n^{-1} D_n e\|.$$

The norms of  $(U_n + \Delta U_n)^{-1} D_n e$  and  $U_n^{-1} D_n e$  satisfy

$$\left. \begin{aligned} \|(U_n + \Delta U_n)^{-1} D_n e\| \\ \|U_n^{-1} D_n e\| \end{aligned} \right\} \leq \sqrt{2} \sum_{k=1}^n \frac{\|r_{k-1}\|}{\sigma_k(U_k)} \tag{22}$$

$$\leq \sqrt{2} \|A^{-1}\| \sum_{k=1}^n \frac{\eta_k^{-1} \|r_{k-1}\|}{\sigma_k([q_1, Z_{k-1}])},$$

provided that  $\eta_k \equiv 1 - \text{cuk}(A)\kappa([q_1, Z_{k-1}]) > 0$  for all  $k = 1, \dots, n$ .

*Proof.* Since  $e_k^T D_n e_k = \alpha_k$  and  $|\alpha_k| = \sqrt{\|r_{k-1}\|^2 - \|r_k\|^2} \leq \sqrt{2} \|r_{k-1}\|$ , we have

$$\begin{aligned} \|(U_n + \Delta U_n)^{-1} D_n e\| &\leq \sum_{k=1}^n \|(U_n + \Delta U_n)^{-1} D_n e_k\| \\ &\leq \sqrt{2} \sum_{k=1}^n \frac{\|r_{k-1}\|}{\sigma_k([U_n + \Delta U_n]_{1:k, 1:k})}, \end{aligned} \tag{23}$$

where  $[U_n + \Delta U_n]_{1:k,1:k}$  denotes the principal  $k \times k$  submatrix of  $U_n + \Delta U_n$ . Owing to (12), we can estimate the perturbation of  $[U_n]_{1:k,1:k} = U_k$  as  $\|[\Delta U_n]_{1:k,1:k}\| \leq cu\|U_k\|$ . Perturbation theory of singular values shows that

$$\begin{aligned} \sigma_k([U_n + \Delta U_n]_{1:k,1:k}) &\geq \sigma_k(U_k) - cu\|U_k\| \\ &\geq \sigma_k(A[q_1, Z_{k-1}]) - cu\|A\| \| [q_1, Z_{k-1}] \| \\ &\geq \sigma_N(A)\sigma_k([q_1, Z_{k-1}]) - cu\|A\| \| [q_1, Z_{k-1}] \|, \end{aligned} \quad (24)$$

which, together with (23), concludes the proof of the first inequality. The second inequality is proved analogously.  $\square$

The first estimate given in (22), which involves the minimal singular values of  $U_k$  ( $k = 1, \dots, n$ ), is quite sharp. However, the second estimate relating the minimal singular values of  $U_k$  to those of  $[q_1, Z_{k-1}]$  can be a large overestimate, as also observed in our numerical experiments in Section 3.

Theorems 2.1 and 2.2 indicate that *as soon as the backward error of the approximate solution in the simpler approach gets below  $cu\kappa(A)\kappa([q_1, Z_{n-1}])$ , the difference between the backward errors in the simpler and update approaches may become visible and can be expected to be up to the order of  $\kappa(A)$* . Based on our experience it is difficult to find an example where this difference is significant. Similarly to Sleijpen, van der Vorst and Moderitzki [26], we use here a model example, where  $A = G_1 D G_2^T \in \mathbb{R}^{100 \times 100}$  with  $D = \text{diag}(10^{-8}, 2 \cdot 10^{-8}, 3, 4, \dots, 100)$  and with  $G_1$  and  $G_2$  being Givens rotations over an angle of  $\frac{\pi}{4}$  in the  $(1, 10)$ -plane and the  $(1, 100)$ -plane, respectively; finally,  $b = e$ . The numerical experiments were performed in MATLAB using double precision arithmetic ( $u \approx 10^{-16}$ ), and the zero vector was chosen as the initial guess  $x_0$ . In Figure 1 we have plotted the norm-wise backward errors  $\|b - A\hat{x}_n\| / (\|A\| \|\hat{x}_n\|)$  (solid and dashed lines) and the relative 2-norms of the errors  $\|x - \hat{x}_n\| / \|x\|$  (dash-dotted and dotted lines) for Simpler GMRES and refined ORTHODIR, respectively. The reciprocals of the condition numbers of the basis  $[q_1, Z_{n-1}]$ , the triangular matrix  $U_n$  and the system matrix  $A$  are depicted by dashed, dashed-dotted and dotted lines. The same quantities for SGMRES/RB and refined ORTHOMIN are reported in Figure 2. We see that the actual backward errors are close until where they stagnate: for refined ORTHODIR and refined ORTHOMIN this happens approximately at a level close to  $u\kappa(A)$ , while for Simpler GMRES and SGMRES/RB we have stagnation on the roundoff unit level. In contrast, the 2-norms of the errors stagnate on the  $u\kappa(A)$  level in all schemes considered.

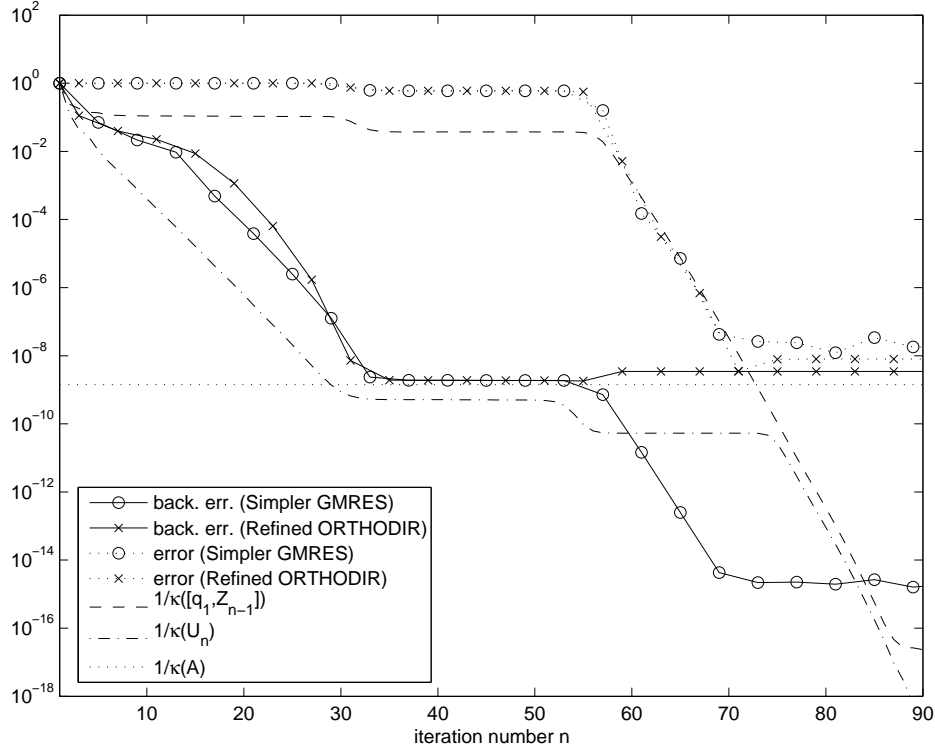


Figure 1: The test problem solved by Simpler GMRES and refined ORTHODIR.

### 3 Choice of basis and stability

In this section we discuss the two main particular choices for the matrix  $Z_{n-1}$  leading to different algorithms for the simpler and update schemes. For the sake of simplicity, we assume exact arithmetic here. First, we choose  $Z_{n-1} = V_{n-1}$ , which leads to the Simpler GMRES method of Walker and Zhou [29] and to the refined version of ORTHODIR by Young and Jea [31], respectively. Hence, we choose  $\{q_1, v_1, \dots, v_{n-1}\}$  as a basis of  $\mathcal{K}_n(A, r_0)$ . To be sure that such a choice is adequate, we state the following simple lemma.

**Lemma 3.1.** *Let  $v_1, \dots, v_{n-1}$  be an orthonormal basis of  $AK_{n-1}(A, r_0)$ ,  $r_0 \notin AK_{n-1}(A, r_0)$ . Then the vectors  $q_1, v_1, \dots, v_{n-1}$  form a basis of  $\mathcal{K}_n(A, r_0)$ .*

*Proof.* It follows from the assumption  $r_0 \notin AK_{n-1}(A, r_0)$  implying that  $q_1 \notin AK_{n-1}(A, r_0) = \text{span}\{v_1, \dots, v_{n-1}\}$ .  $\square$

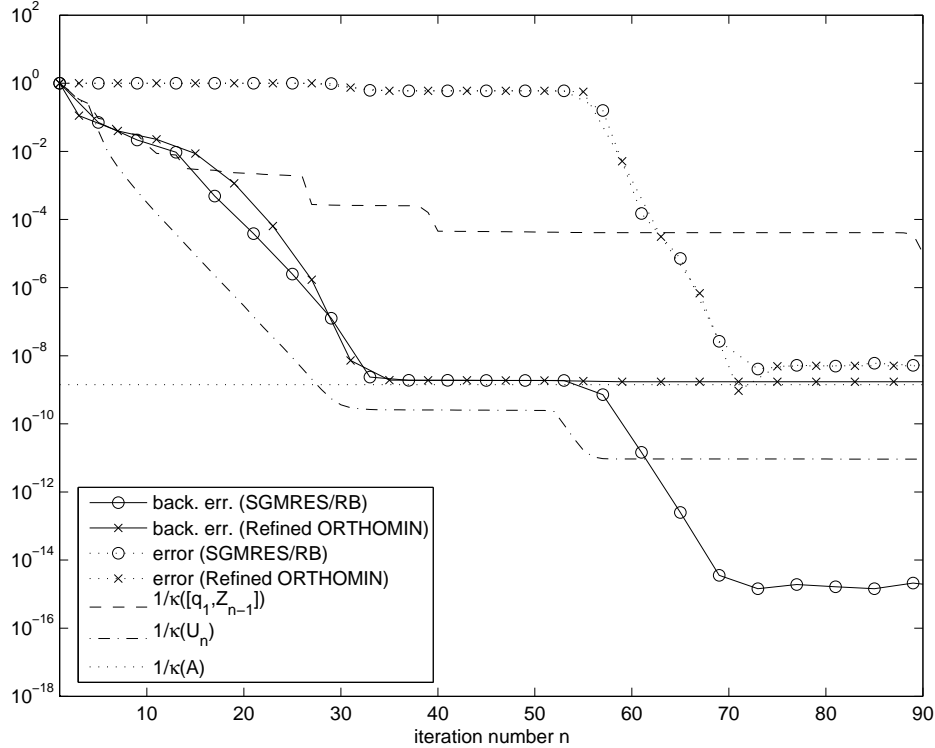


Figure 2: The test problem solved by SGMRES/RB and refined ORTHOMIN.

Note that if  $r_0 \in AK_n(A, r_0)$ , then the condition (2) yields  $x_n = A^{-1}b$ ,  $r_n = 0$ , and any implementation of a minimum residual method will terminate. Lemma 3.1 ensures that it makes sense to build an orthonormal basis  $V_n$  of  $AK_n(A, r_0)$  by the successive orthogonalization of the columns of the matrix  $A[q_1, V_{n-1}]$  via (4). It reflects the fact that, for any initial residual  $r_0$ , both Simpler GMRES and ORTHODIR converge (in exact arithmetic) to the exact solution; see [31]. However, as observed by Liesen, Rozložník and Strakoš [18], this choice of the basis is not very suitable from the stability point of view. This shortcoming is reflected by the unbounded growth of the condition number of  $[q_1, V_{n-1}]$  discussed next. The upper bound we give was also derived in [29].

**Theorem 3.1.** *Let  $r_0 \notin AK_{n-1}(A, r_0)$ . Then the condition number of  $[q_1, V_{n-1}]$*

satisfies

$$\frac{\|r_0\|}{\|r_{n-1}\|} \leq \kappa([q_1, V_{n-1}]) \leq 2 \frac{\|r_0\|}{\|r_{n-1}\|}.$$

*Proof.* Since  $r_{n-1} = (I - V_{n-1}V_{n-1}^T)r_0$ , it is easy to see that  $r_{n-1}$  is the residual of the least squares problem  $V_{n-1}y \approx r_0$ . The statement follows from Theorem 3.2 of [18].  $\square$

The conditioning of  $[q_1, V_{n-1}]$  is thus related to the convergence of the method; in particular, it is inversely proportional to the actual relative norm of the residual. Hence, if the residual is small enough, Simpler GMRES and refined ORTHODIR behave unstably. In practice, this difficulty can be counteracted by frequent restarts.

Now we turn to the second choice,  $Z_{n-1} = [\frac{r_1}{\|r_1\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$ , which leads to SGMRES/RB (which we propose here as a more stable counterpart of Simpler GMRES) and to the refined version of ORTHOMIN by Vinsome [28] known also under the name GCR; see Eisenstat, Elman and Schultz [9, 8]. We have  $[q_1, Z_{n-1}] = R_n B_n^{-1}$ , where  $B_n \equiv \text{diag}(\|r_0\|, \dots, \|r_{n-1}\|)$ , i.e., we choose scaled residuals  $r_0, \dots, r_{n-1}$  as the basis of  $\mathcal{K}_n(A, r_0)$ . To be sure that such a choice is adequate, we state the following result.

**Lemma 3.2.** *Let  $v_1, \dots, v_{n-1}$  be an orthonormal basis of  $AK_{n-1}(A, r_0)$ ,  $r_0 \notin AK_{n-1}(A, r_0)$  and  $r_k = (I - V_k V_k^T)r_0$ , where  $V_k \equiv [v_1, \dots, v_k]$ ,  $k = 1, 2, \dots, n-1$ . Then the following statements are equivalent:*

1.  $\|r_k\| < \|r_{k-1}\|$  for all  $k = 1, \dots, n-1$ ,
2.  $r_0, \dots, r_{n-1}$  are linearly independent.

*Proof.* Since  $r_0 \notin AK_{n-1}(A, r_0) = \mathcal{R}(V_{n-1})$ ,  $r_k \neq 0$  for all  $k = 0, 1, \dots, n-1$ . It is clear that  $\|r_k\| < \|r_{k-1}\|$  if and only if  $\langle r_{k-1}, v_k \rangle \neq 0$ . If that holds for all  $k = 1, \dots, n-1$  the diagonal matrix  $D_{n-1}$  is nonsingular. Using the relation (6) we find that  $R_n[L_{n,n-1}, e_n] = [V_{n-1}D_{n-1}, r_{n-1}]$ . Since  $r_{n-1} \perp V_{n-1}$ , the matrix  $[V_{n-1}D_{n-1}, r_{n-1}]$  has orthogonal nonzero columns, and hence its rank equals  $n$ . Moreover,  $\text{rank}([L_{n,n-1}, e_n]) = n$  and thus  $\text{rank}(R_n) = n$ , i.e.,  $r_0, \dots, r_{n-1}$  are linearly independent. Conversely, from the same matrix relation we find that if  $r_0, \dots, r_{n-1}$  are linearly independent, then  $\text{rank}([V_{n-1}D_{n-1}, r_{n-1}]) = n$ , and hence  $D_{n-1}$  is nonsingular, which proves that  $\|r_k\| < \|r_{k-1}\|$  for all  $k = 1, \dots, n-1$ .  $\square$



Therefore if the method does not stagnate, i.e., if the 2-norms of the residuals  $r_0, \dots, r_{n-1}$  are strictly monotonously decreasing, then  $r_0, \dots, r_{n-1}$  are linearly independent. In this case, we can build an orthonormal basis  $V_n$  of  $AK_n(A, r_0)$  by the successive orthogonalization of the columns of  $AR_nB_n^{-1}$  via (4). If  $r_0 \in AK_{n-1}(A, r_0)$ , we have an exact solution of (1), and the method stops with  $x_{n-1} = A^{-1}b$ .

Several conditions for the non-stagnation of the minimum residual method have been given in the literature. For example, Eisenstat, Elman and Schultz [8, 9] show that GCR (and hence any minimum residual method) does not stagnate if the symmetric part of  $A$  is positive definite, i.e., if the origin is not contained in the field of values of  $A$ . See also Greenbaum and Strakoš [13] for a different proof, and Eiermann and Ernst [7]. Several other conditions can be found in Simoncini and Szyld [25] and the references therein. If stagnation occurs, the residuals are no longer linearly independent, and thus the method prematurely breaks down. In particular, if  $0 \in \mathcal{F}(A)$ , choosing  $x_0$  such that  $r_0 \in \mathcal{F}(A)$  leads to a breakdown in the first step. This was first pointed out by Young and Jea [31] with a simple  $2 \times 2$  example.

However, as shown in the following theorem, when the minimum residual method does not stagnate, the columns of  $R_nB_n^{-1}$  are a reasonable choice for the basis of  $\mathcal{K}_n(A, r_0)$ .

**Theorem 3.2.** *If  $r_0 \notin AK_{n-1}(A, r_0)$ , the condition number of  $R_nB_n^{-1}$  satisfies*

$$1 \leq \kappa(R_nB_n^{-1}) \leq \sqrt{n} \gamma_n, \quad \gamma_n \equiv \sqrt{1 + \sum_{k=1}^{n-1} \frac{\|r_{k-1}\|^2 + \|r_k\|^2}{\|r_{k-1}\|^2 - \|r_k\|^2}}. \quad (25)$$

*Proof.* From (6) it follows that

$$R_nB_n^{-1}[Q_{n,n-1}, e_n] = [V_{n-1}, \frac{r_{n-1}}{\|r_{n-1}\|}], \quad Q_{n,n-1} \equiv B_nL_{n,n-1}D_{n-1}^{-1}.$$

Since  $[V_{n-1}, \frac{r_{n-1}}{\|r_{n-1}\|}]$  is an orthonormal matrix, we have from Theorem 3.3.16 of [16]

$$\begin{aligned} 1 &= \sigma_n([V_{n-1}, \frac{r_{n-1}}{\|r_{n-1}\|}]) \leq \sigma_n(R_nB_n^{-1}) \|[Q_{n,n-1}, e_n]\| \\ &\leq \sigma_n(R_nB_n^{-1}) \|[Q_{n,n-1}, e_n]\|_F. \end{aligned}$$

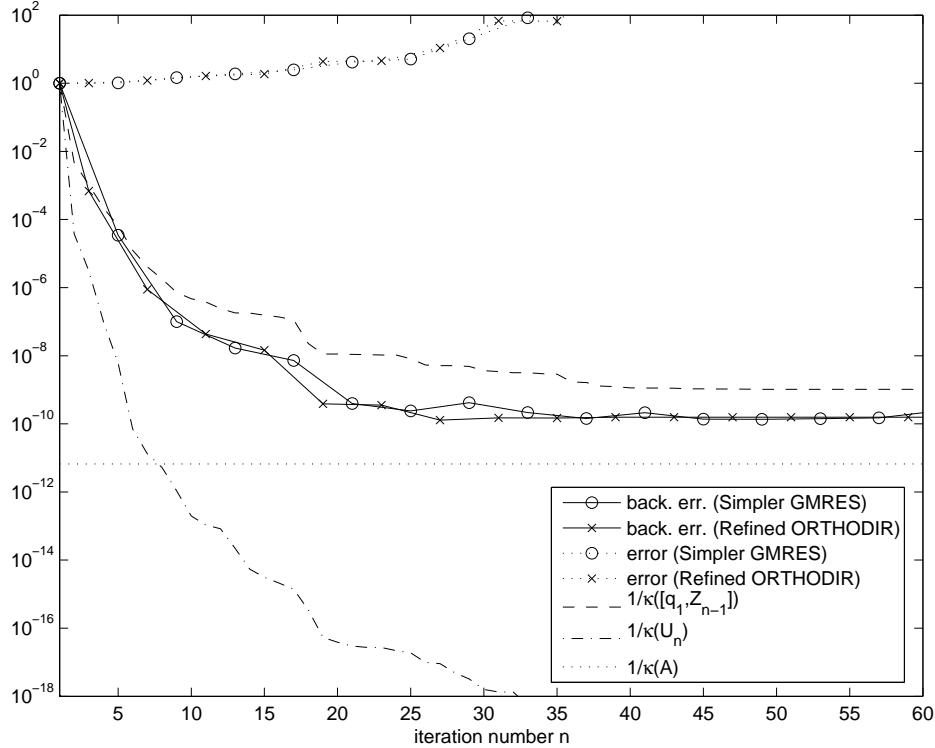


Figure 3: The test problem FS1836 solved by Simpler GMRES and refined ORTHODIR.

The value of  $\| [Q_{n,n-1}, e_n] \|_F$  can be directly computed as

$$\| [Q_{n,n-1}, e_n] \|_F = \sqrt{1 + \sum_{k=1}^{n-1} \frac{\|r_{k-1}\|^2 + \|r_k\|^2}{\|r_{k-1}\|^2 - \|r_k\|^2}} = \gamma_n,$$

since  $\alpha_k^2 = \|r_{k-1}\|^2 - \|r_k\|^2$ . The statement follows using  $\|R_n B_n^{-1}\| \leq \sqrt{n}$ .  $\square$

We define the quantity  $\gamma_n$  in (25) as the *stagnation factor*. The conditioning of  $R_n B_n^{-1}$  is thus related to the convergence of the method, but in contrast to the conditioning of  $[q_1, V_{n-1}]$ , it is related to the intermediate decrease of the residual norms, not to the residual decrease with respect to the initial residual.

We illustrate our theoretical results by a numerical example using the ill-conditioned matrix FS1836 ( $\|A\| \approx 1.18 \cdot 10^9$ ,  $\|A^{-1}\| \approx 1.47 \cdot 10^2$ ) ob-

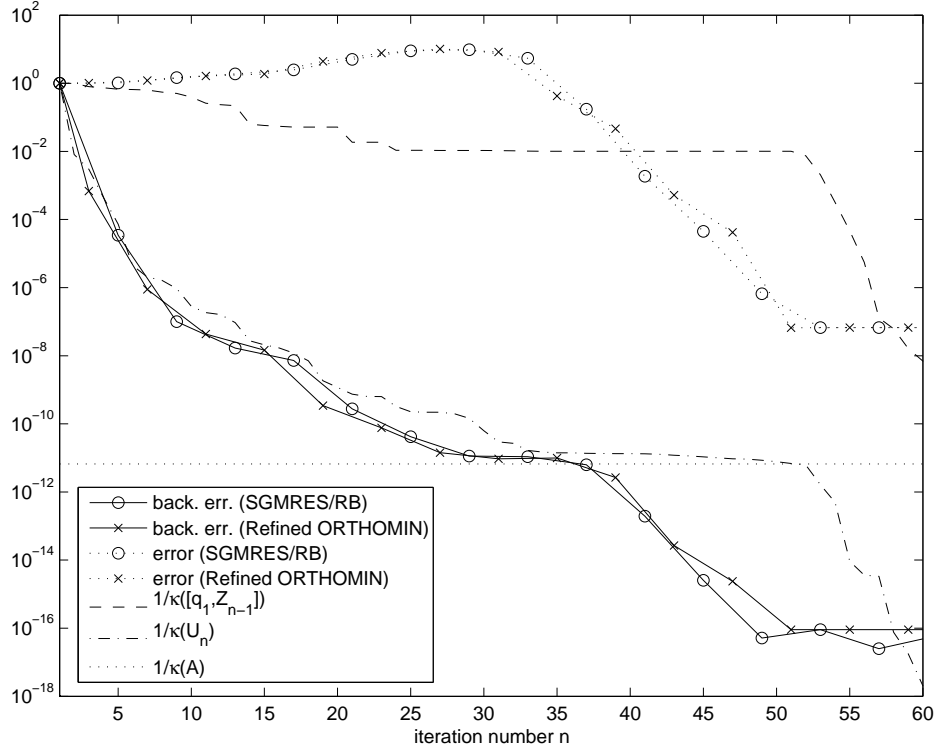


Figure 4: The test problem FS1836 solved by SGMRES/RB and refined ORTHOMIN.

tained from the Matrix Market [1] with the right-hand side  $b = Ae$  (see also the experiments in [18], where the relative residual norms were reported). In Figures 3 and 4, we show again the normwise backward error  $\|b - A\hat{x}_n\|/(\|A\|\|\hat{x}_n\|)$  (solid lines with circles and crosses) and the relative 2-norms of the error  $\|x - \hat{x}_n\|/\|x\|$  (dotted lines with circles and crosses) for the choice  $[q_1, Z_{n-1}] = [q_1, V_{n-1}]$  that corresponds to Simpler GMRES and refined ORTHODIR, and for  $[q_1, Z_{n-1}] = [\frac{r_0}{\|r_0\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$  corresponding to SGMRES/RB and refined ORTHOMIN, respectively. The reciprocals of the condition numbers of the basis  $[q_1, Z_{n-1}]$ , the triangular matrix  $U_n$  and the system matrix  $A$  are depicted by dashed, dashed-dotted and dotted lines. We see that the backward errors and the error norms are almost identical for the simpler and update approaches. This can be observed in most cases leading to practically negligible difference between Simpler GMRES and refined ORTHODIR, and SGMRES/RB and refined ORTHOMIN,

respectively. Figure 3 illustrates our theoretical considerations and shows that, after some initial reduction, *the backward error of Simpler GMRES and refined ORTHODIR may stagnate on a significantly higher level than the backward error of SGMRES/RB or refined ORTHOMIN, which stagnates on a level proportional to the roundoff unit*, as shown in Figure 4. Due to Theorem 3.1, after some initial phase, the norms of the errors start to diverge in Simpler GMRES and refined ORTHODIR, while for SGMRES/RB and refined ORTHOMIN we have a stagnation on a level approximately proportional to  $u\kappa(A)$ . The difference is clearly caused by the choice of the basis  $[q_1, Z_{n-1}]$ , which has an effect on the conditioning of the matrix  $U_n$ . We see that  $[q_1, Z_{n-1}] = [\frac{r_0}{\|r_0\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$  remains well-conditioned up to the very end of the iteration process, while the conditioning of  $[q_1, V_{n-1}]$  is linked to the convergence of Simpler GMRES and may lead to a very ill-conditioned triangular matrix  $U_n$ . Consequently the approximate solution  $\hat{x}_n$  computed from (8) becomes inaccurate and its error starts to diverge. Since the stagnation factor  $\gamma_n \approx 55.8$  (for  $n = 50$ ), the matrix  $U_n$  remains well-conditioned, and this problem does not occur in the SGMRES/RB method.

## 4 Conclusions

In this paper we have studied the numerical behavior of several minimum residual methods mathematically equivalent to GMRES. Two general formulations have been analyzed: the simpler approach that does not require an upper Hessenberg factorization and the update approach which is based on generating a sequence of appropriately computed direction vectors. It has been shown that for the simpler approach our analysis leads to an upper bound for the backward error proportional to the roundoff unit, whereas for the update approach the same quantity can be bounded by a term proportional to the condition number of  $A$ . Although our analysis suggests that the difference between both may be up to the order of  $\kappa(A)$ , in practice they behave very similarly, and it is very difficult to find a concrete example with a significant difference in the limiting accuracy measured by the normwise backward error of the approximate solutions  $x_n$ . Our first test problem displayed in Figures 1 and 2 is such a rare example. Moreover, when looking at the errors, we note that both approaches lead essentially to the same accuracy of  $x_n$ .

We have indicated that the choice of the basis  $[q_1, Z_{n-1}]$  is the most

important issue for the stability of the considered schemes. Our analysis supports the well-known fact that even when implemented with the best possible orthogonalization techniques Simpler GMRES and ORTHODIR are inherently less stable due to the choice  $[q_1, Z_{n-1}] = [q_1, V_{n-1}]$  for the basis. The situation becomes significantly better, when we use the residual basis  $[q_1, Z_{n-1}] = [\frac{r_0}{\|r_0\|}, \dots, \frac{r_{n-1}}{\|r_{n-1}\|}]$ . This choice leads to the popular GCR (ORTHOMIN, GMRESR) method, which is widely used in applications. Assuming some reasonable residual decrease (which happens almost always in finite precision arithmetic), we have shown that this scheme is quite efficient and proposed a conditionally backward stable variant (called SGMRES/RB here). Our theoretical results in a sense justify the use of the GCR method in practical computations. In this paper we studied only the unpreconditioned implementations. The implications for the preconditioned GCR scheme will be discussed elsewhere.

## 5 Acknowledgments

We would like to thank Julien Langou, Yvan Notay and Kees Vuik for valuable discussions during the preparation of the paper.

## References

- [1] *Matrix Market*. URL: <http://math.nist.gov/MatrixMarket>.
- [2] W. E. ARNOLDI, *The principle of minimized iterations in the solution of the matrix eigenvalue problem*, Quart. Appl. Math., 9 (1951), pp. 17–29.
- [3] S. F. ASHBY AND M. H. GUTKNECHT, *A matrix analysis of conjugate gradient algorithms*, in Advances in Numerical Methods for Large Sparse Sets of Linear Systems, Parallel Processing for Scientific Computing, M. Natori and T. Nodera, eds., vol. 9, Yokohama, Japan, 1993, Keio University, pp. 32–47.
- [4] S. F. ASHBY, T. A. MANTEUFFEL, AND P. E. SAYLOR, *A taxonomy for conjugate gradient methods*, SIAM J. Numer. Anal., 27 (1990), pp. 1542–1568.

- [5] A. BJÖRCK, *Solving linear least squares problems by Gram–Schmidt orthogonalization*, BIT, 7 (1967), pp. 1–21.
- [6] J. DRKOŠOVÁ, A. GREENBAUM, M. ROZLOŽNÍK, AND Z. STRAKOŠ, *Numerical stability of GMRES*, BIT, 35 (1995), pp. 309–330.
- [7] M. EIERMANN AND O. ERNST, *Geometric aspects of the theory of Krylov subspace methods*, Acta Numerica, (2001), pp. 251–312.
- [8] 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), pp. 345–357.
- [9] H. C. ELMAN, *Iterative methods for large sparse nonsymmetric systems of linear equations*, PhD thesis, New Haven, 1982.
- [10] D. K. FADDEEV AND V. N. FADDEEVA, *Computational Methods of Linear Algebra*, Fizmatgiz, Moskow, 1960. in russian.
- [11] A. GREENBAUM, *Estimating the attainable accuracy of recursively computed residual methods*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 535–551.
- [12] A. GREENBAUM, M. ROZLOŽNÍK, AND Z. STRAKOŠ, *Numerical behaviour of the modified Gram-Schmidt GMRES implementation*, BIT, 37 (1997), pp. 706–719.
- [13] A. GREENBAUM AND Z. STRAKOŠ, *Matrices that generate the same Krylov residual spaces*, in Recent Advances in Iterative Methods, G. H. Golub, A. Greenbaum, and M. Luskin, eds., New York, 1994, Springer-Verlag, pp. 95–119.
- [14] M. H. GUTKNECHT AND Z. STRAKOŠ, *Accuracy of two three-term and three two-term recurrences for Krylov space solvers*, SIAM J. Matrix Anal. Appl., 22 (2000), pp. 213–229.
- [15] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.
- [16] R. A. HORN AND C. R. JOHNSON, *Topics in Matrix Analysis*, Cambridge University Press, new ed., 1994.

- [17] K. C. JEA AND D. M. YOUNG, *On the simplification of generalized conjugate-gradient methods for nonsymmetrizable linear systems*, Linear Algebra Appl., 52 (1983), pp. 399–417.
- [18] J. LIESEN, M. ROZLOŽNÍK, AND Z. STRAKOŠ, *Least squares residuals and minimal residual methods*, SIAM J. Sci. Comp., 23 (2002), pp. 1503–1525.
- [19] G. MEURANT, *Computer Solution of Large Linear Systems*, North Holland, 1999.
- [20] C. C. PAIGE, M. ROZLOŽNÍK, AND Z. STRAKOŠ, *Modified Gram-Schmidt (MGS), least squares, and backward stability of MGS-GMRES*, SIAM J. Matrix Anal. Appl., 28 (2006), pp. 264–284.
- [21] C. C. PAIGE AND M. A. SAUNDERS, *Solution of sparse indefinite systems of linear equations*, SIAM J. Numer. Anal., 12 (1975), pp. 617–629.
- [22] M. ROZLOŽNÍK AND Z. STRAKOŠ, *Variants of residual minimizing Krylov subspace methods*, in Proceedings of the 6th Summer School Software and Algorithms of Numerical Mathematics, I. Marek, ed., 1995, pp. 208–225.
- [23] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, SIAM, 2nd ed., 2003.
- [24] Y. SAAD AND M. H. SCHULTZ, *GMRES: a generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Stat. Comp., 7 (1986), pp. 856–869.
- [25] V. SIMONCINI AND D. B. SZYLD, *New conditions for non-stagnation of minimal residual methods*, Tech. Rep. 07-4-17, Apr. 2007.
- [26] G. L. G. SLEIJPEN, H. A. VAN DER VORST, AND J. MODERSITZKI, *Differences in the effects of rounding errors in Krylov solvers for symmetric indefinite linear systems*, SIAM J. Matrix Anal. Appl., 22 (2000), pp. 726–751.
- [27] H. A. VAN DER VORST AND C. VUIK, *GMRESR: a family of nested GMRES methods*, Numer. Linear Algebra Appl., 1 (1994), pp. 369–386.

- [28] P. K. W. VINSOME, *Orthomin, an iterative method for solving sparse sets of simultaneous linear equations*, in Proceedings Fourth Symposium on Reservoir Simulation, SPE of AIME, Los Angeles, Feb. 1976.
- [29] H. F. WALKER AND L. ZHOU, *A simpler GMRES*, Numer. Linear Algebra Appl., 1 (1994), pp. 571–581.
- [30] J. H. WILKINSON, *Rounding Errors in Algebraic Processes*, Prentice Hall, Inc., New Jersey, 1963.
- [31] D. M. YOUNG AND K. C. JEA, *Generalized conjugate gradient acceleration of nonsymmetrizable iterative methods*, Linear Algebra Appl., 34 (1980), pp. 159–194.