

# A comparison of some GMRES-like methods

Report 90-22

C. Vuik



Technische Universiteit Delft  
Delft University of Technology

Faculteit der Technische Wiskunde en Informatica  
Faculty of Technical Mathematics and Informatics

ISSN 0922-5641

Copyright © 1990 by the Faculty of Technical Mathematics and Informatics, Delft, The Netherlands.

No part of this Journal may be reproduced in any form, by print, photoprint, microfilm, or any other means without permission from the Faculty of Technical Mathematics and Informatics, Delft University of Technology, The Netherlands.

Copies of these reports may be obtained from the bureau of the Faculty of Technical Mathematics and Informatics, Julianalaan 132, 2628 BL Delft, phone +31 15784568.

A selection of these reports is available in PostScript form at the Faculty's anonymous ftp-site. They are located in the directory /pub/publications/tech-reports at [ftp.twi.tudelft.nl](ftp://ftp.twi.tudelft.nl)

# A comparison of some GMRES-like methods

C. Vuik

Department of Technical Mathematics and Computer Science  
Delft University of Technology

Julianalaan 132

Delft, the Netherlands

and

H.A. Van der Vorst

Mathematical Institute

University of Utrecht

Budapestlaan 6

Utrecht, the Netherlands

April 7, 1994

## Abstract

GMRES and CGS are well-known iterative methods for the solution of certain sparse linear systems with a non-symmetric matrix. These methods have been compared experimentally in many studies and specific observations on their convergence behaviour have been reported. A new iterative method to solve a non-symmetric system is proposed by Eirola and Nevanlinna. The purpose of this paper is to investigate this method and to compare it with GMRES. We have seen problems for which this method is more efficient than GMRES. The original method has as drawbacks that it is not scaling invariant and that it may suffer from numerical instability but it will be shown that these deficiencies can be repaired. A method proposed by Broyden seems to be somehow related to the new method and is therefore included in the comparison.

# Introduction

In this paper we compare the GMRES-method [6], the EN-method [3] and the B-method [1]. Our main motivation to study the EN-method is that it deepens the insight in projection-type methods, which hopefully leads to better iterative methods. Descriptions and some relevant properties of these methods are given in Section 1. In Section 2 we describe numerical experiments for EN, which motivate the theoretical analysis of Section 3. In that section we give a relation between the EN- and GMRES-method. Subsequently we compare the efficiency of both methods. Though in some cases the EN-method is more efficient than the GMRES-method this is not the case in general. In Section 4 we show that the convergence and the stability properties of EN are not scaling invariant, as they are for GMRES and other projection methods and we also show how this can be repaired to the advantage of the EN-method. Furthermore, we describe some problems for which EN diverges GMRES converges. In Section 5 we consider a variant of the EN-method, which is algebraically equivalent to the GMRES-method. This enables us to make a better comparison between GMRES and EN, and it gives more insight in GMRES. Finally, in Section 6 we compare the EN-method with the B-method and a general class of methods given in [2]. Furthermore we compare the efficiency of B and GMRES. From these comparisons it appears that the most efficient and robust method is the implementation of the full GMRES method as described in, e.g., [6] and [7]. However, it appears from experiments that if the iterative methods (EN and GMRES) are restarted then EN can be much more efficient than GMRES. This aspect is subject of further study and is not reported in this paper.

## 1 GMRES, EN and B-method

The GMRES-method is originally proposed in [6]. We use results in [5] for understanding the convergence behaviour of GMRES. Consider the linear system  $Ax = b$  with  $x, b \in R^n$  and  $A \in R^{n \times n}$  is nonsingular. The Krylov subspace  $K^k(A; r_0)$  is defined by  $K^k(A; r_0) = \text{span} \{r_0, Ar_0, \dots, A^{k-1}r_0\}$ . The k-th iterate  $x_k$  is written as  $x_k = x_0 + z_k$  where  $z_k \in K^k(A; r_0)$  and  $r_0 = b - Ax_0$ . In the GMRES-method the vector  $z_k$  is chosen as the vector which solves the linear least-squares problem

$$z_k = \underset{z \in K^k(A; r_0)}{\text{arg min}} \|b - A(x_0 + z)\|_2 . \quad (1)$$

From this definition it follows that

$$\|r_k\|_2 = \underset{z \in K^k(A; r_0)}{\text{min}} \|b - Ax_0 - Az\|_2 = \underset{\alpha_1, \dots, \alpha_k \in R}{\text{min}} \|r_0 + \sum_{i=1}^k \alpha_i A^i r_0\|_2 . \quad (2)$$

In the EN-method we take a different splitting of the matrix in each iteration step:

$$A = H_k^{-1} - R_k ,$$

which leads to the basic iteration method

$$x_k = x_{k-1} + H_k r_{k-1} .$$

The key idea is to improve  $H_k$  from step to step by (cheap) rank-1 updates:

$$H_k = H_{k-1} + u_{k-1}v_{k-1}^T .$$

For the  $k$ -th step this leads to

$$\begin{aligned} r_k &= r_{k-1} - A(H_{k-1} + u_{k-1}v_{k-1}^T)r_{k-1} \\ &= (I - AH_{k-1})r_{k-1} - \mu_{k-1}Au_{k-1} \end{aligned}$$

with  $\mu_{k-1} = v_{k-1}^T r_{k-1}$ .

The ideal choice for  $u_{k-1}$  would have been, such that

$$\begin{aligned} \mu_{k-1}Au_{k-1} &= (I - AH_{k-1})r_{k-1} , \\ \text{or } \mu_{k-1}u_{k-1} &= A^{-1}(I - AH_{k-1})r_{k-1} . \end{aligned}$$

If  $H_{k-1}^{-1}$  defines a suitable splitting for  $A$  then  $A^{-1}$  could be replaced by  $H_{k-1}$  and this motivates the choice for  $u_{k-1}$ :

$$u_{k-1} = H_{k-1}(I - AH_{k-1})r_{k-1} .$$

The choice for  $v_{k-1}$  now follows by minimizing  $\| r_k \|_2$  as a function of  $\mu_{k-1}$ :

$$\mu_{k-1} = (Au_{k-1})^T(I - AH_{k-1})r_{k-1} / \| Au_{k-1} \|_2^2$$

so that

$$v_{k-1} = \frac{1}{\| Au_{k-1} \|_2^2}(I - AH_{k-1})^T Au_{k-1}$$

is an obvious choice.

This leads to the following algorithm ([3]: p.512,513):

1. given  $x_0, H_0$ , compute  $r_0$  and take  $k = 0$ ,
2.  $E_k = I - AH_k$ ,  $u_k = H_k E_k r_k$ ,  $v_k = E_k^T Au_k / \| Au_k \|_2^2$ ,
3.  $H_{k+1} = H_k + u_k v_k^T$ ,  $x_{k+1} = x_k + H_{k+1} r_k$ ,  $r_{k+1} = b - Ax_{k+1}$ ,
4. stop if  $\| r_{k+1} \|_2$  is small enough, otherwise  $k := k + 1$  and return to step 2.

The only difference between EN and GMRES is the choice of  $u_k$ . By taking  $u_k = H_k r_k$ , instead of  $u_k = H_k E_k r_k$ , we obtain an iterative method algebraically equivalent to GMRES.

The following equalities and definitions will be used in our analysis:

$$c_k \equiv Au_k / \| Au_k \|_2, \tag{3}$$

$$E_{k+1} = (I - P_k)E_0, \quad P_k = \sum_{i=0}^k c_i c_i^T \quad \text{and} \quad c_i^T c_j = 0 \quad \text{for } i \neq j, \tag{4}$$

$$r_{k+1} = E_{k+1} r_k. \tag{5}$$

Equation (4) only holds if all  $H_k$  are nonsingular. Therefore, in the case that  $H_{k+1}$  is singular whereas  $H_k$  is nonsingular we take  $H_{k+1} = H_k$  (see [3]: p.518). The following property may be

used to check whether  $H_{k+1}$  is singular.

Property 1.6 ([3]: p.518, Proposition 2.1)

Assume  $H_k$  is nonsingular. Then  $H_{k+1}$  is singular if and only if  $c_k^T E_0 r_k = 0$ .

The description of the algorithm given above is suitable for analysis, however in order to save computational work we prefer the following implementation given in ([3]: p.519):

at step  $k$   $x_k, r_k, u_0, \dots, u_{k-1}, c_0, \dots, c_{k-1}$  are known, then compute

1.  $\alpha_i = c_i^T (r_k - AH_0 r_k)$  for  $i = 0, \dots, k-1$ ,  $\eta = H_0 r_k + \sum_{i=0}^{k-1} \alpha_i u_i$ ,  $\xi = r_k - A\eta$ ,
2.  $\beta_i = c_i^T (\xi - AH_0 \xi)$  for  $i = 0, \dots, k-1$ ,  $u_k = \tau \left( H_0 \xi + \sum_{i=0}^{k-1} \beta_i u_i \right)$ ,  $c_k = Au_k$ ,  
where  $\tau$  is such that  $\|c_k\|_2 = 1$ ,
3.  $x_{k+1} = x_k + \eta + u_k c_k^T \xi$ ,  $r_{k+1} = \xi - c_k c_k^T \xi$ .

In the sequel, EN1 denotes the given implementation.

In another implementation given in ([3]: p.519),  $\xi$  and  $c_k$  are computed as follows:  $\xi = r_k - AH_0 r_k - \sum_{i=0}^{k-1} \alpha_i c_i$  and  $c_k = \tau \left( AH_0 \xi + \sum_{i=0}^{k-1} \beta_i c_i \right)$ . This implementation is used, in situations where it is more efficient to compute a linear combination of  $k+1$  vectors instead of multiplying one vector by  $A$ . Note that  $\xi$  is the component of  $r_k - AH_0 r_k$  orthogonal to span  $\{c_0, \dots, c_{k-1}\}$ . Hence  $\beta_i$  is equal to  $-c_i^T AH_0 \xi$  which implies that  $c_k$  is the normalized component of  $AH_0 \xi$  orthogonal to span  $\{c_0, \dots, c_{k-1}\}$ . In this implementation the vectors  $\xi$  and  $c_k$  are made orthogonal by the Gram Schmidt process. For stability reasons we propose the following implementation (EN2) based on the modified Gram Schmidt process:

1.  $\xi^{(0)} = (I - AH_0) r_k$ ,  $\eta^{(0)} = H_0 r_k$ ,  
 $\alpha_i = c_i^T \xi^{(i)}$ ,  $\xi^{(i+1)} = \xi^{(i)} - \alpha_i c_i$ ,  $\eta^{(i+1)} = \eta^{(i)} + \alpha_i u_i$ ,  $i = 0, \dots, k-1$ ,
2.  $c_k^{(0)} = AH_0 \xi^{(k)}$ ,  $u_k^{(0)} = H_0 \xi^{(k)}$ ,  
 $\beta_i = -c_i^T c_k^{(i)}$ ,  $c_k^{(i+1)} = c_k^{(i)} + \beta_i c_i$ ,  $u_k^{(i+1)} = u_k^{(i)} + \beta_i u_i$ ,  $i = 0, \dots, k-1$ ,  
 $c_k = c_k^{(k)} / \|c_k^{(k)}\|_2$ ,  $u_k = u_k^{(k)} / \|c_k^{(k)}\|_2$ ,
3.  $x_{k+1} = x_k + \eta^{(k)} + u_k c_k^T \xi^{(k)}$ ,  $r_{k+1} = (1 - c_k c_k^T) \xi^{(k)}$ .

In our experiments the stability properties of EN1 and EN2 have appeared to be more or less equivalent.

In the B-method also a nonsingular matrix  $H_0 \in R^{n \times n}$  must be specified which again is viewed as an approximation to the inverse of  $A$ .

The algorithm runs as follows ([1]: p.94):

1. given  $x_0, H_0$ , compute  $r_0$  and take  $k = 0$ ,

2.  $p_k = H_k r_k$ ,  $x_{k+1} = x_k + p_k$ ,  $r_{k+1} = b - Ax_{k+1}$ ,
3.  $y_k = r_k - r_{k+1}$ ,
4.  $H_{k+1} = H_k - (H_k y_k - p_k) p_k^T H_k / p_k^T H_k y_k$ ,  $k := k + 1$  and return to step 2.

## 2 Numerical experiments

In order to get some idea of the convergence behaviour of the EN-method we report on some numerical experiments. The numerical experiments have been carried out on a HP9000-845 computer in double precision arithmetic (about 15 decimal places). Our test matrices and right-hand sides are taken from ([5]: p.16,17). These matrices are of the form  $A = SBS^{-1}$  with  $A, S, B \in R^{100 \times 100}$ . We have selected  $S$  to be equal to

$$S = \begin{pmatrix} 1 & \beta & & \emptyset \\ & 1 & \ddots & \\ & & \ddots & \ddots & \beta \\ \emptyset & & & & 1 \end{pmatrix}.$$

The system  $Ax = b$  is solved for right-hand sides, such that  $x = (1, \dots, 1)^T$  (experiments with other choices of  $x$  show more or less the same convergence behaviour). In these experiments we take  $H_0 = I$  and  $x_0 = (0, \dots, 0)^T$ . The matrices in our testset are as follows (the numbering refers to the problems in ([5]: p.17]):

Problem P6:

$$B = \begin{pmatrix} 1 & & & & \\ & 1 + \alpha & & & \\ & & 3 & & \emptyset \\ & & & 4 & \\ & & & & \ddots \\ \emptyset & & & & & \\ & & & & & & 100 \end{pmatrix} : (\text{double eigenvalue for } \alpha \rightarrow 0).$$

Problem P7:

$$B = \begin{pmatrix} 1 & \alpha & & & \\ -\alpha & 1 & & & \emptyset \\ & & 3 & & \\ & & & 4 & \\ & & & & \ddots \\ \emptyset & & & & & \\ & & & & & & 100 \end{pmatrix} : (\text{conjugate eigenpair } 1 \pm \alpha i).$$





equal to  $1/30$  (EN is applied to the system multiplied by  $450/(\gamma/60 + 1)$ ). Starting with  $x_0 = (0, \dots, 0)^T$  gives the following results:

method	EN	GMRES
$\gamma$		
0	38	65
30	44	84
60	35	70
300	86	150
3000	408	455

Table 1. Number of iteration steps, for which  $\|r_i\|_2 / \|b\|_2 \leq 10^{-12}$

Except for the choice  $\gamma = 3000$ , it appears from Table 1 that roughly  $2k$  steps of GMRES are comparable with  $k$  steps of EN (see also the Figures 9 and 10).

### 3 A comparison of EN and GMRES

In this section we will show that the space spanned by the vectors  $c_k$ , generated by EN, is contained in a Krylov subspace. Furthermore, we will compare the norms of the residuals in EN and GMRES. Then by estimating the required amount of work and memory we will be able to compare the efficiency of both methods.

First we will show that the vectors  $c_k$  which are generated by the EN-method are elements of a Krylov subspace.

#### Theorem 3.1

If  $H_k$  is not singular and  $E_k r_k \neq 0$  then:

$$r_k = r_0 + \sum_{i=1}^{2k} \alpha_{ki} (AH_0)^i r_0 \text{ and } \text{span} \{c_0, \dots, c_k\} \subset \text{span} \{(AH_0) r_0, \dots, (AH_0)^{2k+2} r_0\}.$$

#### Proof

In order to simplify relations, we redefine  $c_k$ , in this proof, as:

$$c_k = Au_k, \tag{6}$$

(note that only the direction of  $c_k$  is relevant).

We prove the theorem by an induction argument in  $k$ . From (6) it follows that  $c_0 = AH_0 E_0 r_0 = AH_0 (I - AH_0) r_0$ , so that  $c_0 \in \text{span} \{(AH_0) r_0, (AH_0)^2 r_0\}$ . This implies the theorem to be true for  $k = 0$ .

Combination of (4) and (5) gives

$$r_{k+1} = E_{k+1} r_k = (I - P_k)(I - AH_0) r_k = (I - AH_0) r_k - P_k E_0 r_k.$$

Since  $P_k$  is the orthogonal projection onto  $\text{span} \{c_0, \dots, c_k\}$  it follows by induction that

$$r_{k+1} = r_0 + \sum_{i=1}^{2(k+1)} \alpha_{k+1,i} (AH_0)^i r_0. \tag{7}$$

Furthermore, from (6) we obtain  $c_{k+1} = AH_{k+1}E_{k+1}r_{k+1} = (I - E_{k+1})E_{k+1}r_{k+1}$  .  
 Together with (4) this gives:

$$c_{k+1} = (I - (I - P_k)(I - AH_0))E_{k+1}r_{k+1} = (AH_0 + P_kE_0)E_{k+1}r_{k+1} .$$

Another application of (4) leads to:

$$c_{k+1} = P_kE_0E_{k+1}r_{k+1} + AH_0(I - P_k)(I - AH_0)r_{k+1}$$

and hence

$$c_{k+1} = P_kE_0E_{k+1}r_{k+1} - AH_0P_kE_0r_{k+1} + AH_0(I - AH_0)r_{k+1} .$$

Since  $P_k$  is the orthogonal projection onto  $\text{span}\{c_0, \dots, c_k\}$  it follows by induction and (7) that  $c_{k+1} \in \text{span}\{(AH_0)r_0, \dots, (AH_0)^{2(k+1)+2}r_0\}$ , which completes the proof.  $\square$

The following definition is used for the comparison of the residuals of EN and GMRES.

Definition 3.2

$r_k^{EN}$  is the residual in the  $k$ -th step of EN.  $r_k^G$  is the residual in the  $k$ -th step of GMRES applied to the postconditioned linear system  $AH_0y = b$  where  $H_0$  is the same matrix in both methods (note that  $x = H_0y$  solves the system  $Ax = b$ ).

From Theorem 3.1 and (2) we obtain the following inequality

$$\| r_k^{EN} \|_2 \geq \| r_{2k}^G \|_2 . \tag{8}$$

This inequality supports our earlier observation made in the numerical experiments, reported in Section 2.

In order to compare the efficiency of EN and GMRES we need an estimate for the amount of work and memory in each method. For obvious reasons we have listed in Table 2 the amount of work and memory requirements for  $k$  steps of EN and  $2k$  steps of GMRES.

method	steps	multiplications with		inner products	vectorupdates	memory
		$H_0$	$A$			
EN1	$k$	$2k$	$4k$	$k^2$	$k^2$	$2kn$
EN2	$k$	$2k$	$2k$	$k^2$	$2k^2$	$2kn$
GMRES	$2k$	$2k$	$2k$	$2k^2$	$2k^2$	$2kn (+2k^2)$

Table 2. Amount of work and memory for different methods.

The inner products in EN1 can be computed in parallel. Furthermore in EN2 the vectorupdates, used to form  $\eta$  and  $\xi$  (or  $u_k$  and  $c_k$ ), can be computed in parallel. The inner products and vectorupdates in the implementation of GMRES as given in [7] can not be computed in parallel. This might be a disadvantage for GMRES in a parallel computing environment. Since in most of our numerical experiments  $\| r_k^{EN} \|_2$  and  $\| r_{2k}^G \|_2$  differ considerably, we also give estimates for the amount of work and memory requirements for the following experiment. The solution of Problem P10 with  $\gamma = 300$  is computed with the EN-method and the GMRES-method. The results are plotted in Figure 11. Note that EN requires more multiplications with

$H_0$  and  $A$  than GMRES to obtain the same accuracy. Choosing  $eps = 10^{-12}$  it appears that  $\| r_{86}^{EN} \|_2 / \| b \|_2 \leq eps$  and  $\| r_{150}^G \|_2 / \| b \|_2 \leq eps$ . The amount of work and memory requirements to obtain this accuracy are listed in Table 3.

method	steps	multiplications with		inner products	vectorupdates	memory
		$H_0$	$A$			
EN1	86	172	344	7396	7396	154800
EN2	86	172	172	7396	14792	154800
GMRES	150	150	150	11250	11250	135000(+11250)

Table 3. Amount of work and memory for different methods.

In practical situations the order of the linear system  $n$  will be much larger than the required number of iterations. In such cases the term  $2k^2$  in the required amount of memory for the GMRES-method is relatively negligible.

We conclude that when

$$\| r_k^{EN} \|_2 \approx \| r_{2k}^G \|_2$$

then the EN2-method is more efficient than the GMRES-method in terms of flops-counts. However the given experiment has shown that there are problems for which

$$\| r_k^{EN} \|_2 \approx \| r_j^G \|_2$$

with  $j < 2k$ .

In the following section we will give more evidence for such situations. In such cases it is less clear which method is preferable in terms of flops-counts. With respect to the memory requirements we note that GMRES is preferable.

## 4 Some specific properties of EN

In this section we will show that the convergence and stability properties of the EN-method are not scaling invariant. Subsequently we will provide some examples where the EN-method does not converge. Finally we will show that Property 1.6 is useless from a practical point of view.

### 4.1 The convergence behaviour of EN with respect to scaling

From its construction it follows that GMRES is scaling invariant, which means that when the method is applied to the system  $\rho Ax = \rho b$  then the iterates are the same for every choice of  $\rho \neq 0$ . One might expect from the foregoing that EN has the same property. However, from our experiments it follows that EN is not scaling invariant. This is well illustrated by the results for Problem P6 (with  $\alpha = 10^{-5}$  and  $\beta = 0.9$ ). In our first experiment we take  $H_0 = \rho I$  as an approximation for  $A^{-1}$ . Obvious choices for  $\rho$  are  $\rho = \frac{1}{\lambda_1}$ ,  $\rho = \frac{2}{\lambda_1 + \lambda_n}$ , and  $\rho = \frac{1}{\lambda_n}$ , where  $\lambda_1 = 1$  is the smallest and  $\lambda_n = 100$  the largest eigenvalue of  $A$ . We obtain  $\| r_{100}^{EN} \|_2 / \| r_0 \|_2 = 10^{65}$  for  $\rho = \frac{1}{\lambda_1} = 1$ ,  $\| r_{38}^{EN} \|_2 / \| r_0 \|_2 \leq 10^{-12}$  for  $\rho = \frac{2}{\lambda_1 + \lambda_2} = \frac{2}{101}$ , and  $\| r_{40}^{EN} \|_2 / \| r_0 \|_2 \leq 10^{-12}$  for  $\rho = \frac{1}{\lambda_n} = \frac{1}{100}$ .

So the convergence behaviour of EN strongly depends on the choice of  $\rho$ .

As a second experiment we apply EN to Problem P6 with  $B := \rho B$  for  $\rho = 10^{-1}, 10^{-2}, 10^{-3}$  and  $10^{-4}$  and  $H_0 = I$ . The method is terminated as soon as  $\|r_i\|_2 / \|b\|_2 \leq 10^{-12}$ . The number of iteration steps, for different choices of  $\rho$ , is given in Table 4.

$\rho$	$10^{-1}$	$10^{-2}$	$10^{-3}$	$10^{-4}$
iterates	78	40	64	66

Table 4. Number of iteration steps, for which  $\|r_i\|_2 / \|b\|_2 \leq 10^{-12}$  (for P6).

The convergence behaviour is displayed in Figure 12. In this figure, each curve is plotted at the right scale. For  $\rho = 10^{-1}$  we notice that initially the residuals increase. For  $\rho = 10^{-2}$  the curve is identical to the corresponding curve in Figure 1. Note that the curves for  $\rho = 10^{-3}$  and  $\rho = 10^{-4}$  are nearly the same. Furthermore, these curves show a striking resemblance with the corresponding curve for GMRES in Figure 2.

A possible explanation for this might come from the observation that for  $\rho = 10^{-4}$  we have that  $E_0 = I - AH_0 \approx I$ . This together with (4) implies

$$E_k \approx (I - P_{k-1}).$$

Using this expression and (5) it follows from

$$u_k = H_k E_k r_k = H_k E_k^2 r_{k-1} \approx H_k E_k r_{k-1} = H_k r_k,$$

that  $u_k \approx H_k r_k$ . This explains the resemblance of the curves, since the choice  $u_k = H_k r_k$  leads to a method algebraically equivalent to GMRES (see [3]: p.513 and also the following section).

In our example the choice  $\rho = 10^{-2}$  is obviously preferable. We will call this value  $\rho_{opt}$  for our experiment. However, in general we know of no criterium which could be used for defining a priori an optimal  $\rho$ . Hence  $\rho_{opt}$  has to be determined experimentally. Furthermore, for this example we observe for  $\rho = 10\rho_{opt}$  the speed of convergence is halved, whereas for  $\rho = 0.1\rho_{opt}$  the speed of convergence is approximately the same as for GMRES. Taking into account the amount of work and memory for both methods (see Table 2, Section 3) we conclude that we need a fairly good guess for  $\rho_{opt}$  if we want EN be more efficient than GMRES.

From these experiments it seems attractive that the spectral radius of  $(I - AH_0)$  has to be less than one (compare Section 4.2). This conjecture is confirmed by the following experiment. We take  $\Omega$  to be the unit square and consider the pde

$$\Delta u = 0 \text{ on } \Omega \text{ and } u|_{\partial\Omega} \text{ is given.}$$

Using the standard five point central finite difference approximation over an equidistant rectangular grid we obtain a symmetric linear system. For  $H_0$  we take an average of the incomplete Choleski (IC) and a modified incomplete Choleski matrix (MIC) see ([8]: Section 3). The IC matrix corresponds with  $\theta = 0$ , whereas the MIC matrix corresponds with  $\theta = 1$ . Taking 200 points in  $x$ - and  $y$ -direction, and  $x_0 = 0$  we obtain the results as given in Table 5.

$\theta$	0	0.5	0.9	0.95	0.96	0.97	0.98	0.99	1
iterates	21	17	14	13	14	13	17	46	*

Table 5. Number of iteration steps, for which  $\|r_i\|_2 / \|r_0\|_2 \leq 10^{-6}$

Note that EN converges rather fast for the choices  $0 \leq \theta \leq 0.98$  but diverges for the choice  $\theta = 1$  which corresponds with the MIC preconditioner. This seems to be quite in line with similar experiments reported for preconditioned cg in ([8]: Section 3). However, if we apply EN to  $0.1 * AH_0$  and  $\theta = 1$  then we obtain  $\|r_{23}^{EN}\|_2 / \|r_0\|_2 \leq 10^{-6}$ . Therefore we believe that these experiments confirm our conjecture, since the spectral radius of  $(I - AH_0)$  with the IC matrix is less than one, whereas with the MIC matrix the spectral radius is much larger than one (see [4]). This result suggests that the divergence for  $\theta = 1$  in the previous experiment is not caused only by a loss of independence among the Krylov subspace basis vectors for this value of  $\theta$  (which is the reason for slow convergence of cg in this case ([8]: Section 3)). We conclude that the convergence behaviour of EN depends not only on the choice of  $H_0$  but also on the scaling parameter  $\rho_{opt}$ . We expect good convergence if the spectral radius of  $(I - \rho_{opt}AH_0)$  is less than one.

Our experiments show that EN is not invariant with respect to a general transformation of coordinates. Note that this conclusion is not in contradiction with ([3]: Proposition 2.2), which states that EN is invariant under unitary transformations.

## 4.2 The stability of EN with respect to scaling

From Figure 12 it appears that initially the residuals increase for  $\rho = 10^{-1}$ . To illustrate this phenomenon we will describe some experiments for  $\rho$  in the vicinity of 0.1. The results are given in Table 6 where  $i$  is the smallest value such that  $\|r_i\|_2 / \|b\|_2 \leq 10^{-12}$  and  $imax$  is defined by  $\|r_{imax}\|_2 = \max_{1 \leq j \leq i} \|r_j\|_2$ .

$\rho$	0.09	0.10	0.11	0.13	0.15
$imax$	1	32	42	53	59
$\ r_{imax}\ _2 / \ b\ _2$	1	$1.9 \times 10^1$	$1.8 \times 10^3$	$4.5 \times 10^7$	$9.4 \times 10^{11}$
$i$	74	78	80	84	91
$\ r_i\ _2 / \ b\ _2$	$9 \times 10^{-13}$	$2.6 \times 10^{-13}$	$4.3 \times 10^{-13}$	$3.3 \times 10^{-13}$	$4.4 \times 10^{-13}$
$\ b - Ax_i\ _2 / \ b\ _2$	$9 \times 10^{-13}$	$4.3 \times 10^{-13}$	$2.9 \times 10^{-11}$	$1.4 \times 10^{-6}$	$2.2 \times 10^{-2}$

Table 6.  $\|r_{imax}\|_2$  for different values of  $\rho$ .

This table shows that initial residuals increase fast for  $\rho \geq 0.1$  and that the inequality  $\|b - Ax_i\|_2 / \|b\|_2 \leq 10^{-12}$  does not hold for  $\rho \geq 0.10$ , as it should in exact arithmetic. For a possible explanation of the increase of the residuals we make use of the equality  $r_{k+1} = (I - P_k)E_0 r_k$ . The right-hand side consists of two parts: firstly a multiplication with  $E_0$  and secondly a multiplication with the orthogonal projection  $(I - P_k)$ . Since  $\sigma(E_0) \subset [1 - 100\rho, 1 - \rho]$  it follows that when  $\rho > 2 \times 10^{-2}$ ,  $\|E_0 r_k\|_2$  can be larger than  $\|r_k\|_2$ . For the second part

we always have  $\| (I - P_k) E_0 r_k \|_2 \leq \| E_0 r_k \|_2$ . From this it appears that for  $\rho \in (0, 0.02)$  the residual decreases in both parts. For  $\rho \in [0.02, 0.09]$  the increase in the first part is cancelled by the decrease in the second part. For  $\rho \in (0.09, \infty)$  initially the increase in the first part dominates whereas after a number of iteration steps ( $imax$ ) the decrease in the second part dominates.

Note that in exact arithmetic  $r_i = b - Ax_i$ . However, for  $\rho \geq 0.1$  this is clearly violated in EN and hence the reliability of  $r_i$  given by EN depends on the value  $\rho$ . To explain this we assume that  $r_i$  and  $x_i$  denote the exact values and  $\hat{r}_i$  and  $\hat{x}_i$  denote the numerically computed values. Now define  $z_i = r_{imax} - r_i$ ,  $\hat{z}_i = \hat{r}_{imax} - \hat{r}_i$ , and suppose that  $\| \hat{r}_{imax} - r_{imax} \|_2 / \| r_{imax} \|_2 \approx \epsilon$  and  $\| \hat{z}_i - z_i \|_2 / \| z_i \|_2 \approx \epsilon$ , where  $\epsilon$  is a modest multiple of the machine precision. For  $\rho = 0.15$  this implies  $\| \hat{r}_i - r_i \|_2 = \| \hat{r}_{imax} - r_{imax} - (\hat{z}_i - z_i) \|_2 \approx \epsilon \{ \| r_{imax} \|_2 + \| z \|_2 \} \approx 2 \times 10^{12} \| b \|_2 \epsilon$  and  $\| \hat{r}_i - (b - A\hat{x}_i) \|_2 = \| \hat{r}_i - r_i + (b - Ax_i) - (b - A\hat{x}_i) \|_2 \approx 2 \times 10^{12} \| b \|_2 \epsilon$ . This implies that due to rounding errors it is possible for  $\rho = 0.15$  that

$$\| r_i \|_2 / \| b \|_2 \leq 10^{-12}$$

whereas

$$\| b - A\hat{x}_i \|_2 / \| b \|_2 \approx 10^{12} \epsilon$$

(note that  $\kappa_2(A) = 100$ ).

We conclude that the stability of the EN-method depends on  $\rho$ . In the given experiment the EN-method is quite stable for  $\rho \leq 0.09$  and rather unstable for  $\rho \geq 0.1$ . It is, in general, not known for which  $\rho$  EN is stable. These results do not support the stability properties claimed in ([3]: p.516).

### 4.3 Some examples where the EN-method does not converge

In this subsection we give some examples for which EN fails to converge. In order to identify such problems we look for nonsingular matrices  $A$  and  $H_0$  such that  $H_1$  is singular. Taking  $H_0 = I$  it follows from (3) that if  $E_0 r_0 \neq 0$  then  $c_0 = \gamma A E_0 r_0$  with  $\gamma = 1 / \| A E_0 r_0 \|_2$ . Using Property 1.6 it follows that  $H_1$  is singular if and only if  $c_0^T E_0 r_0 = \gamma (A E_0 r_0)^T E_0 r_0 = 0$ . Thus  $A$  should be such that  $(Av)^T v = 0$  for  $v \in R^n$  which means that  $Av$  and  $v$  are orthogonal. A simple matrix with this property is  $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ .

#### Example 1

We apply EN to  $Ax = b$  with  $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ ,  $H_0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ ,  $x = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$  and  $b = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ . Starting with  $x_0 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  gives  $r_0 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ . Since  $E_0 = I - AH_0 = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$  we obtain  $E_0 r_0 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$  and  $c_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ , which implies that  $c_0^T E_0 r_0 = 0$  ( $H_1$  singular). Continuing the method with  $H_1 = H_0$  yields  $E_1 = E_0$  and  $c_1^T E_0 r_1 = 0$  ( $H_2$  singular). After  $k$  iteration steps

we obtain  $H_k = H_0$ ,  $E_k = E_0$  and  $r_k = E_0^k r_0$ . The eigenvalues of  $E_0$  are  $1 + i$  and  $1 - i$  so that

$$r_k = P \begin{bmatrix} (1+i)^k & 0 \\ 0 & (1-i)^k \end{bmatrix} P^{-1} r_0 \text{ and } \|r_k\|_2 \rightarrow \infty \text{ for } k \rightarrow \infty.$$

Thus, for this example the EN-method is clearly divergent.

This example shows that EN does not converge for each given linear system. It is known that GMRES converges slowly for this type of matrices. In ([5]: p.23) it is shown that when GMRES is applied to  $Ax = b$  with  $A \in R^{n \times n}$  given by

$$A = \begin{pmatrix} 0 & & & 0 & 1 \\ 1 & 0 & & & 0 \\ 0 & 1 & 0 & & 0 \\ & & & \ddots & \\ \emptyset & & & & 1 & 0 \end{pmatrix}, \quad b = (1, 0, \dots, 0)^T \text{ and } x_0 = (0, \dots, 0)^T,$$

then  $x_i = x_0$ ,  $0 \leq i \leq n - 1$  and  $x_n = x$ .

In our following example EN converges slowly, whereas GMRES converges very fast.

### Example 2

Take

$$H_0 = I, \quad A = \rho \begin{pmatrix} 0 & -10^4 & & & \\ 10^4 & 0 & & & \emptyset \\ & & 1.03 & & \\ & & & 1.04 & \\ & & \emptyset & & \ddots \\ & & & & & 2 \end{pmatrix} \text{ and } x = (1, \dots, 1)^T.$$

Starting with  $x_0 = (0, \dots, 0)^T$  we obtain  $\|r_{100}^{EN1}\|_2 / \|b\|_2 \geq 10^{-7}$  for  $\rho = 10^{-3}$ , whereas  $\|r_{15}^G\|_2 / \|b\|_2 \leq 10^{-12}$ . The rather bizarre convergence behaviour of EN in dependence on the scaling parameter  $\rho$  is nicely illustrated by the fact that  $\|r_{14}^{EN1}\|_2 / \|b\|_2 \leq 10^{-12}$  for  $\rho = 10^{-4}$  but  $\|r_{100}^{EN1}\|_2 / \|b\|_2 \geq 10^{-5}$  for  $\rho = 10^{-5}$ . Using the EN2 implementation we obtain for the updated residual  $\|r_{14}^{EN2}\|_2 / \|b\|_2 \leq 10^{-12}$  for  $\rho = 10^{-3}, 10^{-4}$  and  $10^{-5}$  whereas the exact residual  $\|Ax_{14} - b\|_2 / \|b\|_2$  equals  $3 \times 10^{-4}$ ,  $2 \times 10^{-6}$  and  $5 \times 10^{-7}$  for  $\rho$  respectively  $10^{-3}, 10^{-4}$  and  $10^{-5}$ .

In Section 4.1 we have seen, and explained, that for  $\rho$  small enough application of EN to  $\rho Ax = \rho b$  gives  $\|r_i^{EN}\|_2 \approx \|r_i^G\|_2$  for some problems. Example 2 shows that there are also linear systems where this equivalence does not hold.

## 4.4 The practical relevance of Property 1.6

In this subsection we consider the application of EN1 to Example 2 for  $\rho = 10^{-3}$ . Taking into account the similarity between Examples 1 and 2 we expect that in Example 2,  $H_1$  is nearly singular. By computation it follows that  $\|E_1 r_1\|_2 = 1.4 \times 10^2$  and  $\|H_1 E_1 r_1\|_2 = 2.6 \times 10^{-10}$  so  $\|H_1^{-1}\|_2$  is very large. It appears that the computed vector  $c_1 = \gamma A H_1 E_1 r_1$  has a large relative

error and that  $c_0^T c_1$  equals  $1.5 \times 10^{-4}$  instead of 0. This explains the bad convergence behaviour EN1 in Example 2.

The large relative error in  $c_1$  is also predicted by ([3]: Theorem 2.2) if we use that

$$c_H = \| E_1 \|_2 / \| AH_1 E_1 r_1 \|_2 \geq 3 \times 10^8.$$

This experience motivates us to investigate the practical applicability of Property 1.6. We note the following drawbacks:

- if  $c_k^T E_0 r_k = 0$  then it is possible of course that the computed value of  $c_k^T E_0 r_k \neq 0$ ,
- if  $c_k^T E_0 r_k \neq 0$  then it is still possible that  $H_{k+1}$  is nearly singular.

To get around these difficulties we could replace condition  $c_k^T E_0 r_k = 0$  by

$$| c_k^T E_0 r_k | / \| E_0 r_k \|_2 \leq \epsilon \text{ for } \epsilon \geq 0. \quad (9)$$

If inequality (9) holds we take  $H_{k+1} = H_k$ . However this condition has certain disadvantages too. First of all it is not clear which value of  $\epsilon$  is feasible. Secondly implementation of this condition does not help much in Example 2. In this case we have  $| c_0^T E_0 r_0 | / \| E_0 r_0 \|_2 = 1.8 \times 10^{-12}$ . If we take  $\epsilon \leq 1.8 \times 10^{-12}$  then we obtain the same results as without this condition, whereas  $\epsilon > 1.8 \times 10^{-12}$  leads to  $H_k = H_0$  for  $0 \leq k \leq 100$  and  $\| r_{100} \|_2 / \| b \|_2 = 10^{100}$ . Hence, for Example 2 there is no value of  $\epsilon$  such that the EN1-method combined with (9) is convergent.

This indicates that implementing Property 1.6 in this way is useless from a practical point of view.

From the given examples it follows that EN is not attractive if some of the matrices  $H_k$  are (nearly) singular. Therefore, it is important to know a priori when the matrices  $H_k$  are (nearly) singular. In ([3] p.516, Theorem 2.3) the following "safe" case is stated: if  $AH_0$  is positive (negative) definite then  $\| (AH_k)^{-1} \|_2 \leq 1/\mu$  where

$$\mu = \inf_{\| x \|_2 = 1} | (AH_0 x)^T x | / \left( \| x \|_2^2 + \| AH_0 x \|_2^2 \right)^{1/2}.$$

The following theorem states that if  $AH_0$  is neither positive nor negative definite then it is possible to obtain a singular matrix  $H_k$ .

#### Theorem 4.2

If  $AH_0$  is neither positive nor negative definite on  $Im(E_0)$  then there exists a right-hand side vector  $b$  such that  $H_1$  is singular.

#### Proof

The condition on  $AH_0$  implies that there is a vector  $v \in Im(E_0)$  such that  $(AH_0 v)^T v = 0$ . Since  $v \in Im(E_0)$  we can find  $b \in R^n$  such that  $E_0 b = v$ . Applying EN to this system with  $x_0 = (0, \dots, 0)^T$  yields  $c_0 = \gamma AH_0 E_0 b$  with  $\gamma = 1 / \| AH_0 E_0 b \|_2$ . From Property 1.6 and the equations

$$c_0^T E_0 r_0 = \gamma (AH_0 E_0 b)^T E_0 b = \gamma (AH_0 v)^T v = 0$$

it follows that  $H_1$  is singular.

Note that if there is a vector  $v$  such that  $(Av)^T v = 0$  then  $1/\mu$  is infinite.  $\square$

Our conclusion is that it is only "safe" to apply the EN-method if  $AH_0$  is positive or negative definite.



## 4.5 A scaling invariant version of the EN method

In Section 4.1 and 4.2 we have shown that the convergence and stability properties of EN are not scaling invariant. As a consequence of this one should estimate a parameter  $\rho_{opt}$  such that the spectral radius of  $(I - \rho_{opt}AH_0)$  is less than one. In this section we modify the EN method such that parameter estimation is not longer required.

In Section 1 we have shown that  $r_{k+1} = (I - AH_k)r_k - \mu_k Au_k$ . Combination with  $u_k = H_k(I - AH_k)r_k$  gives  $r_{k+1} = (I - \mu_k AH_k)(I - AH_k)r_k$ .

Since  $E_k = (I - AH_k) = (I - P_{k-1})(I - AH_0)$ ,  $r_{k+1}$  can also be written as  $r_{k+1} = (I - \mu_k AH_k)(I - P_{k-1})(I - AH_0)r_k$ .

Note that it is the multiplication with  $(I - AH_0)$  which makes EN not scaling invariant. Using this observation we modify EN such that  $r_{k+1}$  obtained with the modified EN method can be written as follows:

$$r_{k+1} = (I - \mu_k AH_k)(I - P_{k-1})(I - \gamma_k AH_0)r_k$$

where the constant  $\gamma_k = (AH_0 r_k)^T r_k / \|AH_0 r_k\|_2^2$  minimizes  $\|(I - \gamma AH_0)r_k\|_2$ . To implement the modified method (EN3) the first step of the implementation EN2 should be changed as follows:

$$\gamma = (AH_0 r_k)^T r_k / \|AH_0 r_k\|_2^2, \quad \xi^{(0)} = (I - \gamma AH_0)r_k, \quad \eta^{(0)} = \gamma H_0 r_k, \quad \alpha_i = c_i^T \xi^{(i)}, \quad \xi^{(i+1)} = \xi^{(i)} - \alpha_i c_i, \quad \eta^{(i+1)} = \eta^{(i)} + \alpha_i u_i, \quad i = 0, \dots, k-1.$$

It is easy to show that EN3 is scaling invariant, which is confirmed by our numerical experiments. Application of EN3 to Problem P6 (with  $\alpha = 10^{-5}$  and  $\beta = 0.9$ ) with  $B := \rho B$  gives  $\|r_{35}^{EN}\|_2 / \|r_0\|_2 \leq 10^{-12}$  for all choices of  $\rho$ . Finally we apply EN3 to the pde problem given in Section 4.1. The results are given in Table 7.

$\theta$	0	0.5	0.9	0.95	0.96	0.97	0.98	0.99	1
iterates	21	18	12	11	11	10	10	9	22

Table 7. Number of iteration steps, for which  $\|r_i^{EN3}\|_2 / \|r_0\|_2 \leq 10^{-6}$ .

Note that EN3 converges also for the choice  $\theta = 1$ . Furthermore the optimal number of iterates of EN2 in Table 5 equals 13 whereas the optimal number of iterates of EN3 in Table 7 equals 9. Thus in this example we observe that the convergence of EN3 is approximately 1.5 times as good as the convergence of EN2.

## 5 Another formulation of the GMRES-method

In ([3]: p.513) it is noted without proof that, when choosing

$$u_k = H_k r_k, \tag{10}$$

instead of  $u_k = H_k E_k r_k$ , the EN-method leads to an algorithm algebraically equivalent to GMRES. In this section we first prove this equivalence under the assumption that the matrices  $H_k$  are nonsingular. Subsequently we give a slight modification of the choice (10) such that the method remains equivalent to GMRES even if the matrices  $H_k$  are singular. A suitable

implementation of this method arises if an orthonormal basis for the Krylov subspace is generated by the modified Gram Schmidt process.

First we will show that the vectors  $c_k$  form an orthonormal basis for

$$\text{span}\{(AH_0)r_0, \dots, (AH_0)^{k+1}r_0\}.$$

Since (4) is only valid for the choice  $u_k = H_k E_k r_k$  we use the equality

$$E_{k+1} = \left(I - c_k c_k^T\right) E_k \quad (11)$$

(cf [3]: p.512) .

#### Theorem 5.1

Let  $u_k$  be chosen as  $u_k = H_k r_k$  in EN . When  $H_k$  is not singular and  $r_k \neq 0$  , then

$$r_{k+1} = (I - P_k) r_0 \text{ where } P_k = \sum_{i=0}^k c_i c_i^T \text{ is the orthogonal projection onto } \text{span}\{(AH_0)r_0, \dots, (AH_0)^{k+1}r_0\}.$$

#### Proof

Similar as in the proof for Theorem 3.1 we take

$$c_k = A u_k \quad (12)$$

We prove the theorem by an induction argument in  $k$ . Using (10) and (12) we obtain  $c_0 = (AH_0)r_0$ . Combination of (5) and (11) gives  $r_1 = E_1 r_0 = \left(I - c_0 c_0^T\right) E_0 r_0$ . Since  $c_0 = (AH_0)r_0$  it follows that  $r_1 = \left(I - c_0 c_0^T\right) r_0$ . This implies the theorem to be true for  $k = 0$ .

It follows from (10) and (12) that  $c_{k+1} = AH_{k+1} r_{k+1} = (I - E_{k+1}) r_{k+1}$ .

Equation (11) implies  $E_{k+1} = (I - P_k) E_0$  and  $c_{k+1} = (I - (I - P_k)(I - AH_0))(I - P_k) r_0$  by induction. The last equation can also be written as

$$c_{k+1} = (I - P_k) AH_0 (I - P_k) r_0 = (I - P_k) \left( c_0 - \sum_{i=0}^k (AH_0) c_i c_i^T r_0 \right).$$

By induction it follows that  $(AH_0) c_i \in \text{span}\{c_0, \dots, c_k\}$ , for  $i = 0, \dots, k - 1$ , hence

$$c_{k+1} = -(I - P_k) AH_0 c_k c_k^T r_0 \quad (13)$$

Since  $c_k$  has a non-zero component in the direction of  $(AH_0)^{k+1} r_0$ ,  $H_{k+1}$  is nonsingular and  $r_{k+1} \neq 0$  it follows that  $c_{k+1}$  has a non-zero component in the direction of  $(AH_0)^{k+2} r_0$ . Using (13) it follows by induction that  $c_i^T c_{k+1} = 0, i = 0, \dots, k$ . Thus  $\{c_0, \dots, c_{k+1}\}$  is an orthonormal basis for  $\text{span}\{(AH_0)r_0, \dots, (AH_0)^{k+2}r_0\}$ . Combining (5), (11) and (12) we obtain

$$r_{k+2} = E_{k+2} r_{k+1} = \left(I - c_{k+1} c_{k+1}^T\right) (I - AH_{k+1}) r_{k+1} ,$$

so that

$$r_{k+2} = \left(I - c_{k+1} c_{k+1}^T\right) (r_{k+1} - c_{k+1}) = \left(I - c_{k+1} c_{k+1}^T\right) r_{k+1} .$$

By induction it follows that

$$r_{k+2} = \left( I - P_k - c_{k+1}c_{k+1}^T \right) r_0,$$

which concludes our proof.  $\square$

From this theorem we conclude that the method converges if the matrices  $H_k$  are nonsingular. Since GMRES leads to the solution within a finite number of iteration steps we look for a modification of (10) such that the condition on  $H_k$  can be dropped. To this end we note that it follows from (13) that  $c_{k+1}$  is a unit vector in the direction of  $(I - P_k)(AH_0)c_k$ . Choose the vector  $u_k$  as follows

$$\begin{cases} u_0 = H_0 r_0 \\ u_k = u_{k-1} - H_k c_{k-1}, k \geq 1 \end{cases} \quad \text{which implies} \quad \begin{cases} c_0 = AH_0 r_0 \\ c_k = -(I - P_{k-1}) AH_0 c_{k-1} \end{cases} \quad (14)$$

it can be proven that  $r_{k+1} = (I - P_k)r_0$  where  $P_k = \sum_{i=0}^k c_i c_i^T$  is the orthogonal projection onto  $\text{span} \{(AH_0)r_0, \dots, (AH_0)^{k+1}r_0\}$ . Furthermore, it is easy to show if  $c_k \neq 0$  and  $c_{k+1} = 0$  then  $r_{k+1} = 0$ .

From this remark and (2) it follows that EN with  $u_k$  as in (14) is equivalent to GMRES applied to the postconditioned linear system  $AH_0 y = b$ .

An implementation of this method is:

1.  $u_0 = H_0 r_0 / \| AH_0 r_0 \|_2, c_0 = Au_0, k = 0,$   
 $x_1 = x_0 + u_0 c_0^T r_0$  and  $r_1 = r_0 - c_0 c_0^T r_0,$
2. while  $\| r_{k+1} \|_2 > \text{eps}$  do  $k := k + 1,$   
 $c_k^{(0)} = AH_0 c_{k-1}, u_k^{(0)} = H_0 c_{k-1},$   
 $\alpha_i = c_i^T c_k^{(i)}, c_k^{(i+1)} = c_k^{(i)} - \alpha_i c_i, u_k^{(i+1)} = u_k^{(i)} - \alpha_i u_i, i = 0, \dots, k-1,$   
 $c_k = c_k^{(k)} / \| c_k^{(k)} \|_2, u_k = u_k^{(k)} / \| c_k^{(k)} \|_2,$
3.  $x_{k+1} = x_k + u_k c_k^T r_k$  and  $r_{k+1} = r_k - c_k c_k^T r_k.$

Note that the vectors  $c_k$  are made mutually orthogonal by the modified Gram-Schmidt process. In this implementation of GMRES,  $2k$  iteration steps involve  $2k$  multiplications with  $A$  and  $H_0$ ,  $2k^2$  inner products,  $4k^2$  vectorupdates and  $4kn$  memory space. Comparing this with GMRES in Table 2 it follows that this implementation requires  $2k^2$  vectorupdates and  $2kn$  memory space extra.

Using the choice (14) the GMRES-method is formulated in the same way as the EN-method. This correspondence gives some theoretical insight, for instance that also for GMRES a matrix  $H_k$  can be formed, which approximates the inverse of  $A$ . With respect to flop-counts and memory requirements we prefer the implementation of GMRES given in [6] and [7]

## 6 A comparison of the EN- and the B-methods

In this section we compare the EN-method with the B-method described in [1]. The B-method is mostly used to solve nonlinear systems but it can also be used to solve a linear system. The description of the B-method indicates certain similarities with the EN-method. However a further investigation reveals essential differences. The main difference is that  $r_k = r_0 + \sum_{i=1}^{2k} \alpha_{ki} (AH_0)^i r_0$  for the EN-method, whereas

$r_k = r_0 + \sum_{i=1}^k \beta_{ki} (AH_0)^i r_0$  for the B-method . We conclude this section with a comparison of the B- and the GMRES-methods.

From the descriptions of the EN- and B-methods, in Section 1, we note the following correspondence: in both methods rank-one updates are used to construct a matrix  $H_k$  which is an approximation to the inverse of  $A$  (compare [1] and [2]).

In order to make a more detailed comparison we use the following vectors:

### Definition 6.1

$$u_k = -(H_k A H_k - H_k) r_k, \quad v_k = H_k^T H_k r_k / r_k^T H_k^T H_k A H_k r_k, \quad c_k = A u_k .$$

Note that  $u_k, v_k$ , etc., are different for the different methods. Since only the residuals for the methods will be compared we have chosen to identify them by a superscript like  $r_k^B$  (for Broyden), where necessary.

From the description of the B-method, Definition 6.1 and the equations  $p_k + H_k r_k$  and  $y_k = r_k - r_{k+1} = r_k - b + A(x_k + H_k r_k) = A H_k r_k$  we deduce that

$$H_{k+1} = H_k + u_k v_k^T . \quad (15)$$

### Theorem 6.2

If  $r_k^T H_k^T H_k A H_k r_k \neq 0$  then  $r_k = r_0 + \sum_{i=1}^k \alpha_{ki} (AH_0)^i r_0$ , where  $\text{span} \{c_0, \dots, c_k\} \subset \text{span} \{(AH_0) r_0, \dots, (AH_0)^{k+2} r_0\}$  .

### Proof

We prove the theorem by an induction argument in  $k$  . From Definition 6.1 it follows that  $c_0 = A u_0 = -(AH_0)^2 r_0 + (AH_0) r_0$ , hence  $\text{span} \{c_0\} \subset \text{span} \{(AH_0) r_0, (AH_0)^2 r_0\}$ . This implies the theorem to be true for  $k = 0$ .

Since  $x_{k+1} = x_k + H_k r_k$  we have  $r_{k+1} = (I - AH_k) r_k$ . This together with Definition 6.1 and (15) yields

$$r_{k+1} = r_k - A \left( H_0 + \sum_{i=0}^{k-1} u_i v_i^T \right) r_k = r_k - A H_0 r_k - \sum_{i=0}^{k-1} c_i v_i^T r_k .$$

Now, it follows by induction that

$$r_{k+1} = r_0 + \sum_{i=1}^{k+1} \alpha_{k+1,i} (AH_0)^i r_0 . \quad (16)$$

By Definition 6.1 we have that  $c_{k+1} = -(AH_{k+1})^2 r_{k+1} + (AH_{k+1}) r_{k+1}$ . Since  $H_{k+1} = H_0 + \sum_{i=0}^k u_i v_i^T$  the following equation holds

$$c_{k+1} = - \left( AH_0 + \sum_{i=0}^k c_i v_i^T \right)^2 r_{k+1} + \left( AH_0 + \sum_{i=0}^k c_i v_i^T \right) r_{k+1} .$$

This implies that

$$c_{k+1} = -(AH_0)^2 r_{k+1} + (AH_0) r_{k+1} + \sum_{i=0}^k c_i v_i^T \{ -AH_0 - \sum_{i=0}^k c_i v_i^T + 1 \} r_{k+1} - \sum_{i=0}^k (AH_0) c_i v_i^T r_{k+1} .$$

Using (16) it follows by induction that  $c_{k+1} \in \text{span} \{ (AH_0) r_0, \dots, (AH_0)^{k+3} r_0 \}$ .  $\square$

Theorem 6.2 together with (2) yields the following inequality  $\| r_k^B \|_2 \geq \| r_k^G \|_2$ . From the numerical experiments given in Section 2 it follows that  $\| r_k^G \|_2$  can be much larger than  $\| r_k^{EN} \|_2$ . Hence, the B- and EN-methods can not be equivalent.

In [2] a generalization of the B-method is given. In this method, the BG-method, the update of  $H_k$  is as follows:

$$H_{k+1} = H_k - (H_k y_k - p_k) q_k^T / q_k^T y_k ,$$

where  $p_k = H_k r_k$ ,  $y_k = AH_k r_k$  and  $q_k$  is an arbitrary vector subject only to the restriction that  $q_k^T y_k \neq 0$ . Note that with the choice  $q_k = H_k^T p_k$  the BG-method is equal to the B-method. In the same way as for the B-method it follows that  $\| r_k^{BG} \|_2 \geq \| r_k^G \|_2$ , hence there is no choice of  $q_k$  such that BG and EN are equivalent.

It can be shown that for  $q_k = E_k^T A u_k$ , BG is algebraically equivalent to GMRES, which starts with  $x_0 + H_0(b - Ax_0)$ . Furthermore, it appears that BG with  $q_k = E_k^T A u_k$  is a secant method [2]. However, the specific method given in ([2]: p.373), is different from GMRES.

In order to estimate the efficiency of the BG-method we make a comparison with the GMRES-method. With respect to the amount of work and memory for an implementation of the BG-method we note that the k-th step costs at least one multiplication with  $H_0$  and  $A$  together with  $k$  inner products and  $2k$  vector updates, whereas  $2k$  vectors of length  $n$  should be stored in memory. This, in combination with the inequality  $\| r_k^{BG} \|_2 \geq \| r_k^G \|_2$ , yields that for every choice of  $q_k$  the BG-method is less efficient than GMRES applied to the postconditioned system  $AH_0 y = b$ .

## 7 Conclusions

In this paper we have compared the methods GMRES, EN and B. From this comparison it appears that in some numerical experiments EN takes less work than GMRES. However, a theoretical investigation shows that the efficiency of EN can be at most only slightly better than that of GMRES. Furthermore, the numerical experiments show that the convergence and stability of EN are not scaling invariant. However we specify a new version of the EN-method, which is scaling invariant. The convergence behaviour of this version seems to be better than that of the original EN-method. Subsequently we gave a formulation of GMRES in the same spirit

as the EN-method. This correspondence gives theoretical insight, but in practical situations we prefer the implementation of the GMRES-method as given in [6] and [7]. Since the class of BG-methods, proposed in [2], seems to be related to EN this class is included in our comparison. We show that the EN-method (with  $uk = H_k E_k r_k$ ) is not equivalent to any BG-method. With respect to GMRES, a BG-method is specified, which is algebraically equivalent to GMRES.

## References

- [1] Broyden, G.C. *A new method of solving nonlinear simultaneous equations*. The Computer J. 12 pp.94-99, 1969.
- [2] Broyden, G.C. *The convergence of single-rank quasi-Newton methods*. Math. Comp. 24 pp.365-382, 1970.
- [3] Eirola, T. Nevanlinna, O. *Accelerating with rank-one updates*. L.A.A. 121 pp.511-520, 1989.
- [4] Gustafsson, I. *A class of first order factorization methods*. BIT 18 pp.142-156, 1978.
- [5] Huang, Y. Van der Vorst, H.A. *Some observations on the convergence behavior of GMRES*. Delft University of Technology, Report 89-09, 1989.
- [6] Saad, Y. Schultz, M.H. *GMRES: a generalized minimal residual algorithm for solving non symmetric linear systems*. SIAM. J. Sci. Stat. Comput. 7 pp.856-869, 1986.
- [7] Van der Vorst, H.A. *The convergence behavior of some iterative solution methods*. Proceedings of the fifth International Symposium on Numerical Methods in Engineering 1 pp.61-72. Editors: Gruber, R. Periaux, J. Shaw, R.P. Springer Verlag , Berlin, 1989.
- [8] Van der Vorst, H.A. *The convergence behaviour of preconditioned CG and CG-S in the presence of rounding errors*. in: O.Axelsson and L.Yu.Kolotilina (Eds.), *Preconditioned Conjugate Gradient Methods*, Proceedings, Nijmegen, 1989, Lecture Notes in Mathematics 1457, Springer Verlag, Berlin, 1990

$\alpha = 1$   
 $\alpha = 10^{-1}$   
 $\alpha = 10^{-3}$   
 $\alpha = 10^{-5}$   
 $\alpha = 0$

Figure 1: Problem P6,  $\beta = 0.9$ , EN

$\alpha = 1$   
 $\alpha = 10^{-1}$   
 $\alpha = 10^{-3}$   
 $\alpha = 10^{-5}$   
 $\alpha = 0$

Figure 2: Problem P6,  $\beta = 0.9$ , GMRES

$\alpha = 10^4$   
 $\alpha = 10^2$   
 $\alpha = 1$   
 $\alpha = 10^{-2}$   
 $\alpha = 10^{-4}$

Figure 3: Problem P7,  $\beta = 0.9$ . EN

$\alpha = 10^4$   
 $\alpha = 10^2$   
 $\alpha = 1$   
 $\alpha = 10^{-2}$   
 $\alpha = 10^{-4}$

Figure 4: Problem P7,  $\beta = 0.9$ , GMRES



$\alpha = 1$   
 $\alpha = 10^{-1}$   
 $\alpha = 10^{-3}$   
 $\alpha = 10^{-5}$   
 $\alpha = 0$

Figure 5: Problem P8,  $\beta = 0.9$ , EN

$\alpha = 1$   
 $\alpha = 10^{-1}$   
 $\alpha = 10^{-3}$   
 $\alpha = 10^{-5}$   
 $\alpha = 0$

Figure 6: Problem P8,  $\beta = 0.9$ , GMRES

$\beta = 0.95$   
 $\beta = 0.9$   
 $\beta = 0.5$   
 $\beta = 0.2$   
 $\beta = 0$

Figure 7: Problem P9, EN

$\beta = 0.95$   
 $\beta = 0.9$   
 $\beta = 0.5$   
 $\beta = 0.2$   
 $\beta = 0$

Figure 8: Problem P9, GMRES

$\gamma = 3000$   
 $\gamma = 300$   
 $\gamma = 60$   
 $\gamma = 30$   
 $\gamma = 0$

Figure 9: Problem P10, EN

$\gamma = 3000$   
 $\gamma = 300$   
 $\gamma = 60$   
 $\gamma = 30$   
 $\gamma = 0$

Figure 10: Problem P10, GMRES

$$10_{\log(r_i^{EN})}$$

$$10_{\log(r_{2i}^G)}$$

Figure 11 Problem P10

- I  $\rho = 10^{-1}$
- II  $\rho = 10^{-2}$
- III  $\rho = 10^{-3}$
- IV  $\rho = 10^{-4}$

Figure 12 Problem P6