

Backward stability of FGMRES

(i.e. Using FGMRES to recover backward stability in mixed precision)

M. Arioli, I. S. Duff,

http://www.numerical.rl.ac.uk/people/marioli/marioli.html

Sparse Days, Cerfacs, Toulouse, June 2008 - p.1/21



Outline

Mixed precision

- Iterative refinement, GMRES and Flexible GMRES
- Roundoff error analysis
- Numerical experiments



Linear system

We wish to solve large sparse systems

Ax = b $A \in \mathbb{R}^{N \times N}$

to high accuracy using mixed precision arithmetic. For example, we might want to achieve double precision accuracy while using a single precision factorization of the matrix A. We will use this lower accuracy factorization as a preconditioner for FGMRES.



Very fast 32-bit arithmetic unit



- Very fast 32-bit arithmetic unit
- We use 32-bit arithmetic for factorization and triangular solves M is the fl(LU) of A and $||M A|| \le c(N)\sqrt{\varepsilon}||A||$ $(\varepsilon = 2.2 \times 10^{-16})$



- Very fast 32-bit arithmetic unit
- We use 32-bit arithmetic for factorization and triangular solves M is the fl(LU) of A and $||M A|| \le c(N)\sqrt{\varepsilon}||A||$ $(\varepsilon = 2.2 \times 10^{-16})$
- If $\kappa(A)\sqrt{\varepsilon} > 1$ then Iterative Refinement may not converge. FGMRES does



- Very fast 32-bit arithmetic unit
- We use 32-bit arithmetic for factorization and triangular solves M is the fl(LU) of A and $||M A|| \le c(N)\sqrt{\varepsilon}||A||$ $(\varepsilon = 2.2 \times 10^{-16})$
- If $\kappa(A)\sqrt{\varepsilon} > 1$ then Iterative Refinement may not converge. FGMRES does
- FGMRES backward stable



- Very fast 32-bit arithmetic unit
- We use 32-bit arithmetic for factorization and triangular solves M is the fl(LU) of A and $||M A|| \le c(N)\sqrt{\varepsilon}||A||$ $(\varepsilon = 2.2 \times 10^{-16})$
- If $\kappa(A)\sqrt{\varepsilon} > 1$ then Iterative Refinement may not converge. FGMRES does
- FGMRES backward stable
- GMRES is not always backward stable



GMRES and FGMRES

Let $r_0 = b - Ax_0$ and $\mathcal{K}_k(A, r_0)$ be the usual Krylov space GMRES

$$\min_{x \in x_0 + \mathcal{K}_k(A, r_0)} ||r_0 - Ax||_2 \qquad r_0 - Ax_k \perp A\mathcal{K}_k(A, r_0)$$



GMRES and FGMRES

Let $r_0 = b - Ax_0$ and $\mathcal{K}_k(A, r_0)$ be the usual Krylov space GMRES

$$\min_{x \in x_0 + \mathcal{K}_k(A, r_0)} ||r_0 - Ax||_2 \qquad r_0 - Ax_k \perp A\mathcal{K}_k(A, r_0)$$

GMRES Right preconditioning

$$AM^{-1}y = b \begin{cases} (AM^{-1}, r_0) \longrightarrow (A, r_0) \\ \mathcal{K}_k(AM^{-1}, r_0) \longrightarrow \mathcal{K}_k(A, r_0) \\ x_k = M^{-1}y_k \\ AM^{-1}V_k = V_{k+1}H_k \end{cases}$$



GMRES and FGMRES

Let $r_0 = b - Ax_0$ and $\mathcal{K}_k(A, r_0)$ be the usual Krylov space GMRES

$$\min_{x \in x_0 + \mathcal{K}_k(A, r_0)} ||r_0 - Ax||_2 \qquad r_0 - Ax_k \perp A\mathcal{K}_k(A, r_0)$$

GMRES Right preconditioning

$$AM^{-1}y = b \quad \begin{cases} (AM^{-1}, r_0) \longrightarrow (A, r_0) \\ \mathcal{K}_k(AM^{-1}, r_0) \longrightarrow \mathcal{K}_k(A, r_0) \\ x_k = M^{-1}y_k \\ AM^{-1}V_k = V_{k+1}H_k \end{cases}$$

Flexible GMRES Right preconditioning

$$Z_k \longrightarrow \mathcal{K}_k(A, r_0) \ x_k = x_0 + Z_k y_k \quad AZ_k = V_{k+1} H_k$$
$$Z_k = span(r_0, AM_1^{-1} r_0, \dots, \left(\prod_{j=0}^{k-1} AM_j^{-1}\right) r_0)$$

Sparse Days, Cerfacs, Toulouse, June 2008 – p.5/21



Right preconditioned GMRES and Flexible GMRES

procedure
$$[x] = \operatorname{right_Prec_GMRES(A,M,b)}$$

 $x_0 = M^{-1}b, r_0 = b - Ax_0 \text{ and } \beta = ||r_0||$
 $v_1 = r_0/\beta; k=0;$
while $||r_k|| > \mu(||b|| + ||A|| ||x_k||)$
 $k = k + 1;$
 $z_k = M^{-1}v_k; w = Az_k;$
for $i = 1, ..., k$ do
 $h_{i,k} = v_i^T w;$
 $w = w - h_{i,k}v_i;$
end for;
 $h_{k+1,k} = ||w||;$
 $v_{k+1} = w/h_{k+1,k};$
 $V_k = [v_1, ..., v_k];$
 $H_k = \{h_{i,j}\}_{1 \le i \le j+1}; 1 \le j \le k;$
 $y_k = \arg \min y ||\beta e_1 - H_k y||;$
 $x_k = x_0 + M^{-1}V_k y_k$ and $r_k = b - Ax_k;$
end while;
end procedure.

procedure [x] =FGMRES(A, M_i , b) $x_0 = M_0^{-1} b, r_0 = b - A x_0 \text{ and } \beta = ||r_0||$ $v_1 = r_0 / \beta; k = 0;$ while $||r_k|| > \mu(||b|| + ||A|| ||x_k||)$ k = k + 1; $z_k = M_k^{-1} v_k; w = A z_k;$ for $i = 1, \ldots, k$ do $h_{i,k} = v_i^T w;$ $w = w - h_{i,k} v_i;$ end for: $h_{k+1,k} = ||w||;$ $v_{k+1} = w/h_{k+1,k};$ $Z_{k} = [z_{1}, \ldots, z_{k}]; V_{k} = [v_{1}, \ldots, v_{k}];$ $H_k = \{h_{i,j}\}_{1 \le i \le j+1; 1 \le j \le k};$ $y_k = \arg\min_{y} ||\beta e_1 - H_k y||;$ $x_k = x_0 + Z_k y_k$ and $r_k = b - A x_k$; end while ; end procedure.



The roundoff error analysis of both FGMRES and GMRES can be done in three stages:



The roundoff error analysis of both FGMRES and GMRES can be done in three stages:

 Error analysis of the Arnoldi-Krylov process (Giraud and Langou, Björck and Paige, and Paige, Rozložník, and Strakoš). MGS applied to

$$z_{1} = M_{1}^{-1} r_{0} / ||r_{0}||, \quad z_{j} = M_{j}^{-1} v_{j}$$

$$C^{(k)} = (r_{0}, Az_{1}, Az_{2}, \dots, Az_{k}) = V_{k+1} R_{k}$$

$$\begin{bmatrix} ||r_{0}|| & H_{1,1} & \dots & H_{1,k} \\ 0 & H_{2,1} & \dots & H_{2,k} \\ 0 & 0 & \dots & H_{3,k} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & H_{k+1,k} \end{bmatrix}$$



The roundoff error analysis of both FGMRES and GMRES can be done in three stages:

- 1. Error analysis of the Arnoldi-Krylov process (Giraud and Langou, Björck and Paige, and Paige, Rozložník, and Strakoš).
- 2. Error analysis of the Givens process used on the upper Hessenberg matrix H_k in order to reduce it to upper triangular form.



The roundoff error analysis of both FGMRES and GMRES can be done in three stages:

- 1. Error analysis of the Arnoldi-Krylov process (Giraud and Langou, Björck and Paige, and Paige, Rozložník, and Strakoš).
- 2. Error analysis of the Givens process used on the upper Hessenberg matrix H_k in order to reduce it to upper triangular form.
- 3. Error analysis of the computation of x_k in FGMRES and GMRES.



The roundoff error analysis of both FGMRES and GMRES can be done in three stages:

- 1. Error analysis of the Arnoldi-Krylov process (Giraud and Langou, Björck and Paige, and Paige, Rozložník, and Strakoš).
- 2. Error analysis of the Givens process used on the upper Hessenberg matrix H_k in order to reduce it to upper triangular form.
- 3. Error analysis of the computation of x_k in FGMRES and GMRES.

The first two stages of the roundoff error analysis are the same for both FGMRES and GMRES. The last stage is specific to each algorithm.

Roundoff error analysis of FGMRES

Theorem 1. If we apply **FGMRES** to solve Ax = b, using finite-precision arithmetic conforming to IEEE standard with relative precision ε and under the hypotheses:

$$2.12(n+1)arepsilon < 0.01$$
 and $c_0(n)arepsilon\kappa(C^{(k)}) < 0.1 \; orall k$

where

cience & Technology Facilities Council

Rutherford Appleton Laboratorv

$$c_0(n) = 18.53n^{\frac{3}{2}}$$

and

$$|\bar{s}_k| < 1 - \varepsilon, \ \forall k,$$

where \bar{s}_k are the sines computed during the Givens algorithm applied to H_k in order to compute \bar{y}_k , then there exists \hat{k} , $\hat{k} \leq n$ such that, $\forall k \geq \hat{k}$, we have

$$||b - A\bar{x}_{k}|| \leq c_{1}(n,k)\varepsilon \left(||b|| + ||A|| ||\bar{x}_{0}|| + ||A|| ||\bar{x}_{k}|| |\bar{y}_{k}|| + ||A\bar{z}_{k}|| ||\bar{y}_{k}|| \right) + \mathcal{O}(\varepsilon^{2}).$$



The symmetric indefinite case

A particular and important case arises in saddle-point problems where the coefficient matrix is of the form

$$A = \begin{bmatrix} H & B \\ B^T & 0 \end{bmatrix}$$



Why FGMRES for symmetric case?

The computed values of Gaussian factorization $\hat{L} \hat{D}$ are affected by roundoff: $M = \hat{L}\hat{D}\hat{L}^T$ and $||E|| = ||M - A|| \le c(n)\varepsilon||A||$ with $E \ne E^T$



Why FGMRES for symmetric case?

- The computed values of Gaussian factorization $\hat{L} \hat{D}$ are affected by roundoff: $M = \hat{L}\hat{D}\hat{L}^T$ and $||E|| = ||M - A|| \le c(n)\varepsilon||A||$ with $E \ne E^T$
- Thus $M^{-1}A \neq AM^{-1}$ and the preconditioned matrix is non symmetric



Why FGMRES for symmetric case?

- The computed values of Gaussian factorization $\hat{L} \hat{D}$ are affected by roundoff: $M = \hat{L}\hat{D}\hat{L}^T$ and $||E|| = ||M - A|| \le c(n)\varepsilon||A||$ with $E \ne E^T$
- Thus $M^{-1}A \neq AM^{-1}$ and the preconditioned matrix is non symmetric FGMRES is then the only way



Multifrontal approach: HSL_MA57

In order to reduce the fill-in during the LDL^T factorization

• We scale and reorder the entries of A



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold
- However, also this can be unsatisfactory: the numerical pivot strategy is still disrupting the ordering we have chosen and increases the fill-in



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold
- However, also this can be unsatisfactory: the numerical pivot strategy is still disrupting the ordering we have chosen and increases the fill-in
- An ALTERNATIVE is to use Static Pivoting, by replacing the pivot a_k failing the test by $a_k + \tau$ and CONTINUE.



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold
- However, also this can be unsatisfactory: the numerical pivot strategy is still disrupting the ordering we have chosen and increases the fill-in
- An ALTERNATIVE is to use Static Pivoting, by replacing the pivot a_k failing the test by $a_k + \tau$ and CONTINUE.

• We thus have factorized $A + E = LDL^T = M$ where $|E| \le \tau I$



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold
- However, also this can be unsatisfactory: the numerical pivot strategy is still disrupting the ordering we have chosen and increases the fill-in
- An ALTERNATIVE is to use Static Pivoting, by replacing the pivot a_k failing the test by $a_k + \tau$ and CONTINUE.
- We thus have factorized $A + E = LDL^T = M$ where $|E| \le \tau I$
- Several codes use (or have an option for) this device:



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold
- However, also this can be unsatisfactory: the numerical pivot strategy is still disrupting the ordering we have chosen and increases the fill-in
- An ALTERNATIVE is to use Static Pivoting, by replacing the pivot a_k failing the test by $a_k + \tau$ and CONTINUE.
- We thus have factorized $A + E = LDL^T = M$ where $|E| \le \tau I$
- Several codes use (or have an option for) this device:
 - SuperLU (Demmel and Li)



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold
- However, also this can be unsatisfactory: the numerical pivot strategy is still disrupting the ordering we have chosen and increases the fill-in
- An ALTERNATIVE is to use Static Pivoting, by replacing the pivot a_k failing the test by $a_k + \tau$ and CONTINUE.
- We thus have factorized $A + E = LDL^T = M$ where $|E| \le \tau I$
- Several codes use (or have an option for) this device:
 - SuperLU (Demmel and Li)
 - PARDISO (Gärtner and Schenk)



- We scale and reorder the entries of A
- We weaken the numerical pivot strategy by using a threshold
- However, also this can be unsatisfactory: the numerical pivot strategy is still disrupting the ordering we have chosen and increases the fill-in
- An ALTERNATIVE is to use Static Pivoting, by replacing the pivot a_k failing the test by $a_k + \tau$ and CONTINUE.
- We thus have factorized $A + E = LDL^T = M$ where $|E| \le \tau I$
- Several codes use (or have an option for) this device:
 - SuperLU (Demmel and Li)
 - PARDISO (Gärtner and Schenk)
 - MA57 (Duff and Pralet)



GMRES error bounds depend on $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$. (Arioli, Duff, Gratton, and Pralet SISC 2007)



GMRES error bounds depend on $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$. (Arioli, Duff, Gratton, and Pralet SISC 2007) For sparse matrices $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$ can be much larger than ||A||.



- **GMRES error bounds depend on** $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$. (Arioli, Duff, Gratton, and Pralet SISC 2007)
- For sparse matrices $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$ can be much larger than ||A||.
- For the static pivot the growth can be dramatic.



- **GMRES** error bounds depend on $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$. (Arioli, Duff, Gratton, and Pralet SISC 2007)
- For sparse matrices $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$ can be much larger than ||A||.
- For the static pivot the growth can be dramatic.
- Theorem 1 shows that FGMRES does not depend on $|| |\hat{L}| |\hat{D}| |\hat{L}^T| ||$.



Roundoff error right preconditioned GMRES

Theorem 2

We assume of applying Iterative Refinement for solving $M(\bar{x}_k - \bar{x}_0) = \bar{V}_k \bar{y}_k$ at last step.

Under the Hypotheses of Theorem 1 and $|c(n) \varepsilon \kappa(M) < 1|$

$$\exists \hat{k}, \qquad \hat{k} \le n$$

such that, $\forall k \geq \hat{k}$, we have

 $\begin{aligned} ||b - A\bar{x}_{k}|| &\leq c_{1}(n,k)\varepsilon\left\{ ||b|| + ||A|| \, ||\bar{x}_{0}|| + ||A|| \, ||\bar{Z}_{k}|| \, ||M(\bar{x}_{k} - \bar{x}_{0})|| + \\ &||AM^{-1}|| \, ||M|| \, ||\bar{x}_{k} - \bar{x}_{0}|| + \\ &||AM^{-1}|| \, |||\hat{L}| \, |\hat{D}| \, |\hat{L}^{T}| \, |||M(\bar{x}_{k} - \bar{x}_{0})|| \right\} + \mathcal{O}(\varepsilon^{2}). \end{aligned}$



Test Problems

	n	nnz	Description
CONT_201	80595	239596	KKT matrix Convex QP (M2)
CONT_300	180895	562496	KKT matrix Convex QP (M2)
TUMA_1	22967	76199	Mixed-Hybrid finite-element

Test problems



$|| \, |\hat{L}| \, |\hat{D}| \, |\hat{L}^{T}| \, || \, { m vs} \, 1/ au$



Sparse Days, Cerfacs, Toulouse, June 2008 - p.15/21



Numerical experiments



FGMRES on CONT-201 test example



Numerical experiments using mixed precision



GMRES on CONT-201 test example



 $\blacksquare A = QDW$ with Q and W random orthogonal matrices $D = \text{diag}\{d_i\}$

$$d_i = 10^{-c\left(\frac{i-1}{n-1}\right)^{\gamma}}$$



 $\blacksquare A = QDW$ with Q and W random orthogonal matrices $D = \text{diag}\{d_i\}$

$$d_i = 10^{-c\left(\frac{i-1}{n-1}\right)^{\gamma}}$$

The singular values lie between 1 and 10^{-c} , the condition number is 10^{c} , and the distribution can be skewed by altering γ . γ equal to 1 gives a log-linear uniform distribution, values of γ greater than 1 skew towards 1 and values of γ less than 1 towards 10^{-c} .



• A = QDW with Q and W random orthogonal matrices $D = \text{diag}\{d_i\}$ $d_i = 10^{-c(\frac{i-1}{n-1})^{\gamma}}$

Selected Sparse Matrices



• A = QDW with Q and W random orthogonal matrices $D = \text{diag}\{d_i\}$ $d_i = 10^{-c(\frac{i-1}{n-1})^{\gamma}}$

Selected Sparse Matrices

Forward and backward substitution



• A = QDW with Q and W random orthogonal matrices $D = \text{diag}\{d_i\}$ $d_i = 10^{-c(\frac{i-1}{n-1})^{\gamma}}$

- Selected Sparse Matrices
- Forward and backward substitution
 - the vector \bar{z}_k is computed using the forward and backward substitution algorithm in single precision on the single precision conversion of vector \bar{v}_k ,
 - the vector \overline{z}_k is computed using the forward and backward substitution algorithm in double precision on \overline{v}_k after we converted the factors L and U to double precision.



Random dense matrices. $n = 200, c = 8.2, \gamma = 1$

Single Precision				Double Precision					
Total It	Inner it	SR	$ Aar{Z}_{\hat{k}} $	$ {ar Z}_{\hat k}^{} {ar y}_{\hat k}^{} $	Total It	Inner it	SR	$ Aar{Z}_{\hat{k}} $	$ {ar Z}_{\hat k}^{} {ar y}_{\hat k}^{} $
26	26	2.5e-16	7.4e+00	1.9e+02	20	20	1.7e-16	8.0e+00	7.3e+01
27	27	6.6e-16	4.2e+00	4.7e+02	20	20	2.0e-16	3.9e+00	5.9e+01
25	25	1.7e-16	3.3e+00	5.9e+01	20	20	2.7e-16	3.5e+00	4.0e+01
52	52	3.9e-15	4.6e+01	3.0e+03	20	20	1.1e-15	4.5e+01	4.7e+02
88	36	1.1e-16	4.6e+01	6.0e-04	25	5	1.5e-16	4.8e+01	1.5e-05
24	24	1.3e-16	2.0e+00	3.8e+01	20	20	2.6e-16	2.2e+00	2.8e+01
31	31	2.5e-16	8.8e+00	1.7e+02	20	20	1.9e-16	1.1e+01	8.4e+01
24	24	2.0e-16	3.5e+00	1.2e+02	20	20	2.0e-16	3.9e+00	6.9e+01
24	24	1.8e-16	2.7e+00	8.8e+01	20	20	6.2e-16	3.0e+00	5.8e+01
26	26	2.7e-16	3.2e+00	1.5e+02	20	20	3.2e-16	3.5e+00	3.6e+01
44	44	5.7e-16	1.9e+01	5.9e+02	20	20	4.0e-16	2.0e+01	1.6e+02

$$(SR = \frac{||b - A\overline{\mathbf{x}}_{\hat{k}}||}{(||A|| \, ||\overline{\mathbf{x}}_{\hat{k}}|| + ||b||)})$$



MA57 sparse tests using mixed precision

Matrix Id	n	Iterative r	refinement	FGMRES					
		Total It	SR	Total It	Inner it	SR	$ A \bar{Z}_{\hat{k}} $	$ ar{Z}_{\hat{k}} ar{y}_{\hat{k}} $	
bcsstk20	485	30	2.1e-15	2	2	1.4e-11	1.7e+00	4.6e+02	
$\kappa(A) = 5.10^9$				4	2	3.4e-14	1.6e+00	3.8e-01	
				6	2	7.2e-17	1.6e+00	5.6e-04	
bcsstm27	1224	22	1.6e-15	2	2	5.8e-11	1.7e+00	2.7e+01	
$\kappa(A) = 5.10^9$				4	2	1.8e-11	6.3e-01	1.3e+00	
				6	2	6.0e-13	2.0e+00	7.6e-02	
				8	2	1.5e-13	1.7e+00	1.0e-02	
				10	2	1.2e-14	1.7e+00	1.9e-03	
				12	2	2.6e-15	1.8e+00	1.7e-04	
				14	2	1.8e-16	1.6e+00	4.3e-05	
s3rmq4m1	5489	16	2.2e-15	2	2	3.5e-11	1.0e+00	8.6e+01	
$\kappa(A) = 4.10^9$				4	2	2.1e-13	1.1e+00	3.2e-01	
				6	2	4.5e-15	1.7e+00	6.4e-03	
				8	2	1.1e-16	1.6e+00	1.3e-04	
s3dkq4m2	90449	53	1.1e-10	10	10	6.3e-17	1.2e+00	1.2e+03	
$\kappa(A) = 7.10^{10}$									

Sparse matrices results (
$$SR = \frac{||b - A\overline{\mathbf{x}}_{\hat{k}}||}{(||A|| ||\overline{\mathbf{x}}_{\hat{k}}|| + ||b||)}$$
)





IR (PLAN A) does not always work with mixed precision GMRES is also sensitive and not robust



- GMRES is also sensitive and not robust
- **FGMRES** is robust and less sensitive



- GMRES is also sensitive and not robust
- **FGMRES** is robust and less sensitive
- Gains from restarting



- GMRES is also sensitive and not robust
- **FGMRES** is robust and less sensitive
- Gains from restarting
- PLAN B is working