the Generalized Minimal Residual (GMRES) algorithm for solving non-symmetric linear systems. *

this method was introduced by Saad and Schultz in 1986 to solve iteratively matix equations

$$A x = b,$$

where the matrix $A$ is not symmetric.

① Motivation: the method of conjugate gradients as a Galerkin method. ‒

We have seen that if the $k$-th iterate of the CG method is $x^k$, the error, $e^k = x - x^k$ satisfies the equality

$$w^T A e^k = 0 \qquad \forall w \in K_k,$$

where

$$K_k = \text{span} \left\{ A^i r_0 \right\}_{i=0}^{k-1}, \quad r_0 = A x_0 - b.$$

* From: Iterative methods for sparse linear systems, Y. Saad.

We immediately see that the above equation states that the $k$-th residual,

$$r_k = A x_k - b,$$

is orthogonal to $K_k$ since

$$
\begin{aligned}
w^T r_k &= w^T (A x_k - b) \\
&= w^T (A x_k - A x) \\
&= -w^T A e_k \\
&= 0 \qquad\qquad \forall\, w \in K_k.
\end{aligned}
$$

Now, let $\{v_i\}_{i=1}^{k}$ be an $l_2$-orthonormal basis of $K_k$ and set

$$V_k = (v_1, v_2 \ldots v_k).$$

$V_k$ is an $N \times k$ matrix, where $N = $ order of $A$. With this notation, we can write that

$$x_k = x_0 + V_k y_k,$$

where $y_k \in \mathbb{R}^k$ is the solution of

$$(V_k^T A V_k)\, y_k = - V_k^T r_0$$

Note that the above equation defines a Galerkin approximation, $V_k Y_k$, in $K_k$ to the solution of

$$A z = -r_0.$$

Of course, we have that $X = Y_0 + Z$! In other words, the iterates of the CG method are of the form $X_k = X_0 + Z_k$, where $Z_k$ is the A-projection of $Z$ into the space $K_k$.

Note that if we set

$$v_1 = r_0 / \| r_0 \|,$$

we have

$$X_k = X_0 + V_k Y_k,$$
$$Y_k = - \| r_0 \| H_k^{-1} e_1,$$
$$H_k = V_k^T A V_k.$$

This is a rewriting of the CG method that is possible to use for matrices that are not symmetric and positive definite. Note that if $k$ is small, the matrix $H_k$ is easy to invert!

4

<u>Arnoldi's orthogonalization.</u> (of $K_k$).

We can obtain the basis $\{v_i\}_{i=1}^k$ of the Krylov subspace $K_k$ as follows:

(1) $\qquad \hat{v}_1 = r_0, \qquad v_1 = \hat{v}_1 / \| \hat{v}_1 \|.$

(2) $\qquad$ For $\quad j = 1, \dots, k-1 :$

(3) $\qquad\qquad h_{ij} = v_i^T A v_j, \qquad i = 1, \dots, i$

$\qquad\qquad \hat{v}_{j+1} = A v_j - \sum_{i=1}^{j} h_{ij} v_i$

(4) $\qquad\qquad h_{j+1,j} = \| \hat{v}_{j+1} \|, \qquad v_{j+1} = \hat{v}_{j+1} / \| \hat{v}_{j+1} \|.$

$\qquad$ endFor

Notice that if $\ell \leq j$,

$$v_\ell^T \hat{v}_{j+1} = v_\ell^T A v_j - \sum_{i=1}^{j} h_{ij} v_\ell^T v_i$$

So, if $\quad v_\ell^T v_i = \delta_{\ell i}$ for $i, \ell \leq j$ then we have that $\quad v_\ell^T \hat{v}_{j+1} = 0$. Hence, if $\hat{v}_{j+1} \neq 0,$ we have that $\quad v_\ell^T v_i = 0 \quad$ for $i, \ell \leq j+1.$

Next, let us show that if $\dim K_k = k$, then the algorithm (1),(2) does <u>not</u> break down.

5

Let us proceed by induction on $k$. For $k = 1$, the result is obvious. Assume the result holds for $k = m$ and let us show it holds for $k = m+1$. Since it holds for $k = m$, we have that

$$K_m = \text{span} \left\{ A^i r_0 \right\}_{i=0}^{m-1}$$

$$= \text{span} \left\{ A^{i-1} v_1 \right\}_{i=1}^{m}$$

$$= \text{span} \left\{ v_1, \ldots, v_m \right\}.$$

Now, this implies that

$$K_{m+1} = \text{span} \left\{ v_1, \ldots, v_m, A^m v_1 \right\}.$$

By (3) and (4), for $j \leq m$,

$$\hat{v}_{j+1} = A v_j - \sum_{i=1}^{j} h_{i,j} v_i$$

$$\Rightarrow \quad A v_j = \sum_{i=1}^{j+1} h_{i,j} v_i$$

$$\Rightarrow \quad A v_j \in \text{span} \left\{ v_1, \ldots, v_{j+1} \right\}$$

$$\Rightarrow \quad A^m v_1 \in \text{span} \left\{ v_1, \ldots, v_{m+1} \right\}.$$

and so, $v_{m+1} \neq 0$. This implies that the algorithm does not break down.

Finally, notice that if $m > j+1$,

$$U_m^T A v_j = \sum_{i=1}^{j+1} h_{ij} v_m^T v_i = 0,$$

by the inductive hypothesis. This shows that $H_k$ is an upper Hessemben matrix which is tridiagonal if $A$ is symmetric.

③ the Full Orthogonalization Method (FOM)

the above discussion motivates the introduction of the FOM:

- Pick the initial guess $x_0$, compute $r_0 = Ax_0 - b$ and set $v_1 = r_0 / \|r_0\|$

- For $j = 1, 2, \ldots, k$ do
$$h_{ij} = v_i^T A v_j \quad , \quad i = 1, 2, \ldots, j$$
$$\hat{v}_{j+1} = A v_j - \sum_{i=1}^{j} h_{ij} v_i$$
$$h_{j+1, j} = \|\hat{v}_{j+1}\|$$
$$v_{j+1} = \hat{v}_{j+1} / h_{j+1, j}$$
endFor

- Set $x_k = x_0 + V_k y_k$ where
$$y_k = -H_k^{-1} \|r_0\| e_1$$

Note that the algorithm breaksdown only if
for some $j \leq k$, $h_{j+1,j} = 0$. Next we show that
this happens only if $r^j = 0$!

The relations (3) can be rewritten in matrix
form as

$$A V_k = V_{k+1} \overline{H}_k .$$

Note that

$$H_k = V_k^T A V_k = V_k^T V_{k+1} \overline{H}_k ,$$

hence

$$\overline{H}_k = \begin{bmatrix} H_k \\ 0\ 0 \dots\ 0\ h_{k+1,k} \end{bmatrix} .$$

Since

$$X_k = X_0 + V_k Y_k ,$$

we get

$$r_k = r_0 + A V_k Y_k$$
$$= r_0 + V_{k+1} \overline{H}_k Y_k$$
$$= r_0 + V_k H_k Y_k + h_{k+1,k} (e_k^T Y_k) v_{k+1}$$

$\Rightarrow$
$$r_k = h_{k+1,k} (e_k^T y_k) v_{k+1}.$$

this means that

$$\|r_k\| = h_{k+1,k} |e_k^T y_k|,$$

and $r_k = 0$ if $h_{k+1,k} = 0$, as claimed.

④ <u>the GMRES method.</u>

the difference, and only difference, between the FOM and the GMRES methods is the way in which $y_k$ is computed. In FOM, $y_k$ is computed by requesting that

$$r^k = A(x_0 + V_k y_k) - b$$
$$= r_0 + A V_k y_k$$

be <u>orthogonal</u> to $K_k$. In GMRES, instead, it is asked that

$$\|r^k\|^2 = \|r^k(y_k)\|^2 =: J(y_k).$$

is a <u>minimum</u>, hence the name of the method!

Next, notice that

$$J(y) = \| r_0 + AV_k \, y \|^2$$
$$= \| r_0 + V_{k+1} \, \bar{H}_k \, y \|^2$$
$$= \| \, \|r_0\| e_1 + \bar{H}_k y \|^2$$

this implies that

$$Y_k = - (\bar{H}_k^T \bar{H}_k)^{-1} \bar{H}_k e_1 \quad \|r_0\|,$$

and shows that this method is different than the FOM.

To actually compute $Y_k$, we can still exploit the structure of $J$. Suppose that we obtain a $(k+1) \times (k+1)$ matrix $Q_k$, the accumulated product of rotation matrices, and an upper triangular matrix $R_k$ such that

$$Q_k \bar{H}_k = R_k$$

(Note that the last row of $R_k$ is zero!). then

$$J(y) = \| Q_k ( \|r_0\| e_1 + \bar{H}_k y ) \|^2$$
$$= \| \underbrace{\|r_0\| Q_k e_1}_{-g_k} + R_k y \|^2$$

Now if we write

$$g_k = \begin{bmatrix} \hat{g}_k \\ g_{k,k+1} \end{bmatrix}, \quad R_k = \begin{bmatrix} \hat{R}_k \\ 0 \dots 0 \end{bmatrix},$$

then $\quad y_k = \hat{R}_k^{-1} \hat{g}_k \quad$ and $\quad J(y_k) = |g_{k,k+1}|$.

## ⑤ Convergence analysis

to analyze the GMRES method, we begin by noting that if the Arnoldi process does not break down, then $\hat{R}_k$ is invertible. To see this, consider the beginning of the algorithm leading to the equation $\varphi_k \bar{H}_k = R_k$:

$$\begin{bmatrix} c & -s & 0 & 0 & 0 \\ s & c & 0 & 0 & 0 \\ 0 & 0 & & & \\ 0 & 0 & & Id & \\ 0 & 0 & & & \end{bmatrix} \begin{bmatrix} r & x & x & x \\ h & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix} = \begin{bmatrix} t & y & y & y \\ 0 & y & y & y \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix}$$

where $\quad c = \dfrac{r}{\sqrt{r^2 + h^2}}, \quad s = -\dfrac{h}{\sqrt{r^2 + h^2}}, \quad t = \sqrt{r^2 + h^2}$

We see that if $h \neq 0$, then $t > 0$. Now notice that in every multiplication by a rotation, $h = h_{j+1,j}$, which is non-zero because the Arnoldi process did not break down for $j \leq k$! This establishes that $\hat{R}_k$ is invertible.

this means that the GMRES method can only break down if for some $j < N$, $h_{j+1,j} = 0$. Next, we prove that this happens if and only if $r_j = 0$, that is, if and only if $A x_j = b$.

Assume that $h_{j+1,j} = 0$. this implies that

$$A V_j = V_j H_j.$$

In this case, we have

$$
\begin{aligned}
J(y) &= \| r_0 + A V_j \, y \|^2 \\
&= \| r_0 + V_j H_j \, y \|^2 \\
&= \| V_j ( \|r_0\| e_1 + H_j y ) \|^2 \\
&= \| \|r_0\| e_1 + H_j y \|^2 \\
&= 0
\end{aligned}
$$

if

$$ y = - \|r_0\| H_j^{-1} e_1 .$$

this means that $r_j = 0$, as claimed.

Assume now that $r_j = 0$ and that $r_i \neq 0$ $i < j$. then

$$
\| r_j \| = \underset{\underset{(-\|r_0\| \, \Theta_{j+1} \, e_1)_{j+1}}{\|}}{| g_{j,j+1} |} = |s_j| \underset{\underset{(-\|r_0\| \, \Theta_j \, e_1)_j}{\|}}{| g_{j-1,j} |} = |s_j| \| r_{j-1} \|
$$

this implies $\quad 0 = s_j \Rightarrow \quad h_{j+1,j} = 0$, as claimed.

We have thus proven the following result.

<u>theorem</u>. The GMRES method converges in at most $N$ iterations, where $N$ = order of $A$.

We can also obtain an error estimate of the residual in some cases.

<u>theorem</u>. Let $A$ be a diagonalizable matrix. Then

$$\| r_m \|_2 \leq K \in_m \| r_0 \|_2,$$

where

$$K = \| X \|_2 \, \| \vec{X} \|_2 \, ,$$
$$\in_m = \min_{\substack{p \in p^m \\ p(0) = 1}} \max_{1 \leq i \leq N} | p(\lambda_i) |$$

and $\quad \text{diag} \{ \lambda_i, i=1,...,N \} = \Lambda = \vec{X} A X.$

<u>Proof</u> Since

$$X_m = X_0 + V_m Y_m,$$

we have

$$r_m = r_0 + A V_m Y_m$$

and since $\qquad V_m = K_m ,$

$$r_m = p(A) \, r_0 \qquad , \quad \text{degree } p = m, \quad p(0) = 1.$$

Since $x_m$ <u>minimizes the residual</u> over $x_0 + K_m$, we set

$$\| r_m \|_2 = \min_{\substack{p \in p^m \\ p(0) = 1}} \| p(A) \, r_0 \| .$$

But

$$p(A) = X \, p(\Lambda) \, X^{-1}$$

and so

$$\| p(A) \|_2 \leq K \max_{1 \leq i \leq N} | p(\lambda_i) | .$$

This completes the proof. $\qquad \qquad \square$

If all the eigenvalues of $A$ are inside the ellipse of center $(c, 0)$ focal distance $d$ and major semi-axis $a$