



# Samenvatting

## Scientific Computing - Collegejaar 2015-2016

*College aantekeningen*

### Disclaimer

De informatie in dit document is afkomstig van derden. W.I.S.V. 'Christiaan Huygens' betracht de grootst mogelijke zorgvuldigheid in de samenstelling van de informatie in dit document, maar garandeert niet dat de informatie in dit document compleet en/of accuraat is, noch aanvaardt W.I.S.V. 'Christiaan Huygens' enige aansprakelijkheid voor directe of indirekte schade welke is ontstaan door gebruikmaking van de informatie in dit document.

De informatie in dit document wordt slechts voor algemene informatie in dit documentdoeleinden aan bezoekers beschikbaar gesteld. Elk besluit om gebruik te maken van de informatie in dit document is een zelfstandig besluit van de lezer en behoort uitsluitend tot zijn eigen verantwoordelijkheid.

# Samenvatting Scientific Computing 2015/2016

## Chapter 1

### Simulation

- Setting up a model (system of diff. eq's)
- Analytical treatment (properties of solution)
- Numerical treatment (discretizing)
- Implementation
- Embedding
- Visualization
- Validation

### Errors

Errors:	Absolute	$ u - \hat{u}_h $
	Relative	$\frac{ u - \hat{u}_h }{ u }$

- Modelling errors
- Error in data
- Solution errors
- Round-off errors

## Chapter 2

- $\sigma(A) = \{\lambda : \lambda \text{ eigenvalue of } A\}$  is the spectrum of  $A$ .
- Eigenvalues of a symmetric matrix are real.
- A matrix is positive (semi-)definite if  
 $u \in \mathbb{R}^n \setminus \{0\}: u^T A u > 0 \quad (u^T A u \geq 0)$
- A symmetric positive definite (spd) matrix  $A$  has  
 $\sigma(A) \subset \mathbb{R}^+$
- $\langle u, v \rangle_A := u^T A v$  is an inner product. If  $\langle u, v \rangle_A = 0$ ,  $u$  and  $v$  are  $A$ -conjugate.
- $\|R\|_p := \max_{\substack{\|u\|_p=1 \\ u \in \mathbb{R}^n \setminus \{0\}}} \|Ru\|_p = \max_{u \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ru\|_p}{\|u\|_p}$
- If  $A$  is spd, then  $\|A\|_2 = \lambda_{\max}(A)$ .
- In general,  $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$  max. abs. column sum  
 $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$   
 $\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$  max. abs. row sum
- The condition number  $K_p(A)$  controls  $\|u\|$  in terms of  $\|f\|$  in the system  $Au = f$ .  
 $K_p(A) = \|A\|_p \cdot \|A^{-1}\|_p \Rightarrow K_2(A) = \frac{\sqrt{\lambda_{\max}(A^T A)}}{\sqrt{\lambda_{\min}(A^T A)}}$ 
  - $A$  symm:  $\frac{|\lambda_{\max}(A)|}{|\lambda_{\min}(A)|}$
  - $A$  spd:  $\frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$

- Spectral radius  $\rho(A) := \max \{ |\lambda_i| : \lambda_i \in \sigma(A) \}$
- $\rho(A) \leq \|A\|$ : take a random  $\lambda \in \sigma(A)$ . Then  
 $|\lambda| \|u\| = \|\lambda u\| = \|Au\| \leq \|A\| \|u\| \Rightarrow |\lambda| \leq \|A\| \quad \#$
- A matrix  $A$  is irreducible if no permutation  $P$  exists such that  $PAP^T$  is block upper triangular.
- A matrix is an M-matrix iff
  - $a_{ii} > 0 \quad \forall i$
  - $a_{ij} \leq 0 \quad \forall i \neq j$
  - $A$  non-singular
  - $A^{-1} \geq 0$ .

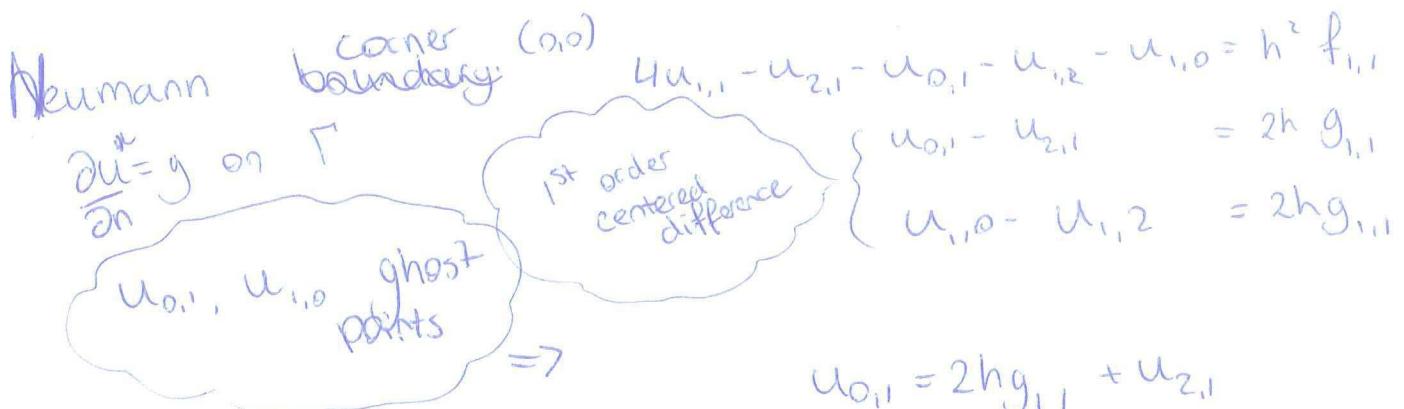
Also, the following are sufficient conditions:

- $a_{ii} > 0 \quad \forall i$
- $a_{ij} \leq 0 \quad \forall i \neq j$
- $A$  is irreducible and diagonally dominant

# Chapter 3

2D

- Dirichlet boundary:  $\frac{1}{h^2} \begin{bmatrix} 0 & -1 & 0 \\ 0 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ ,  $f_{2,j}^h + \frac{1}{h^2} b_{1,j}^h$
- Corner:  $\frac{1}{h^2} \begin{bmatrix} 0 & -1 & 0 \\ 0 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ ,  $f_{2,2}^h + \frac{1}{h^2} b_{1,2}^h + \frac{1}{h^2} b_{2,1}^h$



$$u_{0,1} = 2h g_{1,1} + u_{2,1}$$

$$u_{1,0} = 2h g_{1,1} + u_{1,2}$$

$$\Rightarrow 4u_{1,1} - 2u_{1,2} - 2u_{2,1} = h^2 f_{1,1} + 4h g_{1,1}$$

$$\Rightarrow \frac{1}{h^2} \begin{bmatrix} 0 & -2 & 0 \\ 0 & 4 & -2 \\ 0 & 0 & 0 \end{bmatrix}, \quad h f_{1,1} + \frac{4}{h} g_{1,1}$$

• h-scaled norm:

$$\|u(x) - u^n\|_{h,2} = \sqrt{h \sum_{i=0}^N (u(x_i) - u_i^n)^2} = O(h^2)$$

## • Alternative Grid Orderings

0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

Red/Black

Diagonal

- The resulting matrix is sparse, symmetric, irreducibly diagonal dominant (so pd by Gershgorin  $\rightarrow$  spd)  
So it is also an M-matrix

The matrix with  $\frac{1}{h^2} \begin{bmatrix} -1 & 2 & -1 \end{bmatrix}$  as stencil has eigenvalues

$$\lambda_k^n = \frac{2}{h^2} [1 - \cos(\pi h k)] \quad k=1, \dots, N-1$$

$$v_k^n = \frac{1}{\sqrt{N-1}} \begin{pmatrix} \sin(\pi h x_1) \\ \vdots \\ \sin(\pi h x_{N-1}) \end{pmatrix} = \frac{1}{\sqrt{N-1}} \begin{pmatrix} \sin(\pi h k) \\ \vdots \\ \sin(\pi h (N-1)k) \end{pmatrix} \quad k=1, \dots, N-1$$

eigen vectors

Proof: for nodes  $i$  not connected to the boundary:

Let  $k \in \{1, \dots, N-1\}$ .

$$\begin{aligned} [A^n v_k^n]_i &= \sum_{\alpha=1}^{N-1} a_{i,\alpha} [v_k^n]_\alpha \\ &= \frac{1}{h^2} \frac{1}{\sqrt{N-1}} [-\sin(\pi(i-1)hk) + 2\sin(\pi i hk) - \sin(\pi(i+1)hk)] \\ \text{As } \sin(\alpha+\beta) &\neq \sin(\alpha-\beta) = 2\sin(\alpha)\cos(\beta): \begin{cases} \alpha = \pi i hk \\ \beta = \pi h k \end{cases} \\ &= \frac{1}{h^2} \frac{1}{\sqrt{N-1}} [2\sin(\pi i hk) - 2\sin(\pi i hk)\cos(\pi h k)] \\ &= \frac{2}{h^2} \frac{1}{\sqrt{N-1}} [1 - \cos(\pi h k)] \sin(\pi i hk) \\ &= \frac{2}{h^2} [1 - \cos(\pi h k)] [v_k^n]_i \\ &= \lambda_k^n [v_k^n]_i \quad (2 \leq i \leq N-2) \end{aligned}$$

for  $i=1$ :

Let  $k \in \{1, \dots, N-1\}$

$$\begin{aligned} [A^n v_k^n]_1 &= \sum_{\alpha=1}^{N-1} a_{1,\alpha} [v_k^n]_\alpha \\ &= \frac{1}{h^2} \frac{1}{\sqrt{N-1}} [2\sin(\pi h k) - \sin(\pi \cdot 2 \cdot h k)] \\ \text{as } \alpha = \beta = \pi h k &= \frac{1}{h^2} \frac{1}{\sqrt{N-1}} [2\sin(\pi h k) - 2\sin(\pi h k)\cos(\pi h k)] \\ &= \frac{2}{h^2} [1 - \cos(\pi h k)] [v_k^n]_1 \end{aligned}$$

for  $i=N-1$

$$\begin{aligned} [A^n v_k^n]_{N-1} &= \sum_{\alpha=1}^{N-1} a_{N-1,\alpha} [v_k^n]_{N-1} \\ &= \frac{1}{h^2} \frac{1}{\sqrt{N-1}} [\sin(\pi(N-2)hk) + 2\sin(\pi(N-1)hk)] \end{aligned}$$

$$\alpha = \frac{1}{h^2} \frac{1}{\sqrt{N+1}} [2 \sin(\pi(N-1)hk) \cos(\pi hk) + 2 \sin(\pi(N-1)hk)]$$

$$\beta = \frac{1}{h^2} \frac{1}{\sqrt{N+1}} [1 - \cos(\pi hk)] \sin(\pi(N-1)hk)$$

$$= \lambda_k^n [v_k^n]_{N-1}$$

14

Similarly for  $\frac{1}{h^2} \begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$  we find

$$\lambda_{k,l}^n = \frac{4}{h^2} \left( 1 - \frac{1}{2} \cos(\pi hk) - \frac{1}{2} \cos(\pi lh) \right)$$

$$v_{k,l}^n = \frac{1}{N-1} \begin{pmatrix} \sin(\pi hk) \sin(\pi lh) \\ \sin(\pi h(N-1)h) \sin(\pi lh) \\ \sin(\pi h(N-1)h) \sin(\pi l(N-1)h) \end{pmatrix}$$

Proof: similar, but

$$[A^n v_{k,l}^n]_{i,j} = \sum_{\alpha, \beta=1}^{(N-1)h^2} a_{\alpha, \beta} v_{\alpha, \beta}^n$$

Other discretization tips

$$u''(x_i) \approx \underbrace{\frac{-u(x_{i-1}) + 2u(x_i) - u(x_{i+1})}{h^2}}$$

approx. of second derivative

So in 2D

$$\epsilon u + \frac{\partial^2 u}{\partial x^2} + \epsilon \cdot \frac{\partial^2 u}{\partial y^2} = f(x) \text{ gives stencil}$$

$$\frac{1}{h^2} \begin{bmatrix} 0 & -\epsilon & 0 \\ -1 & 4 + \epsilon h^2 & -1 \\ 0 & -\epsilon & 0 \end{bmatrix}$$

## Chapter 4

The analysis on stability of a system does not take a solution method into account.

Let  $f$  be replaced by  $f + \Delta f$ . Then  $\exists u$ :

$$A(u + \Delta u) = f + \Delta f, \text{ where } \Delta u \text{ solves } A\Delta u = \Delta f.$$

Then  $\frac{\|\Delta u\|}{\|u\|} \leq \|A^{-1}\| \cdot \|A\| \frac{\|\Delta f\|}{\|f\|} = K(f(A)) \cdot \delta$

So if  $K(A)$  is large, a small perturbation can lead to a big difference in the solution vector:

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} u = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \Rightarrow u = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \Rightarrow u = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

$$\text{cond}(f(\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix})) = 3 \cdot 10^8$$

Similarly for perturbations in  $A$ . Due to rounding, such errors will occur, so it is important to have a small condition number to guarantee better approximations.

A way to deal with this is by scaling matrix rows such that they each row has similar  $\|\cdot\|_\infty$ -norm.

A direct way of solving  $Au=f$  is by writing  $A=LU$ .  $L$  is lower triangular and  $U$  upper, and this factorization is unique. Instead of solving  $Au=f$ , we solve first  $Ly=f$  and then  $Uu=y$  is solved for  $u$ . This solve is numerically stable, but the factorization requires some deeper analysis. These solves are easy, as matrices are triangular.

The LU-factorization can be computed by simply bringing  $A^{n \times n}$  into echelon form (in Dutch: "vegan")

Doing this requires  $n-1$  steps. The  $n-1$  steps form an  $L$ , the ~~adjusted~~ adjusted  $A$  is the required  $U$ . Formally, we have

$$M_{n-1} M_{n-2} \dots M_1 A = U \Rightarrow A = \underbrace{(M_{n-1} \dots M_1)}_L^{-1} U,$$

where  $M_k = I - \alpha_k e_k e_k^T = \begin{pmatrix} I & & & \\ & \ddots & & \\ & & I & \\ & & & -b_k/a_{kk} & \\ & & & -b_{k+1}/a_{k+1,k} & \ddots & \\ & & & & \ddots & I \end{pmatrix}$  where  $\alpha = (0, 0, \dots, 0, b_k/a_{kk})$

and  $a_{ij}^{(k)}$  represents the element on index  $i,j$  of the matrix  $A$  after  $k-1$  steps.

The inverse of  $M_k$  is  $M_k^{-1} = I + \alpha_k e_k e_k^T$

Note that this cannot deal with all zeroes on the diagonal.

Another way to prevent the errors having too much impact on the system is by applying partial pivoting instead.

Instead of just subtracting rows from one another, we first swap two rows such that the largest element in that column will be used as a pivot, so  $|a_{k,k}| \geq |a_{i,j}|$ .

Then  $M_{n-1} P_{n-1} \dots M_1 P_1 A = L U \Rightarrow PA = LU$   
where  $P = P_{n-1} \dots P_1$

So we require more flops, but the system is now stable! This is called partial pivoting, and can deal with zeroes.

One can go a step further and also allow column swaps as well. Then we have in step  $k$ :

$$A^{(k)} = \begin{pmatrix} A_{11}^{(k-1)} & A_{12}^{(k-1)} \\ 0 & A_{22}^{(k-1)} \end{pmatrix}, \text{ where we divided } A \text{ into three blocks}$$

Note that  $A_{11}^{(k-1)}$  is upper triangular.

NB Complete pivoting swaps rows and columns such that the largest element in the block  $A_{22}^{(k-1)}$  is now on element  $k,k$ .

### Special cases



- It can be proven using induction that diagonally dominant matrices already have the pivots in the right place, so normal Gaussian elimination is sufficient.
- If  $A$  is spd, the LU factorization is the same as the Cholesky decomposition  $A = C C^T$ , so we can decompose and store the matrices using half the time/memory/flops.

- For band matrices, normal Gaussian elimination guarantees that the matrices L and U inherit the band width q and p respectively.

Partial pivoting does not have this property, and thus should be avoided when possible.

- For tridiagonal systems, a direct ~~work~~ method is possible: let  $A = LDL^T$  and solve

$$a_{ii} = d_{ii}$$

$$a_{k,k-1} = e_{k-1} d_{k-1} \text{ for } 2 \leq k \leq n$$

$$a_{kk} = d_k + e_{k-1}^2, d_k = a_{kk} - e_{k-1} a_{k-1} \text{ for } 2 \leq k \leq n$$

- Smart ways exist to store general sparse matrices, such as the Yale format:

$$\begin{pmatrix} 10 & 20 & 0 & 0 & 0 & 0 \\ 0 & 30 & 0 & 40 & 0 & 0 \\ 0 & 0 & 50 & 60 & 70 & 0 \\ 0 & 0 & 0 & 0 & 0 & 80 \end{pmatrix}$$

elements  
 $AA = (10 \ 20 \dots \ 80)$   
 row switches  
 $IA = (1 \ 3 \ 5 \ 8)$   
 $\Rightarrow$  spots in row  
 $JA = (1 \ 2 \ 2 \ 4 \ 3 \ 5 \ 6)$

However, these methods are slow on large (sparse) matrices A

## Chapter 5

Basic iterative methods start with an initial solution  $u^0$  and keep improving that solution ~~improves~~ by iterating a certain process.

BIMs themselves are not fast enough for real life applications, but they do serve as the building blocks for more powerful methods.

The error after iteration  $k$  is defined as  $e^k = u^* - u^k$

The residual  $r^k = f - Au^k$

Note that knowing  $e^k$  requires  $u^*$  and that is very hard. Also, we know that  $Ae^k = r^k$ .

A splitting of  $A$  are two matrices  $M, N$  such that  $A = M - N$ .

$Au = f$  turns into  $Mu = Nu + f \Rightarrow u^k = M^{-1}Nu^k + M^{-1}f$   
 $= u^k + M^{-1}r^k$

Then:

$$e^k = u^* - u^{k+1} = u^* - M^{-1}Ae^k = \underbrace{(I - M^{-1}A)}_{B, \text{error propagation matrix}} e^k$$

$\therefore$  iteration matrix

$$\text{Similarly, } r^k = \dots = (I - AM^{-1})r^k$$

Two popular in BIMs are Jacobi and Gauss-Seidel.

Jacobi

$$A = \begin{pmatrix} -F & & \\ & D & \\ & & -E \end{pmatrix}$$

$$M_{\text{JAC}} = D, \quad N_{\text{JAC}} = E + F, \quad B_{\text{JAC}} = L + U$$

An update  $u_i^{k+1} = [f_i - \sum_{j=1, j \neq i}^n a_{ij} u_j^k] / a_{ii} \quad \forall i = 1, \dots, n$

is independent of other  $a_{ii}$  elements in  $u^{k+1}$ , so they can be calculated in a parallel fashion. (in other words, the parallel complexity is  $n$ )

Gauss-Seidel

$$M_{\text{GS}} = D - E, \quad N_{\text{GS}} = F, \quad B_{\text{GS}} = (I - L)^{-1} U$$

An update  $u_i^{k+1} = [f_i - \sum_{j=1}^{i-1} a_{ij} u_j^{k+1} - \sum_{j=i+1}^n a_{ij} u_j^k] / a_{ii}$

takes elements of the same iteration into account, and as such can converge more quickly than Jacobi. (parallel complexity: 1).

By changing the ordering to red/black or diagonal, the parallel complexity can be increased.

By calculating entire blocks of elements at once by means of dividing  $A$  into blocks and partitioning  $u$  and  $f$  likewise, solving the subsystems independently can also speed up the process.

Damped Jacobi with parameter  $\omega$  is defined as

$$u^{k+1} = \omega u^{k+1, \text{JAC}} + (1-\omega) u^k$$

This is the same as Richardson ( $M = \tau I, N = \tau F - A$ ) on the scaled system  $D^{-1}Au = D^{-1}f$  with  $\tau = \omega^{-1}$ .

Damping Gauss-Seidel results in a successive overrelaxation method:

$$u_i^{k+1} = (1-\omega) u_i^k + \omega u_i^{k+1, GS} \quad (\text{note that it's defined componentwise})$$

If  $A$  is spd, it is useful to have an spd splitting as well. Jacobi is spd, but for GS, we need to set  $M_1 = \frac{1}{\omega}D - E$ ,  $M_2 = \frac{1}{\omega}D - F$ , called Symmetric Successive Overrelaxation (SSSOR( $\omega$ ))

Note that  $e^k = Be^{k-1} = \dots = B^k e^0$ . That means that

$$\rho(B) < 1 \iff \{e^k\}_{k=1}^{\infty} \text{ converges.}$$

If  $\rho(B) \ll 1$ , convergence is fast.

Any matrix ~~bound~~<sup>norm</sup> is a bound for its  $\rho$ , so we need that

$\|B\|_1, \infty < 1$ ,  $\|B\|_2 < 1$ ,  $\|B\|_{\infty} < 1$ . For Jacobi:

$$\|B\|_1 = \|L + U\|_1 = \max_{1 \leq j \leq n} \sum_{i=1, i \neq j}^n \frac{|a_{ij}|}{|a_{ii}|} < 1 \quad (\text{max row sum})$$

$$\|B\|_{\infty} = \dots \quad (\text{max column sum})$$

$$\|B\|_2 = \dots = \sqrt{\sum_{i,j=1, i \neq j}^n \left(\frac{a_{ij}}{a_{ii}}\right)^2} \quad (\text{square sum criterium})$$

A splitting  $A = M - N$  is regular iff  $M$  non-singular,  $M^{-1} \geq 0$ ,  $N \geq 0$

If  $M-N$  is a regular splitting, then

$$(\rho(B) = \rho(M^{-1}N) < 1) \iff A \text{ is non-singular and } A^T \text{ nonnegative.}$$

If  $A$  is strongly diagonal row dominant, then JAC and GS converge.  $\exists \epsilon$

For SOR( $\omega$ ), we have that  $\rho(B_{SOR(\omega)}) \geq |1-\omega|$ , so we need to have  $0 < \omega < 2$ . If  $A$  is spd,  $0 < \omega < 2$  is also sufficient for guaranteeing convergence.

The rate of convergence/reduction rate is  $\mu \in (0, 1)$  s.t.

$$\|u^k - u^{k+1}\| \leq \mu \|u^{k+1} - u^{k-1}\| \quad \forall k.$$

For this  $\mu$ , we find  $\|e^k\|_2 \leq \frac{\mu^k}{1-\mu} \|e^0\|_2$

$$\boxed{\rho(B) \approx \mu}$$

So if  $\rho(B) \leq 1$ , we have convergence. If  $\rho(B) < 1$ , we have fast convergence

BIMs are such as Jac, GS and Jac( $\omega$ ) have  $\rho(B) \approx 1$ , and  $\rho(B) \rightarrow 1$  as  $h \rightarrow 0$  ( $\rho(B) = 1 - O(h^2)$ )

SSOR( $\omega$ ) has  $\rho(B) = 1 - O(h)$ , which is better.

A good choice of  $\omega^0$  can be made by solving the model on a coarser mesh.

The process can be monitored by tracking  $\|r^k\|_2$ .

A good stopping criterium takes into account the scale and does not depend on  $r^0$ :  $\frac{\|r^k\|}{\|r^0\|} \leq \epsilon$ .

Building better methods is based on:

- Splitting
- Defect correction
- Preconditioning

# Chapter 6

BIMs can deal with problems with low eigenvalues, as they will converge quickly. These eigenvalues (close to 0) have high frequency eigenmodes.

The eigenvalues close to 1 have slowly varying eigenmodes. That means they can be transformed to a smaller grid and then transformed back, without loss of information.

The starting vector is a combination of these modes. BIMs can act as a smoother for multigrid methods to do their magic.

A basic multigrid algorithm estimates the error and looks like this:

$$(\text{pre smoothing}) \quad u^{i+1/3} = S_h^m u_h^i + s_h$$

$$r_h^i = b_h - A_h u_h^{i+1/3}$$

$$r_{2h}^i = I_h^{2h} r_h^i$$

$$e_{2h}^i = (A_{2h})^{-1} r_{2h}^i$$

$$e_h^i = I_{2h}^h e_{2h}^i$$

$$u_h^{i+1/3} = u_h^{i+1/3} + e_h^i$$

$$(\text{post smoothing}) \quad u_h^i = S_h u_h^{i+1/3} + s_h$$

$I_h^{2h}$  transforms a vector in  $\mathbb{R}^{n-1} \rightarrow \mathbb{R}^{n/2-1}$

$I_{2h}^h$  transforms 'back' into  $\mathbb{R}^{n-1}$

$A_{2h} \approx$

So we estimate the exact error on a coarser grid, and bringing that to a fine grid.

Smoothers  $S$  can be any BTM iteration matrix.

$I_h^{zh}$  looks like makes each element a weighted average of corresponding elements in the finer grid.

$$1D: \begin{array}{c} h \\ \downarrow \\ z_h \end{array} \begin{bmatrix} 1 & 1/2 & 1/4 \\ 0 & 1 & 0 \\ 1/2 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1/4 & 1/2 & 1/4 \end{bmatrix}$$

2D: full weighting

$$\begin{bmatrix} 1/16 & 1 & 2 & 1 \\ 2 & 1/4 & 2 & 1 \\ 1 & 2 & 1 & 0 \end{bmatrix}$$

half weighting

$$\begin{bmatrix} 1/6 & 0 & 1 & 0 \\ 1 & 2 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$I_h^{zh}$  brings a smaller vector to a larger grid size:

$$1D: \begin{pmatrix} 1/2 & 0 & 0 \\ 1/2 & 1/2 & 0 \\ 0 & 1/2 & 0 \\ 0 & 1/2 & 1/2 \\ 0 & 0 & 1 \\ 0 & 0 & 1/2 \end{pmatrix}$$

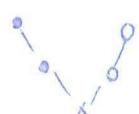
$$2D \quad \begin{bmatrix} 1/4 & 1 & 2 & 1 \\ 2 & 1/4 & 2 & 1 \\ 1 & 2 & 1 & 0 \end{bmatrix} \quad (\text{full weighting})$$

$A_{zh}$  can either be constructed like  $A_h$  but from fewer grid points, or be created by the Galerkin coarsening:

$$A_{zh} = I_h^{zh} A_h I_h^{zh}$$

The algorithm can be improved: multiple coarsening and prolongation steps can be taken.

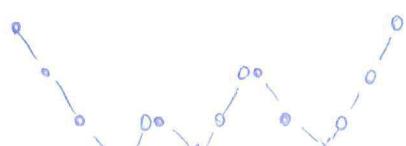
- presmooth    o postsmooth
- \ restriction    / prolongation



(3-grid) V-cycle



W-cycle



(4-grid) F-cycle

# Chapter 7

The next methods are based on the Krylov subspace:

$$K^k(A; r_0) = \text{span} \{r_0, Ar_0, \dots, A^{k-1}r_0\}$$

The first method we'll discuss is the conjugate gradient method. This method is based on the Chebyshev method, and computes

$$u^i \text{ such that } \|u - u^i\|_A = \min_{y \in K^k(A, r_0)} \|u - y\|_A$$

Theoretically, this has some nice properties:

$$1) \cdot \text{span} \{p^1, \dots, p^k\} = \text{span} \{r^0, \dots, r^{k-1}\} = K^k(A, r^0)$$

$$2) \cdot (r^j)^T r^i = 0 \quad 1 \leq i < j \leq k$$

$$\cdot (r^j)^T p^i = 0 \quad 1 \leq i \leq j \leq k$$

$$\cdot (p^j)^T A p^i = 0 \quad 1 \leq i < j \leq k$$

$p$  is search direction  
 $r$  residual

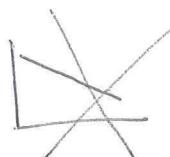
• 1,2:  $r^0, \dots, r^{k-1}$  is an orthogonal basis for  $K^k(A, r^0)$ .

• 1,2:  $r^0, \dots, r^{k-1}$  is an orthogonal basis for  $K^k(A, r^0)$ .

• CG is a finite method:  $K^N(A, r^0) = \mathbb{R}^N \Rightarrow u^N = u^*$ . In practice, this is not reached due to rounding errors, or  $N$  is too large to run  $N$  iterations.

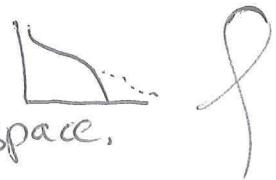
The rate of convergence  $\mu = \frac{\sqrt{k_2(A)} - 1}{\sqrt{k_2(A)} + 1}$

In practice, we have superlinear convergence:



This happens because eigenvectors belonging to big eigenvalues (almost) enter the Krylov subspace,

and thus no longer influence the convergence.



CG can only be applied to Spd matrices  $A$

9/12

By preconditioning the CG method, we can speed up the process: Instead of solving  $Au = b$ , we solve  $M^{-1}A u = M^{-1}b$  instead. We need to find an  $M^{-1}$  such that:

- the eigenvalues of  $M^{-1}A$  are clustered around 1
- $M^{-1}b$  can be obtained easily.

The first condition ensures that  $k_2(A) \leq k_2(M^{-1}A)$ , so it converges faster.

You can also write  $(P^{-1}AP^{-T})(P^T u) = P^{-1}b$ , and let  $M = PPT$ .

Some options are:

- $P$  is diagonal with  $p_{ii} = \sqrt{|a_{ii}|}$  (Diagonal scaling)
- $M^{-1} = M_{GS}^{-1}$  or  $M_{JAC}^{-1}$  (BIM)
- Incomplete Cholesky decomposition: IC6C(m).

The normal Cholesky decomposition always exists for ~~sym~~ spd matrices: ~~ATBTA~~  $A = LL^T$ .

The incomplete forces 0 elements of  $A$  to be 0 elements of  $L$  as well, saving the sparseness of the system. Cholesky decomposition is very stable.

IC6C(m) allows for  $m$  more  $m$  more bands around the diagonal to become non-zero.

MIC6C forces equal rowsums as well, performing better if  $a$  and/or  $b$  are slowly varying vectors.

For general matrices, it is impossible to obtain a method that has the following three properties

- ~~Stability~~  $\dim K^k(A; r^0)$
- Optimality (wrt  $\| \cdot \|_A$  norm)
- Short recurrences.

But for each combination of 2 properties, there are methods.

Opt + short recurrences:



CGNR (Conjugate Gradient normal residual)

Apply CG to  $\underbrace{ATAu}_\text{Spd} = A^T b$

However,  $K_2(ATA) = K_2(A)^2$ , which is usually large, so convergence is slow. Also  $A^T r$  is not always easily obtained.

A reliable improvement is Least Squares QR (LSQR), which

solves  $(\begin{smallmatrix} I & A \\ A^T & 0 \end{smallmatrix}) \begin{pmatrix} r \\ u \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix}$

Krylov + short recurrences



Bi-CG type methods have two series of residuals:

$r^e, \dots, r^{k-1}$  form a basis for  $K^k(A; r^0)$

$\hat{r}^0, \dots, \hat{r}^{k-1}$  form a basis for  $K^k(A, \hat{r}^0)$

~~The bases~~

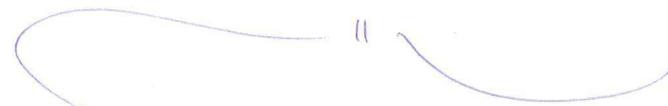
$r^k$  and  $\hat{r}^k$  are orthogonal bases for  $K^k(A; r^0)$  and  $K^k(A, \hat{r}^0)$  respectively.

This method has bad convergence if  $A$  is far from spd ( $A \approx -A^T$ ), doesn't guarantee optimality and requires  $A^T$ . It can also have unlucky (near) break downs.

There are two improvements:

- CGS (CG squared) is based on the fact that  $(r_i, \hat{r}_j) = (\underbrace{P_i(A) P_j(A)}_{\text{abuses this polynomial}} r^0, \hat{r}^0) = 0$  for  $i \neq j$  and converges faster than BiCG usually, but shows very irregular convergence behavior.
- BiCG STAB modifies the polynomials as well

Krylov + ~~whatever~~ optimality



The GMRES type methods require a lot of memory. To deal with this, restart or truncation (only using the last  $\ell$  vectors, for some  $\ell$ ) is required.

GMRES methods (generalized minimized residual) uses a different orthogonal basis for the Krylov subspace. To ensure orthogonality, all earlier vectors must be stored. Lucky breakdowns can occur.

Due to ~~reset~~, only local minimization can be guaranteed and superlinear convergence behavior is lost.

GCR (Generalized Conjugate residual) converges similarly to GMRES, but in some cases it breaks down. However, GCR allows for truncation, so if superlinear convergence is important, GCR is better than GMRES. In general, truncation is better than restarting.

Hybrid methods do exist. They mainly combine BiCG and GMRES type methods.

One method, IDR (Induced Dimension Reduction) has later been overshadowed by CGS and BiCG STAB.

GMRESR (GMRES recursive). This method uses a GCR approach on search directions calculated by GMRES. The search directions  $v_i$  are stored using a little trick:

GMRES(3):  $v_1, v_2, v_3$  restart  $v_1, v_2, v_3$  restart ...

GCR, trunc 3:  $v_1, v_2, v_3 [v_4, v_5, v_6]$  ...

↑↑↑  
↓

GMRESR with  
GMRES(3) as  
inner loop:

$\begin{matrix} \hat{v}_1 & \hat{v}_2 & \hat{v}_3 \\ \downarrow & & \\ v_1 & & \end{matrix}$        $\begin{matrix} \hat{v}_1 & \hat{v}_2 & \hat{v}_3 \\ \downarrow & & \\ v_2 & & \end{matrix}$  ...

So vital information does not get lost.

For GMRESR, truncation and restarting is possible as well.

Deciding which method to use depends on the expected amount of iterations ( $m_g$ ) and the ratio

$$f = \frac{\text{CPU time for prec. matvec}}{\text{CPU time for vector update}}$$

Big  $m_g$ , small  $f$ : BiCG STAB

Small  $m_g$ , big  $f$ : GMRES

In between : GMRESR

These methods can be preconditioned as well.

BLANK PAGE

BLANK PAGE

# Chapter 8



Due to lack of time, we only need to know the Power method.

True Eigenvalues give a lot of practical and theoretical information, so it is important to be able to calculate them. The power method can calculate (multiple) extreme eigenvalues.

Let  $q_0 \in \mathbb{C}$ . Repeat the following process:

$$z_k = A q_{k-1}$$

$$q_k = z_k / \|z_k\|_2$$

$$\lambda^{(n)} = q_{k-1}^T z_k$$

For example, if the matrix  $A$  is  $A = \begin{pmatrix} 10 & 0 \\ 0 & 2 \end{pmatrix}$  and  $q_0 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ ,

$$\text{then } z_0 = \begin{pmatrix} 10 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 10 \\ 2 \end{pmatrix}, \quad q_1 \approx \begin{pmatrix} 0.9806 \\ 0.1961 \end{pmatrix}, \quad \lambda^{(1)} = \begin{pmatrix} 1 & 1 \end{pmatrix} \begin{pmatrix} 10 \\ 2 \end{pmatrix} = 12$$

$$z_1 = \begin{pmatrix} 10 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 0.9806 \\ 0.1961 \end{pmatrix} = \begin{pmatrix} 9.8058 \\ 0.3922 \end{pmatrix}, \quad q_2 \approx \lambda^{(2)} = 9.6923$$

$$\lambda^{(3)} = 9.9872$$

$$\lambda^{(4)} = 9.9997$$

$$\lambda^{(5)} = 10.0000$$

This converges to the largest eigenvalue, as ~~all~~

$$\text{(let } |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|)$$

$$A^k q_0 = A^k (c_1 v_1 + \dots + c_n v_n) = c_1 \lambda_1^k v_1 + \dots + c_n \lambda_n^k v_n, \text{ so this will}$$

so  $\lambda_1$  will have the biggest impact.

$$12/12$$

Note that this fails if  $c_i = 0$  for your choice of  $q_0$ !

The convergence speed depends on  $|\frac{\lambda_2}{\lambda_1}|$ , and so applying the power method to  $A - cI$  depends on  $\left| \frac{\lambda_2^+ - c}{\lambda_1^* - c} \right|$ . \* means new largest eigenvalues So shifting can increase convergence speed.

Taking  $c \approx \lambda_1$  means that  $\lambda_{n-1}$  is now the largest eigenvalue so we can also calculate the smallest eigenvalue. (\*)

If  $|\lambda_1| = |\lambda_2|$ , a fix is possible.

~~Taking  $A^{-1}$  instead of  $A - cI$~~

$A^{-1}$  has eigenvalues  $\frac{1}{\lambda_i}$ , so applying the power method to  $A^{-1}$  will converge to  $\frac{1}{\lambda_1}$ .

Small eigenvalues are generally close to each other, so taking  $A^{-1}$  instead of the shift (\*) will probably lead to a faster convergence.

To compute  $z_k = A^{-1}q_{k-1}$ , one solves  $Az_k = q_{k-1}$ .

To calculate the largest  $p$  eigenvalues, let  $Q_0 \in \mathbb{C}^{n \times p}$  an orthogonal matrix and do

- $z_k = A Q_{k-1}$
- orthonormalize columns of  $z_k$  s.t.  $Q_k R_k = z_k$ ,  $R_k \in \mathbb{R}^{k \times k}$  upp. triag. and  $Q_k^T Q_k = I$