

$$A \vec{x} = \vec{b}$$

**Direct methods** (Gauss elimination, LU decomposition, Cholesky decomposition). Useful for dense or banded (e.g., tridiagonal) matrices.

**Iterative methods.** Useful for large sparse matrices.

In general,  $\vec{x}^{(k+1)} = \mathbf{B}^{(k)} \vec{x}^{(k)} + \vec{g}^{(k)}$   $\vec{x} = \mathbf{B}^{(k)} \vec{x} + \vec{g}^{(k)}$

Hopefully,  $\vec{x}^{(k)} \rightarrow \vec{x}, k \rightarrow \infty$

For the methods we have considered,  $A = A_1 + A_2, \quad A_1 \vec{x} = \vec{b} - A_2 \vec{x}$

$$\vec{x} = A_1^{-1} \vec{b} - A_1^{-1} A_2 \vec{x} \Rightarrow \vec{x}^{(k+1)} = A_1^{-1} \vec{b} - A_1^{-1} A_2 \vec{x}^{(k)}$$

$\mathbf{B} = -A_1^{-1} A_2; \quad \vec{g} = A_1^{-1} \vec{b}$  Both are  $k$ -independent.

**Stationary iterative methods.**

$$\vec{x}^{(k+1)} = \mathbf{B} \vec{x}^{(k)} + \vec{g}$$

Consider the error  $\vec{e}^{(k)} = \vec{x}^{(k)} - \vec{x}$ .

$$\vec{e}^{(k+1)} = \vec{x}^{(k+1)} - \vec{x} = \mathbf{B}^{(k)} \vec{x}^{(k)} + \vec{g}^{(k)} - \mathbf{B}^{(k)} \vec{x} - \vec{g}^{(k)} = \mathbf{B}^{(k)} \vec{e}^{(k)}$$

2

$$\vec{e}^{(k+1)} = \mathbf{B}^{(k)} \vec{e}^{(k)}$$

$$\|\vec{e}^{(k+1)}\| = \|\mathbf{B}^{(k)} \vec{e}^{(k)}\| \leq |\lambda|_{\max} \|\vec{e}^{(k)}\| = \rho(\mathbf{B}) \|\vec{e}^{(k)}\|$$

$\rho(\mathbf{B})$  is called the spectral radius of matrix  $\mathbf{B}$ .  $\|\vec{e}^{(k)}\| \leq [\rho(\mathbf{B})]^k \|\vec{e}^{(0)}\|$

For convergence, need  $\rho(\mathbf{B}) < 1$ . If we want the final error to be at most  $\epsilon$  times the initial error, the minimum number of iterations satisfies  $[\rho(\mathbf{B})]^k \leq \epsilon$

The smaller  $\rho(\mathbf{B})$ , the fewer iterations are needed.

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2 \quad \vec{x}^{(k+1)} = \mathbf{A}_1^{-1} \vec{b} - \mathbf{A}_1^{-1} \mathbf{A}_2 \vec{x}^{(k)} \quad \mathbf{B} = -\mathbf{A}_1^{-1} \mathbf{A}_2$$

For the Jacobi method,  $\mathbf{A}_1 = \mathbf{D}$ ,  $\mathbf{A}_2 = \mathbf{L} + \mathbf{U}$ ,  $\mathbf{B} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})$

$$b_{ij} = \begin{cases} -\frac{a_{ij}}{a_{ii}}, & i \neq j, \\ 0, & i = j. \end{cases} \quad x_i^{(k+1)} = -\frac{1}{a_{ii}} \sum_{j \neq i} a_{ij} x_j^{(k)} + \frac{1}{a_{ii}} b_i$$

Consider an arbitrary eigenvalue of  $\mathbf{B}$ ,  $\lambda$ , and the

corresponding eigenvector  $\vec{y}$ . Assume, without loss of generality, that  $y_i$

is the largest component by absolute value and  $|y_i| = 1$ . Then

$$|\lambda| = |\lambda| |x_i| = \left| \sum_{j \neq i} b_{ij} x_j \right| \leq \sum_{j \neq i} |b_{ij}| = \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|}$$

3  $|\lambda| \leq \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|}$  To ensure convergence, need  $|\lambda| < 1$

This will definitely be the case, if  $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$ . **Diagonal dominance.**

**Gauss-Seidel.** 
$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left( -\sum_{j < i} a_{ij} x_j^{(k+1)} - \sum_{j > i} a_{ij} x_j^{(k)} + \frac{1}{a_{ii}} b_i \right)$$

$$\vec{x}^{(k+1)} = \mathbf{A}_1^{-1} \vec{b} - \mathbf{A}_1^{-1} \mathbf{A}_2 \vec{x}^{(k)} \quad \mathbf{A}_1 = \mathbf{D} + \mathbf{L} \quad \mathbf{A}_2 = \mathbf{U}$$

$\vec{x}^{(k+1)} = \mathbf{B} \vec{x}^{(k)} + \vec{g}$  It can be proved that for symmetric  $\mathbf{A}$  iterations converge **if and only if**  $\mathbf{A}$  is also positive definite (all eigenvalues  $> 0$ ).

E.g., for tridiagonal matrices exactly twice as fast as Jacobi.

Note the result depends on the ordering of the components.

**Successive overrelaxation (SOR)**

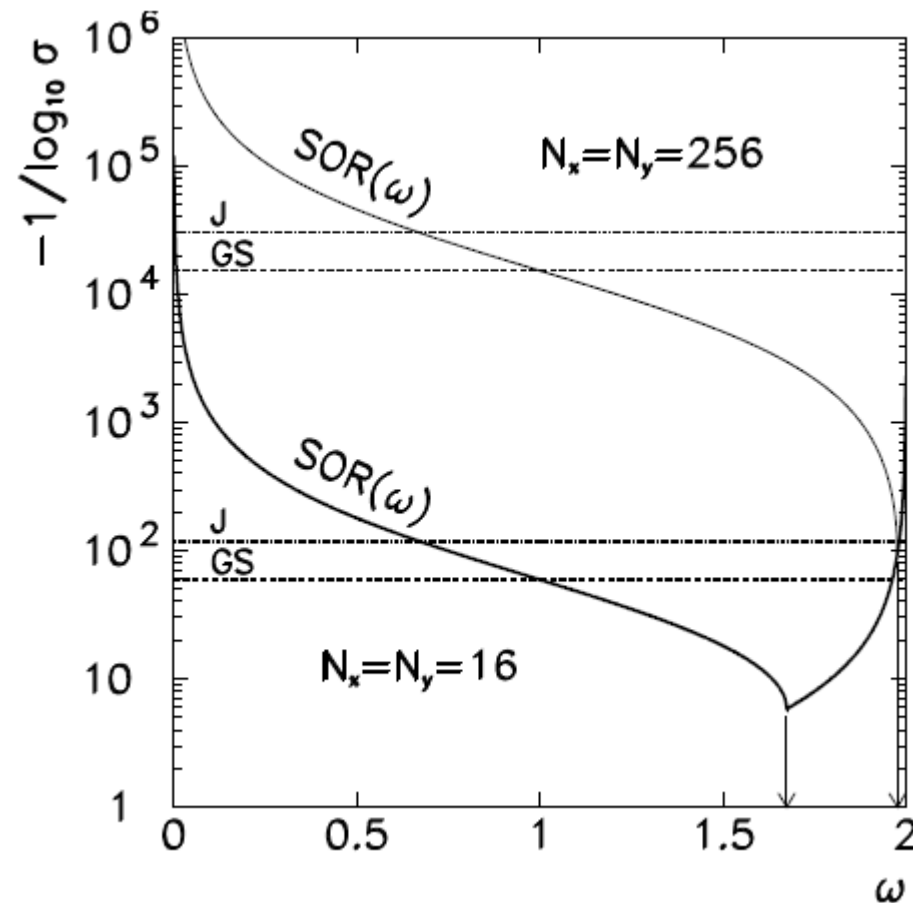
$$x_i^{(n+1)} = (1 - \omega) x_i^{(n)} + \frac{\omega}{a_{ii}} \left( -\sum_{j < i} a_{ij} x_j^{(n+1)} - \sum_{j > i} a_{ij} x_j^{(n)} + \frac{1}{a_{ii}} b_i \right)$$

Converges whenever Gauss-Seidel converges for  $0 < \omega < 2$ .

**Over**relaxation for  $1 < \omega < 2$ .

4

Inverse of the convergence rate for matrices arising from discretization of the Poisson equation in 2D.



Optimal convergence rate for a particular  $\omega$ . The larger the matrix, the larger the gain. Optimal  $\omega$  is known for matrices arising in some problems, but not in general. There are ways to estimate it, but not easy.

## 5

## Non-stationary approaches

$$\vec{x}^{(k+1)} = \mathbf{B}^{(k)} \vec{x}^{(k)} + \vec{g}^{(k)}$$

## Conjugate gradient method

Only works for symmetric positive definite matrices (can be generalized to Hermitian matrices)

Solving a linear system as a **minimization problem**: consider

$$F(\vec{x}) = \frac{1}{2} \vec{x} \cdot \mathbf{A} \cdot \vec{x} - \vec{b} \cdot \vec{x} = \frac{1}{2} a_{ij} x_i x_j - b_i x_i$$

$$\frac{\partial F}{\partial x_i} = \frac{1}{2} (a_{ij} x_j + a_{ji} x_j) - b_i = a_{ij} x_j - b_i = 0$$

Minimum, since the matrix is positive definite.

Start at  $\vec{x}^{(0)}$ . Move in the steepest descent direction

$$\vec{p}^{(0)} = -\nabla F = \vec{b} - \mathbf{A} \vec{x} = \vec{r}^{(0)}. \quad \vec{x}^{(1)} = \vec{x}^{(0)} + \alpha_0 \vec{p}^{(0)}.$$

$$F(\vec{x}^{(1)}) = \frac{1}{2} (\vec{x}^{(0)} + \alpha_0 \vec{p}^{(0)}) \cdot \mathbf{A} \cdot (\vec{x}^{(0)} + \alpha_0 \vec{p}^{(0)}) - \vec{b} \cdot (\vec{x}^{(0)} + \alpha_0 \vec{p}^{(0)})$$

$$\frac{dF}{d\alpha_0} = \alpha_0 \vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)} + \vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{x}^{(0)} - \vec{b} \cdot \vec{p}^{(0)} = \alpha_0 \vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)} - \vec{r}^{(0)} \cdot \vec{p}^{(0)} = 0$$

6

$$\frac{dF}{d\alpha_0} = \alpha_0 \vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)} - \vec{r}^{(0)} \cdot \vec{p}^{(0)} = 0 \Rightarrow \alpha_0 = \frac{\vec{r}^{(0)} \cdot \vec{p}^{(0)}}{\vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)}} = \frac{\vec{r}^{(0)} \cdot \vec{r}^{(0)}}{\vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)}}$$

$$\vec{x}^{(1)} = \vec{x}^{(0)} + \alpha_0 \vec{p}^{(0)} = \vec{x}^{(0)} + \frac{\vec{r}^{(0)} \cdot \vec{r}^{(0)}}{\vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)}} \vec{p}^{(0)}$$

$$\vec{r}^{(1)} = \vec{b} - \mathbf{A} \vec{x}^{(1)} = \vec{b} - \mathbf{A} \vec{x}^{(0)} - \alpha_0 \mathbf{A} \vec{p}^{(0)} = \vec{r}^{(0)} - \alpha_0 \mathbf{A} \vec{p}^{(0)} = \vec{r}^{(0)} - \frac{\vec{r}^{(0)} \cdot \vec{r}^{(0)}}{\vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)}} \mathbf{A} \vec{p}^{(0)}$$

$$\vec{r}^{(1)} \cdot \vec{r}^{(0)} = \vec{r}^{(0)} \cdot \vec{r}^{(0)} - \frac{\vec{r}^{(0)} \cdot \vec{r}^{(0)}}{\vec{p}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)}} \vec{r}^{(0)} \cdot \mathbf{A} \cdot \vec{p}^{(0)} = 0$$

Quite logical: the gradient at the minimum should be orthogonal to the trajectory. Could have used the orthogonality condition instead of the minimization condition to determine  $\alpha_0$ .

What do we do now? Steepest descent again? Not optimal.

$$\vec{p}^{(1)} = \vec{r}^{(1)} + \beta_1 \vec{p}^{(0)} \quad \text{Choose } \beta_1 \text{ so that reach minimum in the}$$

space spanned by  $\vec{r}^{(0)}$  and  $\vec{r}^{(1)}$ .  $\vec{r}^{(2)}$  should be  $\perp \vec{r}^{(0)}, \vec{r}^{(1)}$

7

$$\vec{p}^{(1)} = \vec{r}^{(1)} + \beta_1 \vec{p}^{(0)}$$

$$\vec{x}^{(2)} = \vec{x}^{(1)} + \alpha_1 \vec{p}^{(1)} \quad \vec{r}^{(2)} = \vec{b} - A \vec{x}^{(2)} = \vec{b} - A \vec{x}^{(1)} - \alpha_1 A \vec{p}^{(1)} = \vec{r}^{(1)} - \alpha_1 A \vec{p}^{(1)}$$

$$\vec{r}^{(2)} \cdot \vec{r}^{(0)} = -\alpha_1 \vec{r}^{(0)} \cdot A \cdot \vec{p}^{(1)} = -\alpha_1 \vec{p}^{(0)} \cdot A \cdot \vec{p}^{(1)} = 0$$

So the requirement is

$$\vec{p}^{(0)} \cdot A \cdot \vec{p}^{(1)} = 0 \Rightarrow \vec{p}^{(0)} \cdot A \cdot (\vec{r}^{(1)} + \beta_1 \vec{p}^{(0)}) = 0 \Rightarrow \beta_1 = -\frac{\vec{p}^{(0)} \cdot A \cdot \vec{r}^{(1)}}{\vec{p}^{(0)} \cdot A \cdot \vec{p}^{(0)}}$$

New direction is **conjugate** to the previous one – hence the name.

Next iteration: want to reach the minimum in space spanned by  $\vec{r}^{(0)}, \vec{r}^{(1)}, \vec{r}^{(2)}$ .

Residual  $\vec{r}^{(3)}$  will need to be orthogonal to  $\vec{r}^{(0)}, \vec{r}^{(1)}, \vec{r}^{(2)}$ . It would seem

reasonable to assume that we would now need to choose

$$\vec{p}^{(2)} = \vec{r}^{(2)} + \beta_2 \vec{p}^{(1)} + \gamma_2 \vec{p}^{(0)} \quad \vec{x}^{(3)} = \vec{x}^{(2)} + \alpha_2 \vec{p}^{(2)} \Rightarrow \vec{r}^{(3)} = \vec{r}^{(2)} - \alpha_2 A \vec{p}^{(2)}$$

This way, we would have 3 constants to adjust to satisfy 3 orthogonality conditions. And then we'd have 4 constants, 5, 6, ...

A remarkable thing is that this is **not the case!**  $\gamma_2 = 0!!!$

$$8 \quad \vec{p}^{(2)} = \vec{r}^{(2)} + \beta_2 \vec{p}^{(1)} \quad \vec{x}^{(3)} = \vec{x}^{(2)} + \alpha_2 \vec{p}^{(2)} \Rightarrow \vec{r}^{(3)} = \vec{r}^{(2)} - \alpha_2 \mathbf{A} \vec{p}^{(2)}$$

Adjusting  $\alpha_2$  and  $\beta_2$  to satisfy  $\vec{r}^{(3)} \cdot \vec{r}^{(2)} = 0$  and  $\vec{p}^{(2)} \cdot \mathbf{A} \cdot \vec{p}^{(1)}$ , we

automatically satisfy  $\vec{r}^{(3)} \cdot \vec{r}^{(1)} = 0$ ,  $\vec{r}^{(3)} \cdot \vec{r}^{(0)} = 0$  and  $\vec{p}^{(2)} \cdot \mathbf{A} \cdot \vec{p}^{(0)}$ .

Moreover, this is true for all subsequent iterations as well: if

$$\vec{r}^{(j)} \cdot \vec{r}^{(i)} = 0 \quad \forall i, j < k \quad (i \neq j) \quad \text{and} \quad \vec{p}^{(j)} \cdot \mathbf{A} \cdot \vec{p}^{(i)} = 0 \quad \forall i, j < k-1 \quad (i \neq j),$$

then, defining  $\vec{p}^{(k-1)} = \vec{r}^{(k-1)} + \beta_{k-1} \vec{p}^{(k-2)}$  and  $\vec{x}^{(k)} = \vec{x}^{(k-1)} + \alpha_{k-1} \vec{p}^{(k-1)}$ ,

it is sufficient to ensure  $\vec{r}^{(k)} \cdot \vec{r}^{(k-1)} = 0$  and  $\vec{p}^{(k-1)} \cdot \mathbf{A} \cdot \vec{p}^{(k-2)} = 0$  for

$\vec{r}^{(k)} \cdot \vec{r}^{(j)} = 0$ ,  $j < k-1$  and  $(\vec{p}^{(k-1)})^T \mathbf{A} \vec{p}^{(j)} = 0$ ,  $j < k-2$  to be satisfied

automatically. Indeed, using  $\vec{r}^{(k)} = \vec{b} - \mathbf{A}(\vec{x}^{(k-1)} + \alpha_{k-1} \vec{p}^{(k-1)}) = \vec{r}^{(k-1)} - \alpha_{k-1} \mathbf{A} \vec{p}^{(k-1)}$ ,

$$\vec{r}^{(k)} \cdot \vec{r}^{(j)} = (\vec{r}^{(k-1)} - \alpha_{k-1} \mathbf{A} \vec{p}^{(k-1)}) \cdot \vec{r}^{(j)} = -\alpha_{k-1} \vec{p}^{(k-1)} \cdot \mathbf{A} \cdot \vec{r}^{(j)}$$

$$= -\alpha_{k-1} \vec{p}^{(k-1)} \cdot \mathbf{A} \cdot (\vec{p}^{(j)} - \beta_j \vec{p}_{j-1}) = 0, \quad j < k-1,$$

$$\vec{p}^{(k-1)} \cdot \mathbf{A} \cdot \vec{p}^{(j)} = (\vec{r}^{(k-1)} + \beta_{k-1} \vec{p}^{(k-2)}) \cdot \mathbf{A} \cdot \vec{p}^{(j)} = r^{(k-1)} \cdot \mathbf{A} \cdot \vec{p}^{(j)}$$

$$= r^{(k-1)} \cdot (\vec{r}^{(j)} - \vec{r}^{(j+1)}) / \alpha_j = 0, \quad j < k-2.$$



9

$$\vec{p}^{(k-1)} = \vec{r}^{(k-1)} + \beta_{k-1} \vec{p}^{(k-2)}, \quad \vec{r}^{(k)} = \vec{r}^{(k-1)} - \alpha_{k-1} \mathbf{A} \vec{p}^{(k-1)}$$

$$\vec{r}^{(k-1)} \cdot \vec{r}^{(k)} = \vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)} - \alpha_{k-1} \vec{r}^{(k-1)} \cdot \mathbf{A} \cdot \vec{p}^{(k-1)}$$

$$= \vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)} - \alpha_{k-1} (\vec{p}^{(k-1)} - \beta_{k-1} \vec{p}^{(k-2)}) \cdot \mathbf{A} \cdot \vec{p}^{(k-1)}$$

$$= \vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)} - \alpha_{k-1} \vec{p}^{(k-1)} \cdot \mathbf{A} \cdot \vec{p}^{(k-1)} = 0$$

$$\alpha_{k-1} = \frac{\vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)}}{\vec{p}^{(k-1)} \cdot \mathbf{A} \cdot \vec{p}^{(k-1)}}$$

$$\vec{p}^{(k-2)} \cdot \mathbf{A} \cdot \vec{p}^{(k-1)} = \vec{p}^{(k-2)} \cdot \mathbf{A} \cdot (\vec{r}^{(k-1)} + \beta_{k-1} \vec{p}^{(k-2)}) = 0 \quad \beta_{k-1} = -\frac{\vec{p}^{(k-2)} \cdot \mathbf{A} \cdot \vec{r}^{(k-1)}}{\vec{p}^{(k-2)} \cdot \mathbf{A} \cdot \vec{p}^{(k-2)}}$$

$$\begin{aligned} \vec{p}^{(k-2)} \cdot \mathbf{A} \cdot \vec{r}^{(k-1)} &= \vec{r}^{(k-1)} \cdot \mathbf{A} \cdot \vec{p}^{(k-2)} = \vec{r}^{(k-1)} \cdot (\vec{r}^{(k-2)} - \vec{r}^{(k-1)}) / \alpha_{k-2} \\ &= -\vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)} / \alpha_{k-2} \end{aligned}$$

$$\vec{p}^{(k-2)} \cdot \mathbf{A} \cdot \vec{p}^{(k-2)} = (\vec{r}^{(k-2)} + \beta_{k-2} \vec{p}^{(k-3)}) \cdot \mathbf{A} \cdot \vec{p}^{(k-2)} = \vec{r}^{(k-2)} \cdot \mathbf{A} \cdot \vec{p}^{(k-2)}$$

$$= \vec{r}^{(k-2)} \cdot (\vec{r}^{(k-2)} - \vec{r}^{(k-1)}) / \alpha_{k-2} = \vec{r}^{(k-2)} \cdot \vec{r}^{(k-2)} / \alpha_{k-2}$$

$$\beta_{k-1} = \frac{\vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)}}{\vec{r}^{(k-2)} \cdot \vec{r}^{(k-2)}}$$

## Complete algorithm:

Start with initial guess  $\vec{x}^{(0)}$ .

Calculate the residual  $\vec{r}^{(0)} = \vec{b} - \mathbf{A} \vec{x}^{(0)}$ .

$k = 0$

while  $\|\vec{r}^{(k)}\| > \text{tolerance} \times \|\vec{b}\|$

if  $k = 0$ ,  $\vec{p}^{(0)} = \vec{r}^{(0)}$

else

$$\beta_k = \frac{\vec{r}^{(k)} \cdot \vec{r}^{(k)}}{\vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)}}$$

$$\vec{p}^{(k)} = \vec{r}^{(k)} + \beta_k \vec{p}^{(k-1)}$$

end

$$\alpha_k = \frac{\vec{r}^{(k)} \cdot \vec{r}^{(k)}}{\vec{p}^{(k)} \cdot \mathbf{A} \cdot \vec{p}^{(k)}}$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \vec{p}^{(k)},$$

$$\vec{r}^{(k+1)} = \vec{r}^{(k)} - \alpha_k \mathbf{A} \vec{p}^{(k)},$$

$k = k + 1$

end

After every iteration, the residual gets restricted to space of a smaller dimensionality.

Should vanish completely after N iterations (actually, as many as there are distinct eigenvalues). In this sense, it should behave like an exact method.

Unfortunately, only true in exact arithmetic.

But is useful as an iterative method.

**11** Let's look again at the residuals.

$$\vec{r}^{(1)} = \vec{r}^{(0)} - \alpha_0 \mathbf{A} \vec{p}^{(0)} = \vec{r}^{(0)} - \alpha_0 \mathbf{A} \vec{r}^{(0)} \quad \text{A linear combination of } \vec{r}^{(0)} \text{ and } \mathbf{A} \vec{r}^{(0)}$$

Space spanned by  $\vec{r}^{(0)}$  and  $\vec{r}^{(1)}$  is also space spanned by  $\vec{r}^{(0)}$  and  $\mathbf{A} \vec{r}^{(0)}$

$$\vec{p}^{(1)} = \vec{r}^{(1)} + \beta_1 \vec{p}^{(0)} = \vec{r}^{(1)} + \beta_1 \vec{r}^{(0)}$$

$$\begin{aligned} \vec{r}^{(2)} &= \vec{r}^{(1)} - \alpha_1 \mathbf{A} \vec{p}^{(1)} = \vec{r}^{(0)} - \alpha_0 \mathbf{A} \vec{r}^{(0)} - \alpha_1 \mathbf{A} \vec{r}^{(1)} - \alpha_1 \beta_1 \mathbf{A} \vec{r}^{(0)} \\ &= \vec{r}^{(0)} - \alpha_0 \mathbf{A} \vec{r}^{(0)} - \alpha_1 \mathbf{A} \vec{r}^{(0)} + \alpha_0 \alpha_1 \mathbf{A}^2 \vec{r}^{(0)} - \alpha_1 \beta_1 \mathbf{A} \vec{r}^{(0)} \end{aligned}$$

A linear combination of  $\vec{r}^{(0)}$ ,  $\mathbf{A} \vec{r}^{(0)}$ ,  $\mathbf{A}^2 \vec{r}^{(0)}$

Linear subspace spanned by  $\vec{u}$ ,  $\mathbf{B} \vec{u}$ ,  $\mathbf{B}^2 \vec{u}$ , ...,  $\mathbf{B}^{s-1} \vec{u}$  is called order-s Krylov subspace.

In the conjugate gradient method by construction after s iterations the residual vanishes in order-s Krylov subspace. Such methods are called Krylov subspace methods.

## Convergence rate of the conjugate gradient method

$$\|\vec{x} - \vec{x}^{(k)}\|_{\mathbf{A}} \leq 2 \|\vec{x} - \vec{x}^{(0)}\|_{\mathbf{A}} \left( \frac{\sqrt{K(\mathbf{A})} - 1}{\sqrt{K(\mathbf{A})} + 1} \right)^k,$$

where  $\|\vec{w}\|_{\mathbf{A}} = \sqrt{\vec{w} \cdot \mathbf{A} \cdot \vec{w}}$  and  $K(\mathbf{A})$  is the spectral condition number (the ratio of the largest and smallest eigenvalues). Slow convergence if this ratio is large (ill-conditioned matrix). Sometimes not as bad, e.g., when the extreme eigenvalues are well-separated from the rest. But still, very often an issue, since many matrices are ill-conditioned.

Consider a chain of  $N$  beads connected with springs;  $x_i$  are displacements from equilibrium. Assume all masses are 1.

$\ddot{\vec{x}} = -\mathbf{D} \vec{x}$       Eigenvalues of  $\mathbf{D}$  are squares of the frequencies of the modes.

$$\mathbf{D} \vec{x} = \vec{F} \quad \omega \sim c/\lambda \quad K(\mathbf{D}) \sim N^2$$

**Preconditioning**

$$A \vec{x} = \vec{b} \quad C^{-1} A C^{-1} C \vec{x} = C^{-1} \vec{b} \quad \tilde{A} \tilde{x} = \tilde{b},$$

where  $\tilde{A} = C^{-1} A C^{-1}$ ,  $\tilde{x} = C \vec{x}$ ,  $\tilde{b} = C^{-1} \vec{b}$ .

If  $C$  is symmetric and positive definite, the same is true about  $C^{-1} A C^{-1}$ , so we can apply conjugate gradient to it. Just add tildes.

$$\vec{r} = \tilde{b} - \tilde{A} \tilde{x} = C^{-1} \vec{b} - C^{-1} A C^{-1} C \vec{x} = C^{-1} \vec{r}$$

Let us define  $\tilde{\pi}^{(k)}$  such that  $\tilde{p}^{(k)} = C \tilde{\pi}^{(k)}$ ,  $M = C^2$ ,  $\tilde{z}^{(k)}$  such that  $M \tilde{z}^{(k)} = \vec{r}^{(k)}$

$$\tilde{p}^{(0)} = \vec{r}^{(0)} \Rightarrow C \tilde{\pi}^{(0)} = C^{-1} \vec{r}^{(0)} \Rightarrow \tilde{\pi}^{(0)} = M^{-1} \vec{r}^{(0)} = \tilde{z}^{(0)}$$

$$\beta_k = \frac{\vec{r}^{(k)} \cdot \vec{r}^{(k)}}{\vec{r}^{(k-1)} \cdot \vec{r}^{(k-1)}} = \frac{\vec{r}^{(k)} C^{-1} C^{-1} \vec{r}^{(k)}}{\vec{r}^{(k-1)} C^{-1} C^{-1} \vec{r}^{(k-1)}} = \frac{\vec{r}^{(k)} \tilde{z}^{(k)}}{\vec{r}^{(k-1)} \tilde{z}^{(k-1)}}$$

$$\tilde{p}^{(k)} = \vec{r}^{(k)} + \beta_k \tilde{p}^{(k-1)} \Rightarrow C \tilde{\pi}^{(k)} = C^{-1} \vec{r}^{(k)} + \beta_k C \tilde{\pi}^{(k-1)} \Rightarrow \tilde{\pi}^{(k)} = \tilde{z}^{(k)} + \beta_k \tilde{\pi}^{(k-1)}$$

$$\alpha_k = \frac{\vec{r}^{(k)} \cdot \vec{r}^{(k)}}{\tilde{p}^{(k)} \cdot \tilde{A} \tilde{p}^{(k)}} = \frac{\vec{r}^{(k)} C^{-1} C^{-1} \vec{r}^{(k)}}{\tilde{\pi}^{(k)} C C^{-1} A C^{-1} C \tilde{\pi}^{(k)}} = \frac{\vec{r}^{(k)} \cdot \tilde{z}^{(k)}}{\tilde{\pi}^{(k)} A \tilde{\pi}^{(k)}}$$

$$14 \quad \vec{\tilde{x}}^{(k+1)} = \vec{\tilde{x}}^{(k)} + \alpha_k \vec{\tilde{p}}^{(k)} \Rightarrow \mathbf{C} \vec{\tilde{x}}^{(k+1)} = \mathbf{C} \vec{\tilde{x}}^{(k)} + \alpha_k \mathbf{C} \vec{\tilde{\pi}}^{(k)} \Rightarrow \vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \vec{\pi}^{(k)}$$

$$\begin{aligned} \vec{\tilde{r}}^{(k+1)} = \vec{\tilde{r}}^{(k)} - \alpha_k \tilde{\mathbf{A}} \vec{\tilde{p}}^{(k)} &\Rightarrow \mathbf{C}^{-1} \vec{\tilde{r}}^{(k+1)} = \mathbf{C}^{-1} \vec{\tilde{r}}^{(k)} - \alpha_k \mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-1} \mathbf{C} \vec{\tilde{\pi}}^{(k)} \\ &\Rightarrow \vec{r}^{(k+1)} = \vec{r}^{(k)} - \alpha_k \mathbf{A} \vec{\pi}^{(k)} \end{aligned}$$

Preconditioned algorithm:

Start with initial guess  $\vec{x}^{(0)}$ .

Calculate the residual  $\vec{r}^{(0)} = \vec{b} - \mathbf{A} \vec{x}^{(0)}$ .

$k = 0$

while  $\|\vec{r}^{(k)}\| > \text{tolerance} \times \|\vec{b}\|$

**Solve**  $\mathbf{M} \vec{z}^{(k)} = \vec{r}^{(k)}$

if  $k = 0$ ,  $\vec{p}^{(0)} = \vec{z}^{(0)}$

else

$$\beta_k = \frac{\vec{r}^{(k)} \cdot \vec{z}^{(k)}}{\vec{r}^{(k-1)} \cdot \vec{z}^{(k-1)}}$$

$$\vec{\pi}^{(k)} = \vec{z}^{(k)} + \beta_k \vec{\pi}^{(k-1)}$$

end

$$\alpha_k = \frac{\vec{r}^{(k)} \cdot \vec{z}^{(k)}}{\vec{\pi}^{(k)} \mathbf{A} \vec{\pi}^{(k)}}$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \vec{\pi}^{(k)}$$

$$\vec{r}^{(k+1)} = \vec{r}^{(k)} - \alpha_k \mathbf{A} \vec{\pi}^{(k)}$$

$k = k + 1$

end

15 By construction, equivalent to solving an equation with matrix  $\mathbf{C}^{-1}\mathbf{A}\mathbf{C}^{-1}$ .

Need to choose  $\mathbf{C}$  such that  $\mathbf{C}^{-1}\mathbf{A}\mathbf{C}^{-1}$  is better conditioned. Note  $\mathbf{C}$  does not enter explicitly, only  $\mathbf{M}=\mathbf{C}^2$ .

Need to solve a system with matrix  $\mathbf{M}$  at every step. Obviously, should be easy to solve.

There are many different preconditioners used in practice differing by the amount of work needed to set it up and cost per iteration.

Simplest choice: Jacobi preconditioner – diagonal matrix

$$m_{ij} = \begin{cases} a_{ii}, & i=j \\ 0 & \text{otherwise} \end{cases}$$

Block Jacobi – a block-diagonal matrix. E.g., multiple variables per node, or spatially, or according to division of variables between processors.

Incomplete factorization preconditioners. E.g., Cholesky decomposition keeping zero all matrix elements equal to zero in  $\mathbf{A}$ .

## 16 Indefinite and/or non-symmetric systems

Conjugate gradient is not directly applicable. For non-symmetric systems, there is no short recurrence where only a few recent directions need to be used. For symmetric, but indefinite systems, short recurrences are possible, but the solution no longer corresponds to a minimum. No universally good approach.

Simplest approach: conjugate gradient for normal equations.

$$\mathbf{A} \vec{x} = \vec{b} \quad \text{Multiply by } \mathbf{A}^T: \quad \mathbf{A}^T \mathbf{A} \vec{x} = \mathbf{A}^T \vec{b}$$

Matrix  $\mathbf{A}^T \mathbf{A}$  is both symmetric and positive definite (unless  $\mathbf{A}$  is singular). CG can be applied to it.

Problem: the condition number is the square of that of matrix  $\mathbf{A}$ , so convergence can be rather slow.



**17 MINRES and SYMMLQ:** for symmetric but indefinite matrices. Minimizes the norm of the residual. A sequence of orthogonal vectors in the Krylov subspace can still be generated easily.

**GMRES:** a similar approach, but for non-symmetric systems. No short recurrences, needs to store all vectors generated so far and orthogonalize with respect to all of them. In practice, restarts are done. Reliable but costly in time and memory if many iterations are required.

**Biconjugate gradient (BiCG):** generates two mutually orthogonal sequences of vectors.

$$\begin{aligned} \vec{r}^{(i)} &= \vec{r}^{(i-1)} - \alpha_i A \vec{p}^{(i)}, & \vec{r}'^{(i)} &= \vec{r}'^{(i-1)} - \alpha_i A \vec{p}'^{(i)}, \\ \vec{p}^{(i)} &= \vec{r}^{(i-1)} + \beta_{i-1} \vec{p}^{(i-1)}, & \vec{p}'^{(i)} &= \vec{r}'^{(i-1)} + \beta_{i-1} \vec{p}'^{(i-1)}, \\ \alpha_i &= \frac{\vec{r}'^{(i-1)} \cdot \vec{r}^{(i-1)}}{\vec{p}'^{(i)} \cdot A \cdot \vec{p}^{(i)}} & \beta_i &= \frac{\vec{r}'^{(i)} \cdot \vec{r}^{(i)}}{\vec{r}'^{(i-1)} \cdot \vec{r}^{(i-1)}} & \vec{r}'^{(i)} \cdot \vec{r}^{(j)} &= \vec{p}'^{(i)} \cdot A \cdot \vec{p}^{(j)} = 0, \quad i \neq j \end{aligned}$$

Quasi-minimal residual (QMR), conjugate gradient squared (CGS), biconjugate gradient stabilized (BiCGStab).

18

## Iterative improvement

$$A\vec{x} = \vec{b}$$

Suppose the solution  $\vec{\tilde{x}}$  obtained with a direct method

is approximate due to round-off error. Define  $\delta\vec{x} = \vec{\tilde{x}} - \vec{x}$

$$\vec{r} = \vec{b} - A\vec{\tilde{x}} = A\vec{x} - A\vec{\tilde{x}} = -A\delta\vec{x}$$

Solve for  $\delta\vec{x}$ . Can be repeated, in which case decomposition of A should be done.

## The eigenvalue problem

$$A \vec{x} = \lambda \vec{x}$$

In principle, eigenvalues are the roots of the characteristic polynomial:

$$\det(A - \lambda I) = 0$$

Can use a root-finding procedure. Not a good idea, except perhaps for very small matrices and some special cases. Coefficients are subject to roundoff error and the roots are very sensitive to the values of the coefficients. One important observation, though:

$$\begin{pmatrix} 0 - \lambda & 1 & 0 & \dots & 0 \\ 0 & 0 - \lambda & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 \\ -a_0 & -a_1 & -a_2 & \dots & -a_{n-1} - \lambda \end{pmatrix}$$

The determinant of this is (up to a sign)

$$a_0 + a_1 \lambda + a_2 \lambda^2 + \dots + a_{n-1} \lambda^{n-1} + \lambda^n$$

Thus, for any polynomial there is an eigenvalue problem with the same roots. But there is no exact universal formula for order  $> 4$ . Thus, universal direct methods for eigenvalues cannot exist either – iterations only.

20

Just like for solving linear systems, we can ask the question of “condition” or sensitivity to perturbations in the matrix elements. The relevant condition number is very different from that for solving linear equations. In fact, for symmetric or Hermitian matrices the condition number is 1, so these problems are well-conditioned.

Different methods serve different purposes: first of all, some are preferable when only one or a few eigenvalues are desired (say, the largest or the closest to a particular value) and others are designed to find all of them. Type of matrix matters, too.

The basic approach finds one eigenvalue (largest by absolute value).

Start with an arbitrary vector, multiply with the matrix repeatedly.

Suppose the eigenvalues are  $\lambda_i$  and the corresponding eigenvectors are  $\vec{u}_i$ .

Suppose also that  $\lambda_1$  is the largest one by absolute value.

Start with vector  $\vec{v}^{(0)}$ . Expand it:  $\vec{v}^{(0)} = \sum_{i=1}^N c_i \vec{u}_i$

After  $k$  multiplications,

$$\vec{v}^{(k)} = \mathbf{A}^k \vec{v}^{(0)} = \mathbf{A}^k \sum_{i=1}^N c_i \vec{u}_i = \sum_{i=1}^N \lambda_i^k c_i \vec{u}_i = \lambda_1^k (c_1 \vec{u}_1 + \sum_{i=2}^N c_i (\lambda_i / \lambda_1)^k \vec{u}_i)$$

So the relative weight of eigenvectors other than the first one decays with  $k$

The ratio of the values of a given component from one iteration to the next converges to the eigenvalue.

What can go wrong with this approach?

The starting vector may have no component in the dominant eigenvector. Usually not a problem, because even if this is the case initially, it will be introduced due to round-off errors.

What if there are several vectors with the same largest absolute value of the eigenvalue? Let's suppose there are two. Keep just these terms and neglect everything else.

$$\vec{v}^{(k)} \approx \lambda_1^k (c_1 \vec{u}_1 + c_2 (\lambda_2/\lambda_1)^k \vec{u}_2)$$

If  $\lambda_1 = \lambda_2$ , then this will converge to some linear combination of the vectors.

But if they are different (only equal by absolute value), there will be never-ending oscillations.