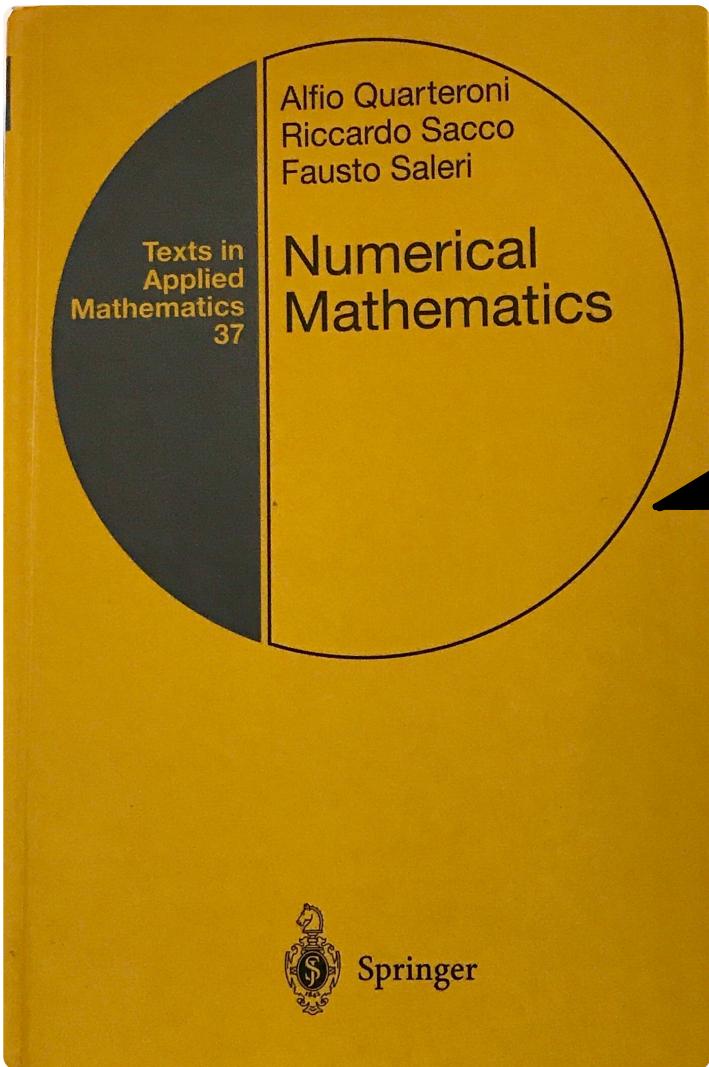


# ME 5311

## Solution of $A\vec{x} = \vec{b}$

Plan for today:

- Solution methods for linear systems
- Our first "grown up" discussion about CFD
  - Multi-dimensional problems
  - Method of solution is important



Chapters 3 and 4

on Husky GT

- Great numerical analysis textbook
- Plenty of examples/problems
- Not really a CFD textbook

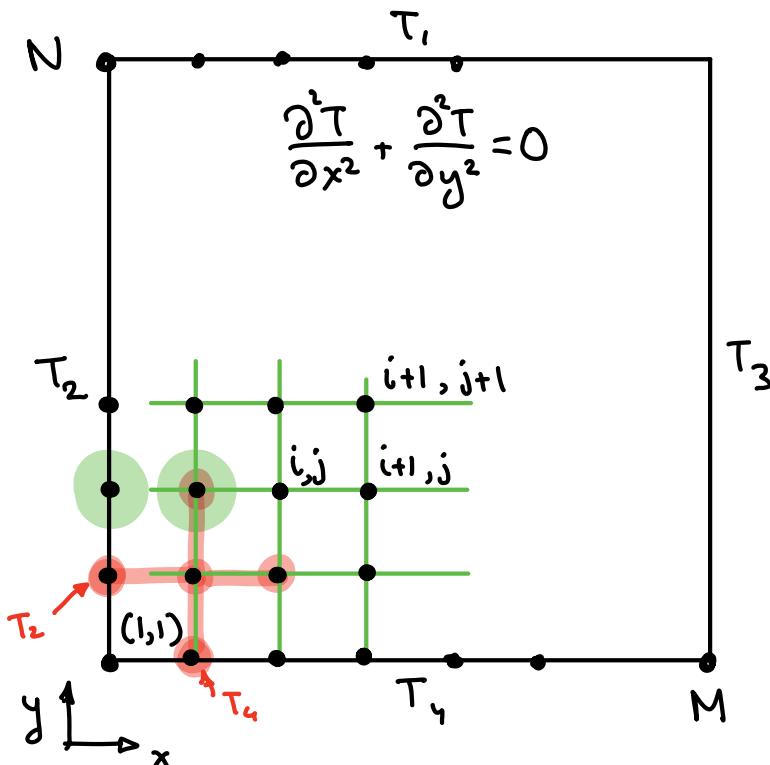
We want to solve  $A \vec{x} = \vec{b}$

$$\begin{matrix} A & \vec{x} & \vec{b} \\ n \times n & n \times 1 & n \times 1 \end{matrix}$$

Where do we find  $A \vec{x} = \vec{b}$ ?

- Implicit method  $\rightarrow$  space  $\frac{\partial}{\partial x}$ , Padé  
 $\rightarrow$  time  $\frac{\partial}{\partial t}$
- Steady  $A \vec{x} = \vec{b}$   
State problems
- Poisson equation for Pressure

# Example A: Steady-state heat equation



$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} \rightarrow (i,j)$$

Find:  $T(x,y)$  steady-state

Steps: 1. Governing eq. (+BC)

2. Discretization

3. Numerical approximation

4. Solution method

→ How many unknowns  $T_{ij}$ ?

$$(M-2) \times (N-2)$$

all the interior  
points

$$\frac{T_{i+1,j} - 2T_{i,j} + T_{i-1,j}}{\Delta x^2} + \frac{T_{i,j+1} - 2T_{i,j} + T_{i,j-1}}{\Delta y^2} = 0 \quad (M-2) \times (N-2)$$



$$i=2, j=2$$

$$\frac{T_{3,2} - 2T_{2,2} + T_{2,1}}{\Delta x^2} + \frac{T_{2,3} - 2T_{2,2} + T_{1,2}}{\Delta y^2} = 0 \quad \leftarrow$$

Form the linear system

$$\begin{aligned} & -\frac{2}{\Delta x^2} - \frac{2}{\Delta y^2} + \frac{1}{\Delta x^2} \\ & \frac{1}{\Delta x^2} \quad \bullet_{i-1,j} \quad \bullet_{i,j} \quad \bullet_{i+1,j} \end{aligned}$$

$$\left[ \begin{array}{c|c} \begin{matrix} T_{2,2} \\ T_{3,2} \\ T_{4,2} \\ \vdots \\ T_{M-1,2} \end{matrix} & \begin{matrix} +\frac{T_2}{\Delta x^2} + \frac{T_4}{\Delta y^2} \\ +\frac{T_2}{\Delta x^2} \end{matrix} \\ \hline \begin{matrix} T_{2,3} \\ T_{3,3} \\ T_{4,3} \\ \vdots \\ T_{M-1,3} \end{matrix} & = -b \end{array} \right]$$

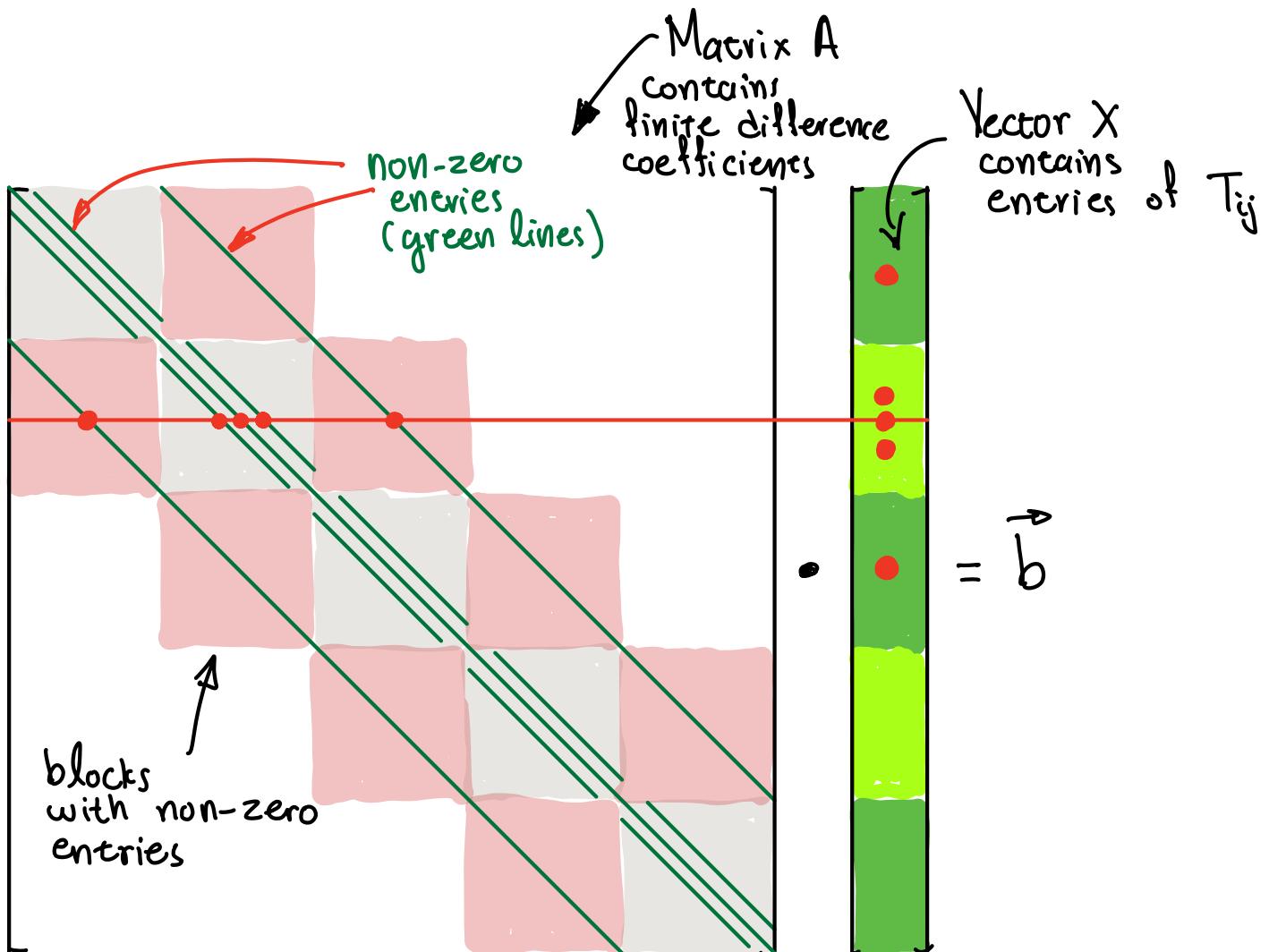
$\rightarrow x$

$$\begin{matrix} \text{number of elements} \\ (M-2)(N-2) \times (M-2)(N-2) \end{matrix} \underbrace{(M-2)(N-2)}_{N_x} \quad \underbrace{(M-2)(N-2)}_{N_y}$$

$\uparrow$   
number of elements :  $(M-2)(N-2)$

## Structure of matrix A:

- Block diagonal



# What is the computational cost?

Memory / storage

# elements  
↓

$$1D \quad T_i : N_x \rightarrow A = N_x^2$$

$$2D \quad T_{i,j} : N_x \cdot N_y \rightarrow A = (N_x \cdot N_y)^2$$

$$3D \quad T_{i,j,k} : N_x \cdot N_y \cdot N_z \rightarrow = (N_x \cdot N_y \cdot N_z)^2$$

Memory

$$N = 256 \quad A : 65 \times 10^3 \times 8 \text{ bytes} \rightarrow 0.5 \text{ MB}$$

$$4.3 \times 10^3 \quad 34 \text{ GB}$$

$$2.8 \times 10^{14}$$

2252 TB

OMG!!

CPU:

Cramer's  $x_j = \frac{D_j}{\det(A)}$  ← det A but j-column n with b

Cost is  $(n+1)!$  operations

$$n = 50 \rightarrow 51! = 1.5 \times 10^{66} \text{ operations}$$

Intel Xeon "Sky lake"

1000 Gflops

1 flop: floating point

operation per sec.

$$\frac{1.5 \times 10^{66}}{\underbrace{\# \text{ operations}}_{\text{flops}}} / \underbrace{1000 \times 10^9}_{\text{flops}} = 1.5 \times 10^{54} \text{ sec}$$

$$= 5 \times 10^{46} \text{ years}$$

# Conclusions:

1. In CFD the boundary between "possible" computation and "not possible" is very sharp

$\Delta x$  becomes  $\frac{\Delta x}{2}$  in 3D  $\Rightarrow \Delta t$  becomes  $\frac{\Delta t}{2}$

$N$  grid points  $\rightarrow 2N$  points in each direction

$N_t$  time steps  $\rightarrow 2N_t$  time steps

Old cost  $N^3 N_t$   $\rightarrow$  new cost  $2^3 N 2 N_t = 16$  (old cost)

2. Good algorithms are important

Classic example is FFT : cost is  $N \log N$  operations instead of  $N^2$

for  $N=1000$  :  $N^2 = 1,000,000$   
 $N \log N \approx 7,000$  factor of about 150 difference

## Very important:

- The matrix is **sparse**

Sparse: most elements are zero

Dense: (not sparse) most elements are not zero

- Do not form/store the entire matrix  $A$
- Data structure stores only the non-zero elements

e.g.: row, column, value

$i$ ,  $j$ ,  $a_{i,j}$

- Typically we only need to find  $\vec{y} = A \vec{x}$ 
  - we can do this without explicitly forming  $A$
  - e.g.

```
function y = matrix_multiply(x)
```

- Do not form the inverse
  - The inverse is dense!
  - $x = A^{-1}b$  is just notation

If you need to construct the matrix use a sparse matrix format!

## Example B: Convection-Diffusion Problem

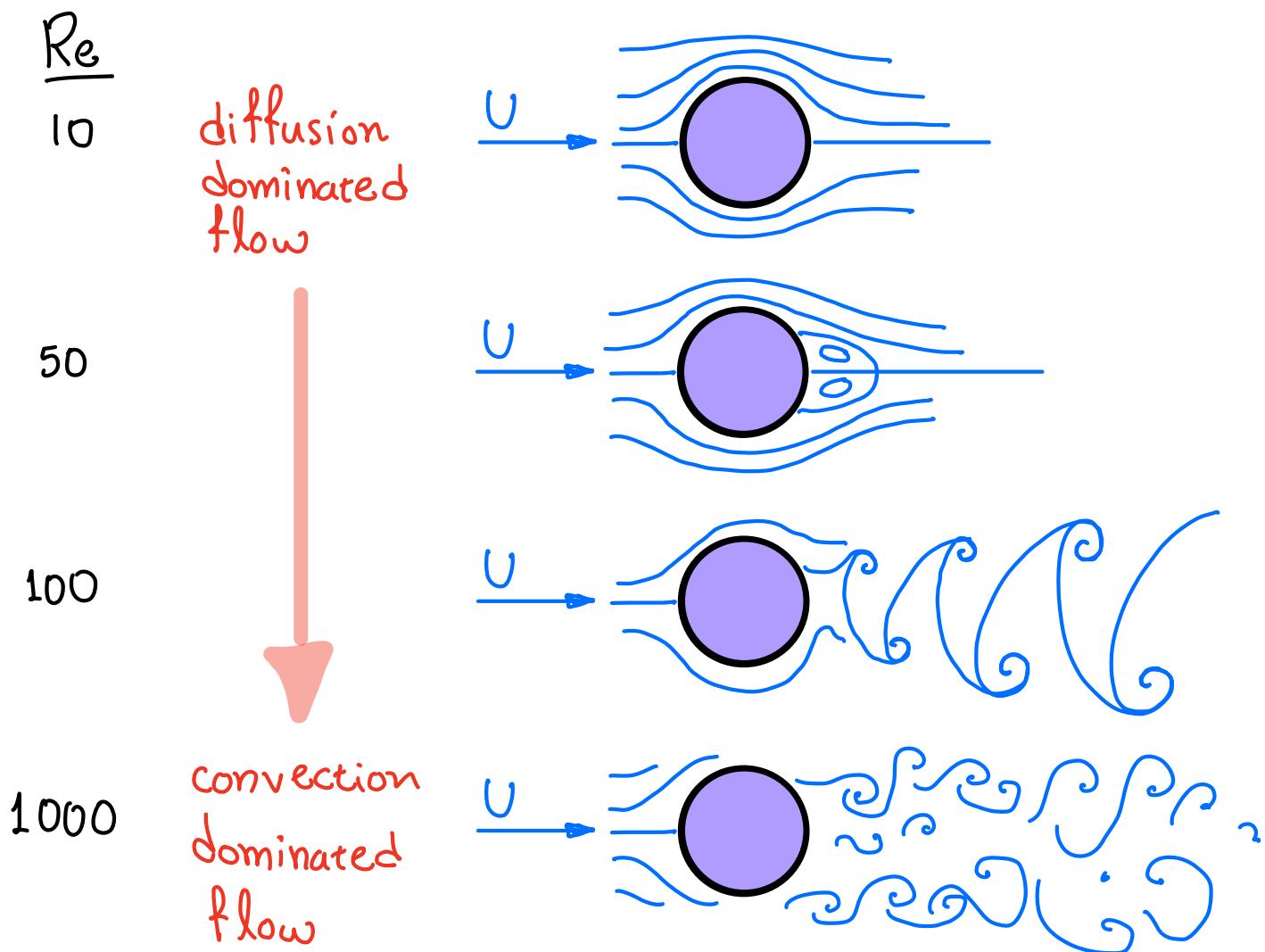
Fluid Dynamics Problem: flow around cylinder

Parameters: fluid velocity  $U$   
(dimensional)

fluid kinematic viscosity  $\nu = \frac{\mu}{\rho}$

cylinder diameter  $D$

Non-dimensional parameter: Reynolds number  $Re = \frac{UD}{\nu}$



# Model Problem: One-dimensional convection-diffusion

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = v \frac{\partial^2 u}{\partial x^2} \quad \text{on periodic domain}$$

$$\frac{d\vec{u}}{dt} = \left[ -\frac{c}{2\Delta x} B_p(-1, 0, 1) + \frac{v}{\Delta x^2} B_p(1-2, 1) \right] \vec{u}$$

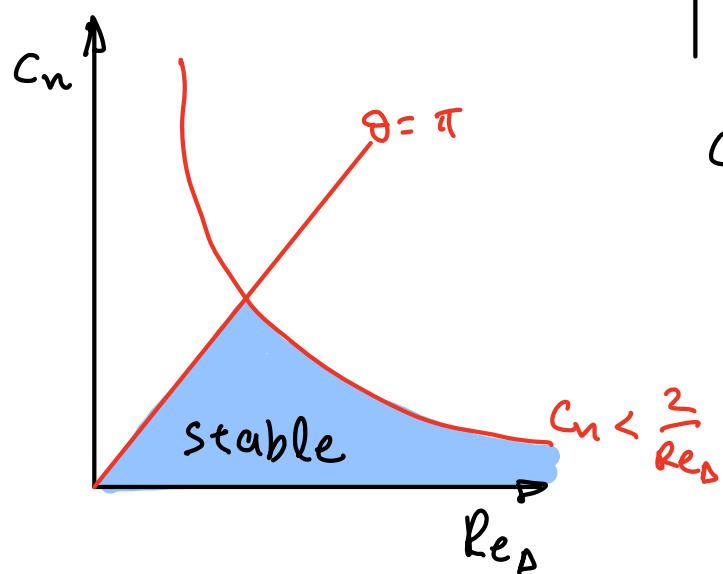
$$\frac{d\vec{u}}{dt} = (A_c + A_d) \vec{u}$$

CFL number  $C_n = \frac{ch}{\Delta x} \quad \Delta t$

Explicit-Explicit

$$\vec{u}^{n+1} = \vec{u}^n + h(A_c + A_d) \vec{u}^n$$

$$|\sigma| = \sqrt{1 + C_n^2 \sin^2 \theta} \left[ 1 - 4 \frac{C_n}{Re_\Delta} \sin^2 \frac{\theta}{2} \right]$$



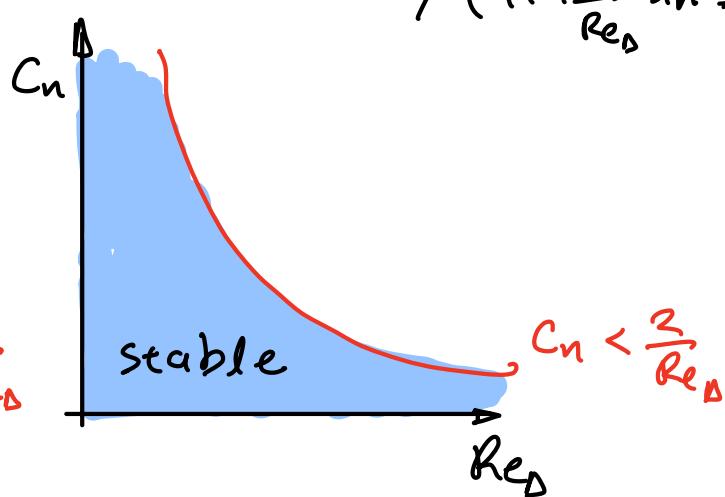
grid Reynolds  $Re_\Delta = \frac{c \Delta x}{v}$

Explicit - Implicit-

$$\frac{u^{n+1} - u^n}{h} = A_c u^n + A_d u^{n+1}$$

$$\underbrace{[I - hA_d]}_A \vec{u}^{n+1} = \underbrace{[I - hA_c]}_A \vec{u}^n \quad \vec{x} = \vec{b}$$

$$|\sigma| = \sqrt{1 + C_n^2 \sin^2 \theta} / \left( 1 + 4 \frac{C_n}{Re_\Delta} \sin^2 \frac{\theta}{2} \right)$$



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

Direct

given  $A, \vec{b}$   $\xrightarrow{\text{seq.}}$   $\vec{x}$

Iterative

start  $\vec{x}^{(0)}$   $\xrightarrow{\text{seq.}}$   $\vec{x}^{(n)} \approx \vec{x}$

Pro: "exact"  $\vec{x}$

Con: computationally  
expensive

Pro: "cheaper"

Con: infinite number  
of iterations

Solution Methods for  $A\vec{x} = \vec{b}$ :

- Direct
- Iterative

**Direct:** given  $A, \vec{b}$  perform a sequence of operations to get  $\vec{x}$ , e.g. Cramer's rule

positive: can get an "exact"  $\vec{x}$

negative: is expensive

**Iterative:** start with an initial guess for  $\vec{x}_0$  and produce a sequence  $\vec{x}_n \rightarrow \vec{x}$

positive: much "cheaper" than direct

negative: infinite number of iterations to get to "exact"  $\vec{x}$ , must accept an error tolerance

When solution  $A\vec{x} = \vec{b}$  exists?  $\det(A) \neq 0$

Condition number:  $K(A) = \|A\| \|A^{-1}\|$

number:  $K = \infty$  if  $A$  is singular

Matlab `cond(A)`

# Solution of Triangular systems

A

$$\begin{bmatrix} a_{11} & & & \\ a_{21} & a_{22} & & 0 \\ & a_{32} & a_{33} & \\ & & \ddots & \\ & & & a_{nn} \end{bmatrix} \quad \vec{x} = \vec{b}$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad \vec{x} = \vec{b}$$

$$L \quad \vec{x} = \vec{b}$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad \vec{x} = \vec{b}$$

$$U \quad \vec{x} = \vec{b}$$

$$x_1 = \frac{b_1}{a_{11}}$$

$$x_2 = \frac{b_2 - a_{21}x_1}{a_{22}}$$

$\vdots$

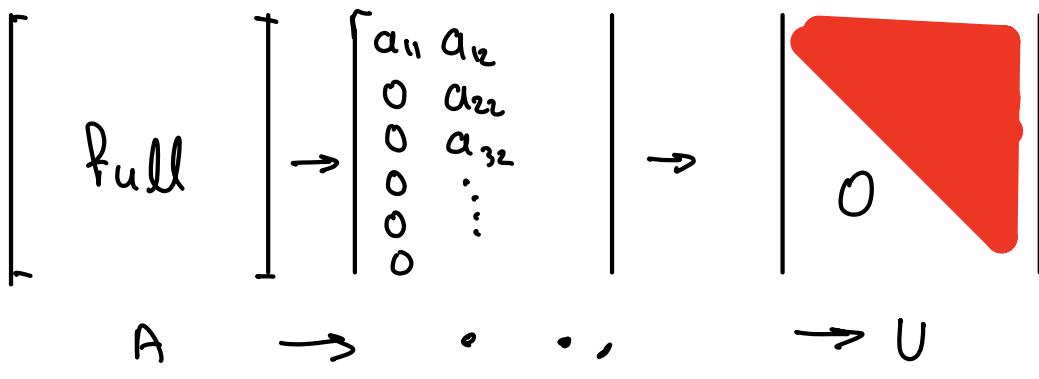
multiplications and division :  $\frac{n(n+1)}{2}$   
 sum and subtractions :  $\frac{n(n-1)}{2}$

$$\text{Total} = n^2$$

A:  $n \times n$

# Gaussian Elimination

$$A \vec{x} = \vec{b} \rightarrow U \vec{x} = \vec{f}$$



Gauss elimination :  $2(n-1)n(n+1)/3 + n(n-1) \sim \frac{2}{3}n^3$   
+  $n^2$  to solve  $U\vec{x} = \vec{b}$

Total :  $\frac{2n^3}{3} + 2n^2 \sim O(n^3)$

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 5 \\ 7 & 8 & 9 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & -1 \\ 0 & -6 & -12 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 0 & -6 & -12 \\ 0 & 0 & -1 \end{bmatrix}$$

↑  
Pivoting

Partial                       $n^2$  searches  
complete                     $\frac{2n^3}{3}$

# LU factorization

$$A \rightarrow A^{(1)} \rightarrow A^{(2)} \rightarrow \dots A^{(n)} = U$$

matrix  $n \times n$

$$A^{(1)} = M_1 A$$

$$A^{(k+1)} = M_k A^{(k)}$$

iteration #

$$M_{n-1} M_{n-2} \dots M_1 A = U$$

$$M_k = \begin{vmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 0 \\ & -m_{2,1}^{-1} & & & \\ & \vdots & & & \\ & & & 1 & \\ & & & & k \end{vmatrix}$$

$\leftarrow$  Gaussian transformation matrix

$$A = \underbrace{M_1^{-1} M_2^{-1} \dots M_{n-1}^{-1}}_{U} U$$

$$A = L U$$

$\leftarrow$  lower triangular    diagonal entries = 1

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

$$\underbrace{L U}_{\vec{y}} \vec{x} = \vec{b}$$

$$\vec{y}$$

$$L \vec{y} = \vec{b} \Rightarrow \text{find } \vec{y} \Rightarrow U \vec{x} = \vec{y} \Rightarrow \text{get } \vec{x}$$

cost :  $O(n^3)$

" $A^{-1} b$ "

→ notation

if we have to compute  $A^{-1}$  ...

$$AA^{-1} = I$$

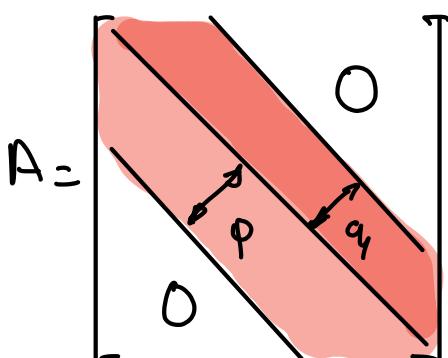
$$A X = I$$

$$A \vec{x}_i = \vec{e}_i \quad e_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{bmatrix} \leftarrow i$$

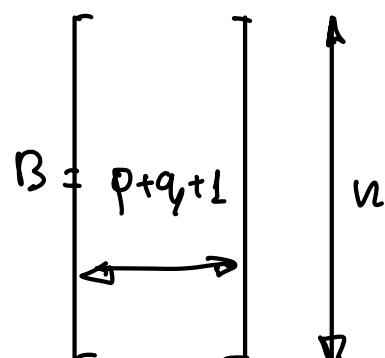
$$A = LU$$

$$\left. \begin{array}{l} LU \vec{x}_i = \vec{e}_i \\ L \vec{y} = \vec{e}_i \\ U \vec{x}_i = \vec{y} \end{array} \right\} \text{Solve } 2n \text{ triangular systems}$$

## Banded systems



→ stored



sparse matrix:  $\begin{bmatrix} i, j, \text{value} \end{bmatrix}$

Tri diagonal

$$A = \begin{bmatrix} a_1 & c_1 & & & \\ b_2 & a_2 & c_2 & & 0 \\ & \ddots & \ddots & & \\ 0 & b_{n-1} & a_{n-1} & c_{n-1} & \\ & b_n & a_n & & \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & & & & \\ b_2 & 1 & 0 & & \\ b_3 & b_2 & 1 & & \\ & \ddots & \ddots & \ddots & \\ 0 & & \ddots & & 1 \end{bmatrix} \quad U = \begin{bmatrix} \alpha_1 & c_1 & & & 0 \\ \alpha_2 & c_2 & & & \\ & \ddots & \ddots & & \\ 0 & & \ddots & \ddots & \alpha_n \end{bmatrix}$$

$$\alpha_1 = a_1 \quad b_i = \frac{b_i}{\alpha_{i-1}} \quad \alpha_i = a_i - b_i c_{i-1}, \quad i=2, \dots, n$$

Thomas       $L \vec{y} = \vec{f}$

Algorithm:

$$y_1 = f_1$$

$$y_i = f_i - b_i y_{i-1}$$

$$U \vec{x} = \vec{y} \quad x_n = \frac{y_n}{\alpha_n}$$

$$x_i = (y_i - c_i x_{i+1}) / \alpha_i$$

Cost:  $8n-7$

# Thomas Algorithm

without divisions :  $A = L D M^T$

$$\begin{bmatrix} \gamma_1^{-1} \\ b_2 \gamma_2^{-1} \\ \vdots \\ - \end{bmatrix} \begin{bmatrix} \gamma_1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \gamma_n \end{bmatrix} \begin{bmatrix} \gamma_1^{-1} c_1 \\ r_2^{-1} \\ \vdots \\ - \end{bmatrix}$$

$$\gamma_i = (a_i - b_i \gamma_{i-1} c_{i-1})^{-1} \quad i = 1, \dots, n$$

$$\gamma_0 = 0, \quad b_1 = 0, \quad c_n = 0$$

$$L \vec{y} = \vec{f} : \quad y_1 = \gamma_1 b_1 \quad y_i = \gamma_i (f_i - b_i y_{i-1})$$

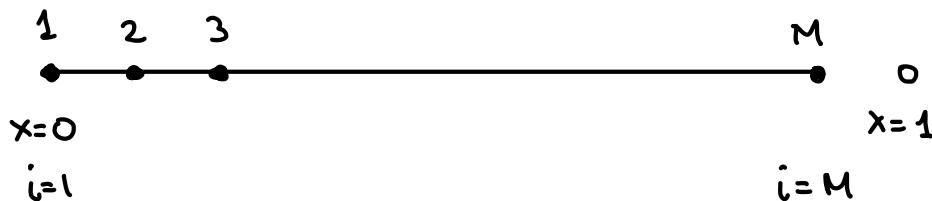
$$U \vec{x} = \vec{y} \quad x_n = y_n \quad x_i = y_i - \gamma_i c_i x_{i+1}$$

## Direct methods: Using Fourier Transform

Solve 1D Poisson equation  $\frac{d^2 p}{dx^2} = d(x)$  on  $x \in [0, L]$

with periodic BCs.

$$\text{Discretization: } \frac{1}{\Delta x^2} (p_{i-1} - 2p_i + p_{i+1}) = d_i \quad (1)$$



$$\Delta x = \frac{L}{M} \text{ uniform grid spacing}$$

$$\text{Note that } d_i = d(x_i) = d(i \Delta x)$$

$$\text{if } d(x) = \cos(x) \rightarrow d_i = \cos(i \Delta x)$$

Solution: need to find  $\vec{p}$  :  $A \vec{p} = \vec{d}$

$$\frac{1}{\Delta x^2} \begin{bmatrix} 1 & & & \\ & -2 & 1 & \\ & 1 & -2 & 1 \\ & & 1 & -2 \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_M \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_M \end{bmatrix}$$

### Solution method:

We will take advantage of the periodicity in  $x$  and utilize a Fourier expansion

$$\text{Remember: } p_i = \sum_{k=0}^{M-1} \hat{p}_k e^{-ik \Delta x i} \quad (2)$$

$\uparrow$  Fourier coefficients from x-direction  
 Discrete Fourier Transform

Use ② in ①:

$$\sum_{k=0}^{M-1} \left[ \frac{1}{\Delta x^2} \hat{P}_k \left( e^{ik\Delta x(i-1)} - 2e^{ik\Delta x i} + e^{ik\Delta x(i+1)} \right) \right] = \sum_{k=0}^{M-1} \hat{d}_k e^{ik\Delta x i}$$

$$\sum_{k=0}^{M-1} \left[ \frac{1}{\Delta x^2} \hat{P}_k \left( e^{-ik\Delta x} - 2 + e^{ik\Delta x} \right) e^{ik\Delta x i} \right] = \sum_{k=0}^{M-1} \hat{d}_k e^{ik\Delta x i}$$

Because of the orthogonality property of the  $e^{ik\Delta x i}$   
we can write separate equations for each  $k$

Thus...

$$\frac{1}{\Delta x^2} \hat{P}_k \left( e^{-ik\Delta x} - 2 + e^{ik\Delta x} \right) = \hat{d}_k \quad \text{for each } k = 0, \dots, M-1$$

$$\frac{1}{\Delta x^2} \hat{P}_k \left( \cos(-k\Delta x) + i \sin(-k\Delta x) - 2 + \cos(k\Delta x) + i \sin(k\Delta x) \right) = \hat{d}_k$$

$$\frac{1}{\Delta x^2} \hat{P}_k \left( 2 \cos(k\Delta x) - 2 \right) = \hat{d}_k$$

$$\frac{1}{\Delta x^2} \hat{P}_k \left( 2 \cos\left(\frac{2\pi k}{M}\right) - 2 \right) = \hat{d}_k \quad k = 0, \dots, M-1$$

$$\hat{P}_k = \Delta x^2 \hat{d}_k \left( 2 \cos\left(\frac{2\pi k}{M}\right) - 2 \right)^{-1}$$

$$p_i = \text{inverse FT}(\hat{P}_k)$$

**Algorithm:** Input: Vector of  $d_i$ ,  $\Delta x$ ,  $M$

1. Compute FT of  $d_i \rightarrow \hat{d}_k$

2. Compute  $\hat{\rho}_k = \Delta x^2 \hat{d}_k \left(2 \cos\left(\frac{2\pi k}{M}\right) - 2\right)^{-1}$

3. Compute  $\rho_i = \text{inverse FT } \hat{\rho}_k$

Note that  $\hat{\rho}_{k=0}$  is undefined

## Fundamentals:

- Spectrum  $\sigma(A)$ : the set of eigenvalues  $\lambda$  of  $A$
- Spectral radius  $\rho(A) = \max |\lambda|, \lambda \in \sigma(A)$
- $\rho(A) \leq \|A\|$   $\leftarrow$  any  $\|\cdot\|$
- Positive definite  $\vec{A}\vec{x} \cdot \vec{x} > 0 \quad \forall \vec{x} \neq 0$   
 $\nwarrow$   
 inner product
- Diagonally dominant
  - by rows  $|a_{ii}| \geq \sum_{j=1, j \neq i}^n |a_{ij}| \quad i=1, \dots, n$
  - by columns  $|a_{ii}| \geq \sum_{i=1, i \neq j}^n |a_{ij}| \quad j=1, \dots, n$
- Convergence
 
$$\vec{x} = \lim_{k \rightarrow \infty} \vec{x}^{(k)}$$

$\uparrow$  sequence of  $x$  vectors

$$\vec{e}^{(k)} = \vec{x}^{(k)} - \vec{x}$$

$\uparrow$  exact

$$\lim_{K \rightarrow \infty} e^{(k)} = 0$$
- Sequence  $\vec{x}^{(k+1)} = B \vec{x}^{(k)} + \vec{f} \quad \vec{x}^{(0)} \text{ initial guess given}$ 
  - $\uparrow$  iteration
  - $\uparrow$   $B \in \mathbb{R}^{n \times n}$
  - $\vec{f} \in \mathbb{R}^n$

- Cost  $B \vec{x} : O(n^2)$

$LU : O(\frac{2}{3}n^3)$

hope that  $n^{0.492}$

- Consistency :  $\vec{f} = (I - B) A^{-1} \vec{b}$       equivalence  
 $\vec{x} = B \vec{x} + \vec{f}$

- Consistency +  $\rho(B) < 1 \Rightarrow$  Convergence

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

$$\vec{x}^{(k)} = e^{(k)} + \vec{x}$$

$$\underbrace{e^{(k+1)} + \vec{x}}_{\vec{x}^{(k+1)}} = B \underbrace{(e^{(k)} + \vec{x})}_{\vec{x}^{(k)}} + \vec{f}$$

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

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

$$e^{(k+1)} = B^k e^{(0)}$$

$$\lim_{k \rightarrow \infty} (B^k e^{(0)}) = 0 \quad \text{iff} \quad \rho(B) < 1$$

or...  $1 > \|B\| \geq \rho(B)$

- $\|B^m\|$  convergence factor  
 $\|B^m\|^{1/m}$  average convergence factor after  $m$ -steps

$$R_m(B) = \frac{1}{m} \log \|B^m\| \quad \text{average convergence rate after } m \text{ steps}$$

$$R(B) = \lim_{k \rightarrow \infty} R_k(B) = -\log \rho(B)$$

$\uparrow$  asymptotic convergence rate

When to stop ???

$$\|e^{(k)}\| \leq \varepsilon \|e^{(0)}\|$$

$\uparrow$  tolerance

$$\# \text{ iterations: } k \geq \frac{\log(\varepsilon)}{\frac{1}{2} \log \|B^k\|} = -\frac{\log(\varepsilon)}{R_k(B)}$$

$$k_{\min} \approx -\frac{\log(\varepsilon)}{R(B)} \quad \uparrow$$

1. Stop based on increment

$$\|\vec{x}^{(k+1)} - \vec{x}^{(k)}\| < \varepsilon$$

$$\frac{\|\vec{x}^{(k+1)} - \vec{x}^{(k)}\|}{\|\vec{b}\|} < \varepsilon$$

2. Stop based on the residual

$$\frac{\|\vec{r}^{(k)}\|}{\|\vec{b}\|} < \varepsilon \quad r^{(k)} = A\vec{x}^{(k)} - \vec{b}$$

# Linear Iterative methods

$$A = P - N$$

↑ Preconditioner

$$P \vec{x}^{(k+1)} = N \vec{x}^{(k)} + \vec{b} \quad k \geq 0$$

$$\vec{x}^{(k+1)} = \underbrace{P^{-1}N}_{B} \vec{x}^{(k)} + \underbrace{P^{-1}\vec{b}}_{\vec{f}}$$

↑ Iteration matrix

$$\begin{aligned} \vec{x}^{(k+1)} &= P^{-1}(P - A)\vec{x}^{(k)} + P^{-1}\vec{b} \\ &= (I - P^{-1}A)\vec{x}^{(k)} + P^{-1}\vec{b} \\ &= \vec{x}^{(k)} + P^{-1}(-A\vec{x}^{(k)} + \vec{b}) \\ &\qquad\qquad\qquad \underbrace{-A\vec{x}^{(k)} + \vec{b}}_{r^{(k)}} \\ \vec{x}^{(k+1)} &= \vec{x}^{(k)} + P^{-1}r^{(k)} \end{aligned}$$

Jacoby:

$$\vec{x}_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} \vec{x}_j^{(k)} \right]$$

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = S(x, y)$$

$$\Downarrow \frac{1}{\Delta y^2} (T_{i,j-1} - 2T_{i,j} + T_{i,j+1}) + \frac{1}{\Delta x^2} (T_{i-1,j} - 2T_{i,j} + T_{i+1,j}) = S_{i,j}$$

$$T_{i,j}^{(k+1)} = \left( \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} \right)^{-1} \left[ \frac{1}{\Delta y^2} (T_{i,j-1}^{(k)} + T_{i,j+1}^{(k)}) + \frac{1}{\Delta x^2} (T_{i-1,j}^{(k)} + T_{i+1,j}^{(k)}) \right] + S_{i,j}$$

$\nwarrow$  diagonal of  $A$

$$P = D \quad \left. \begin{array}{l} \\ \end{array} \right\} \rightarrow B_3 = D^{-1}(D - A) = I - D^{-1}A$$

# Over-relaxation method

$$x^{(k+1)} = \frac{\omega}{a_{ii}} \left[ b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k)} \right] + (1-\omega) x_i^{(k)}$$

↑  
relaxation parameter  $\omega$

$$B_{J\omega} = \omega B_J + (1-\omega) I$$

↑ Jacobi iteration matrix

$$T_{i,j}^{(k+1)} = \omega \left( \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} \right)^{-1} \left[ \frac{T_{i+1,j}^{(k)} + T_{i-1,j}^{(k)}}{\Delta x^2} + \frac{T_{i,j+1}^{(k)} + T_{i,j-1}^{(k)}}{\Delta y^2} \right] + (1-\omega) T_{i,j}^{(k)}$$

under-relaxation  $0 < \omega < 1$

over relaxation  $\omega > 1$

## Gauss - Siedel :

$$T_{i,j}^{(k+1)} = \omega \left( \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} \right)^{-1} \left[ \frac{T_{i+1,j}^{(k)} + T_{i-1,j}^{(k+1)}}{\Delta x^2} + \frac{T_{i,j+1}^{(k+1)} + T_{i,j-1}^{(k+1)}}{\Delta y^2} \right] + (1-\omega) T_{i,j}^{(k)}$$

use available values at  $(k+1)$

$$\text{Residual : } r_{i,j} = \frac{T_{i-1,j} - 2T_{i,j} + T_{i+1,j}}{\Delta x^2} + \frac{T_{i,j-1} - 2T_{i,j} + T_{i,j+1}}{\Delta y^2}$$

then,  $\bullet \max |r_{i,j}|$

or  $\bullet \frac{1}{N_x N_y} \left( \sum_{i,j} |r_{i,j}|^2 \right)^{1/2}$

$$\|r_{i,j}\| = \max_{i,j} |r|$$

Optimal value of relaxation parameter  $\omega$ :

$$\omega_{\text{opt}} = 1 + \left( \rho / (1 + \sqrt{1 - \rho^2}) \right)^2$$

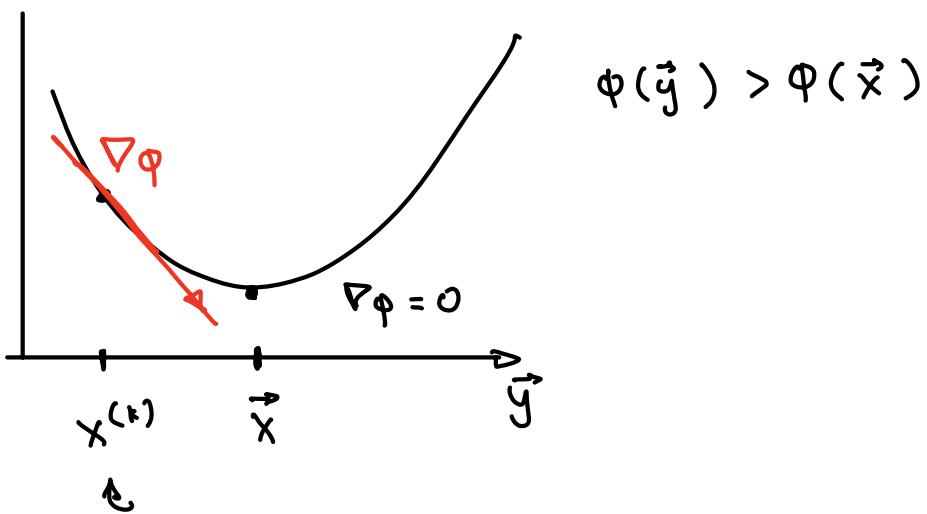
where  $\rho$  is the spectral radius of the Jacobi iteration matrix  $B_J = I - P^{-1}A$

# The Gradient Method

$$\Phi(\vec{y}) = \frac{1}{2} \vec{y}^T A \vec{y} - \vec{y}^T \vec{b}$$

$$\nabla \Phi(\vec{y}) = \frac{1}{2} (A^T + A) \vec{y} - \vec{b} = A \vec{y} - \vec{b}$$

$$\nabla \Phi(\vec{x}) = 0 \quad \text{because } A \vec{x} = \vec{b}$$



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

↑ direction  
↑ scalar

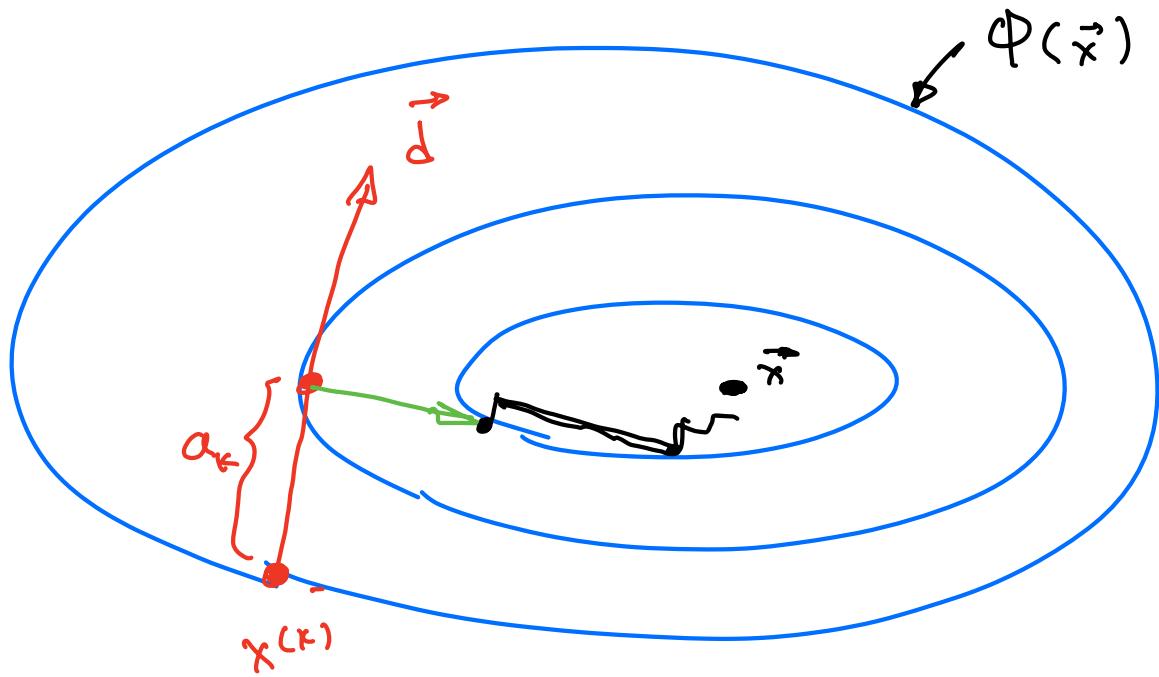
$$\text{To find } \alpha : \nabla_a \Phi(x^{(k+1)}) \rightarrow \alpha_k = \frac{r^{(k)^T} r^{(k)}}{r^{(k)^T} A r^{(k)}}$$

Method: given  $x^{(0)}$

$$r^{(k)} = \vec{b} - A x^{(k)}$$

$$\alpha_k = \frac{r^{(k)^T} r^{(k)}}{r^{(k)^T} A r^{(k)}}$$

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



## The Conjugate Gradient

choose direction  $\vec{p}$

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

$$\phi(\vec{x}^{(k)}) \leq \phi(\vec{x}^{(k)} + \gamma \vec{p})$$

$\uparrow$                    $\uparrow$   
optional      w.r.t.

⋮  
⋮      math

①  $p \perp r$        $\vec{p}$  and  $\vec{q}_j$  are  $A$ -conjugate

②  $x^{(k+1)} = x^{(k)} + \vec{q}_j \Rightarrow \underbrace{\vec{p}^T A \vec{q}_j}_{= 0}$

$$p^{(k+1)} = r^{(k+1)} + b_k p^{(k)}$$

$b_k$  scalar



## Relaxation Method

$$\text{PDE : } \frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = f(x, y)$$

$$\text{Discretization: } \frac{T_{i-1,j} - 2T_{i,j} + T_{i+1,j}}{\Delta x^2} + \frac{T_{i,j-1} - 2T_{i,j} + T_{i,j+1}}{\Delta y^2} = f_{i,j}$$

Iteratively solve for  $T_{i,j}$

$$T_{i,j}^{(k+1)} = \omega \left( \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} \right)^{-1} \left[ \frac{T_{i+1,j}^{(k)} + T_{i-1,j}^{(k+1)}}{\Delta x^2} + \frac{T_{i,j+1}^{(k)} + T_{i,j-1}^{(k+1)}}{\Delta y^2} - f_{i,j} \right] + (1-\omega) T_{i,j}^{(k)}$$

Iteration

use available values at  $(k+1)$

relaxation parameter  $0 < \omega < 2$

Residual

$$r_{i,j} = \frac{T_{i-1,j} - 2T_{i,j} + T_{i+1,j}}{\Delta x^2} + \frac{T_{i,j-1} - 2T_{i,j} + T_{i,j+1}}{\Delta y^2} - f_{i,j}$$

$\uparrow$  this is a 2D array

we can make a scalar by taking  $\max |r_{ij}|$

Optimal value of relaxation parameter  $\omega$ :

$$\omega_{opt} = 1 + \left( \rho / (1 + \sqrt{1 - \rho^2}) \right)^2$$

where  $\rho$  is the spectral radius of the Jacobi iteration

$$\text{matrix } B_J = I - P^{-1}A$$

$\uparrow P = \text{diagonal of } A$