# Numerical Methods I
## MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

# Today

**Last time**

- ▶ Singular value decomposition (SVD)
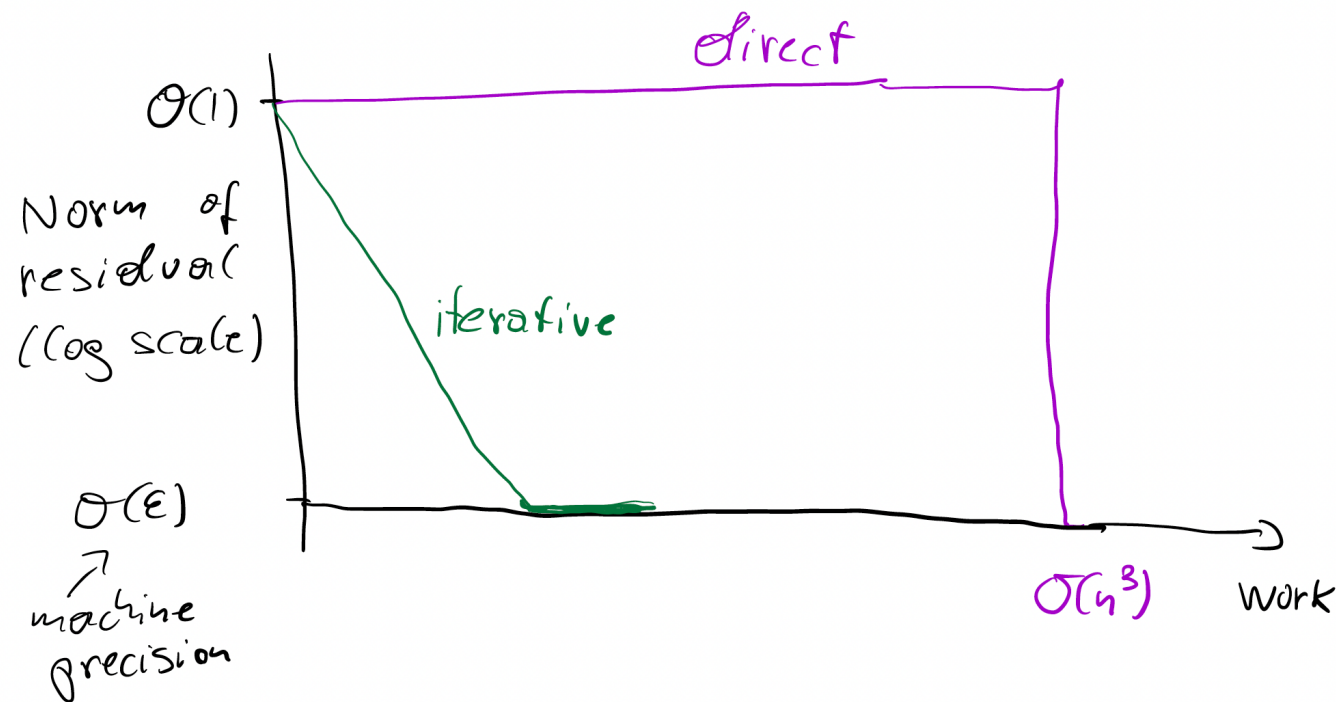- ▶ Iterative methods for systems of linear equations

**Today**

- ▶ Iterative methods
- ▶ Conjugate gradients method

**Announcements**

- ▶ Homework 4 is due Mon, Nov 7, 2022 before class
- ▶ Midterm grades posted

# Recap: Iterative methods



- Iterative methods *can* converge geometrically until residual is below machine precision
- Direct methods make no progress at all until $\mathcal{O}(n^3)$ work is done, and then lead to residual on the order of machine precision

# Recap: Iterative methods

Let $Q$ be invertible, then

$$A\boldsymbol{x} = \boldsymbol{b} \Leftrightarrow Q^{-1}(\boldsymbol{b} - A\boldsymbol{x}) = 0$$
$$\Leftrightarrow (I - Q^{-1}A)\boldsymbol{x} + Q^{-1}\boldsymbol{b} = \boldsymbol{x}$$
$$\Leftrightarrow G\boldsymbol{x} + \boldsymbol{c} = \boldsymbol{x}$$

Leads to fixed-point iteration

$$\boldsymbol{x}_{k+1} = \boldsymbol{G}\boldsymbol{x}_k + \boldsymbol{c}$$

and $\boldsymbol{x} = \boldsymbol{A}^{-1}\boldsymbol{b}$ is stationary point

$$\boldsymbol{G}\boldsymbol{x} + \boldsymbol{c} = (\boldsymbol{I} - \boldsymbol{Q}^{-1}\boldsymbol{A})\boldsymbol{x} + \boldsymbol{Q}^{-1}\boldsymbol{b} = \boldsymbol{x} - \boldsymbol{Q}^{-1}\boldsymbol{b} + \boldsymbol{Q}^{-1}\boldsymbol{b} = \boldsymbol{x}$$

# Recap: Common choices for $Q$

Requirements on a good $Q$ are

- $Q^{-1}$ is representative of $A^{-1}$
- We can numerically solve $Qy = z$ much quicker than $Ax = b$ because in each step we solve

$$Qx_{k+1} = (Q - A)x_k + b$$

Split the matrix $A$ as follows

$$A \quad = \quad L \quad + \quad D \quad + \quad U$$

# Recap: Jacobi method

Theorem: Select now $Q = D$ ... Jacobi method. The Jacobi method converges for any starting point $\mathbf{x}_o$ to the solution of $A\mathbf{x} = \mathbf{b}$ if $A$ is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \qquad \text{for } i = 1, \ldots, n.$$

Notice that $Q = D$ is a good choice in terms of computational costs because we can very quickly solve $Dy = z$ for the diagonal matrix $D$

# Gauss-Seidel method

Theorem: Choose $Q = D + L$ ... Gauss-Seidel method. The Gauss-Seidel method converges for any starting point $x_o$ if $A$ is symmetric positive definite (spd).
$\rightsquigarrow$ board

Gauss- Seidel

$$x_{k+1} = (I - (D+L)^{-1} A) x_k + (D+L)^{-1} b$$

$$= ((D+L)^{-1}(D+L) - (D+L)^{-1}(L + D + R)) x_k + (D+L)^{-1} b$$

$$= -(D+L)^{-1} R x_k + (D+L)^{-1} b$$

For any spd $A$, we have a scalar product

$$\langle x, y \rangle_A = \langle x, A y \rangle \qquad \text{on } \mathbb{R}^n$$

And for any matrix $B \in \mathbb{R}^{n \times n}$, we have adjoint
w.r.t. $\langle \cdot, \cdot \rangle_A$ given by

$$B^* = A^{-1} B^T A$$

so that

$$\langle Bx, y \rangle_A = \langle x, B^* y \rangle_A \qquad \forall x, y \in \mathbb{R}^n$$

$B$ is positive w.r.t. $\langle \cdot, \cdot \rangle_A$ if $\langle Bx, x \rangle_A \geq 0$
$\forall x \neq 0$

Let $G \in \mathbb{R}^{n \times n}$ with adjoint $G^*$ w.r.t. $\langle \cdot, \cdot \rangle$.
Then, if $B = I - G^* G$ is positive w.r.t. $\langle \cdot, \cdot \rangle$,
it follows $\rho(G) < 1$.

In GS, we have

$$G = I - (D+L)^{-1} A$$

$$B = I - G^*G \quad \text{is a positive matrix w.r.l. } \langle \cdot, \cdot \rangle_A.$$

$$G^* = A^{-1} G^T A =$$

$$= I - A^{-1} A^T \underbrace{(D+R)^{-1}}_{(D+L)^{-T}} A$$

$$= I - (D+R)^{-1} A$$

$$B = I - G^*G = \ldots = (D+R)^{-1} D (D+L)^{-1} A$$

show that $B$ is pos.

$$\langle Bx, x \rangle_A = \langle D^{1/2} (D+L)^{-1} Ax, \; D^{1/2} (D+L)^{-1} Ax \rangle$$

$$= \| D^{1/2} (D+L)^{-1} Ax \|_2^2 \; > 0 \qquad \forall x \neq 0$$

$$\Rightarrow \rho(G) < 1$$

# Gauss-Seidel method

Theorem: Choose $Q = D + L$ ... Gauss-Seidel method. The Gauss-Seidel method converges for any starting point $\boldsymbol{x}_o$ if $A$ is symmetric positive definite (spd).
⤳ board

Notice that $Q = D + L$ is a good choice in terms of computational costs because we can very quickly solve $(D + L)y = z$ for the lower triangular matrix $D + L$ (forward substitution)

# Relaxation methods:

Use linear combination between new and previous iterate:

$$\boldsymbol{x}_{k+1} = \omega \underbrace{(G\boldsymbol{x}_k + c)}_{\boldsymbol{x}'_{k+1}} + (1-\omega)\boldsymbol{x}_k = G_\omega \boldsymbol{x}_k + \omega c,$$

where $\omega > 0$ is a damping/relaxation parameter (sometimes, $\omega > 1$ is used, leading to overrelaxation). Target is to choose $\omega$ such that $\rho(G_\omega)$ is minimal.

Def: A fixed point method $\boldsymbol{x}_{k+1} = G\boldsymbol{x}_k + \boldsymbol{c}$ with $G = G(A)$ is called *symmetrizable* if for any symmetric positive definite (spd) matrix $A$, the matrix $I - G$ is similar to an spd matrix, i.e., there is a regular $W$ such that $W(I - G)W^{-1}$ spd.

board

Symmetrizable:

$\hookrightarrow$ Richardson: $\quad G = I - A$

$$\Rightarrow I - G = A \quad \text{is spd, if A spd}$$

$\hookrightarrow$ Jacobi: $\quad G = I - D^{-1}A$ , sel $W = D^{1/2}$

$$D^{1/2}(I-G)D^{-1/2} = I - I + D^{1/2}D^{-1}AD^{-1/2}$$

$$= D^{-1/2}AD^{-1/2}$$

For sym. schemes

$$x_{k+1} = G x_k + c \quad, \quad G = G(A) \quad, \quad A \text{ spd}$$

$$\sigma(G) \subset (-\infty, 1)$$

$$x_{k+1} = w(G x_k + c) + (1-w)x_k$$

$$= G_w x_k + wc$$

with $\quad G_w = wG + (1-w)I$

want find $0 < w \leq 1$ with minimal $\rho(G_w)$

EV of $G_w$ are

$$\lambda_i (G_w) = w \lambda_i (G) + 1 - w$$

$$= 1 - w (1 - \lambda_i (G))$$

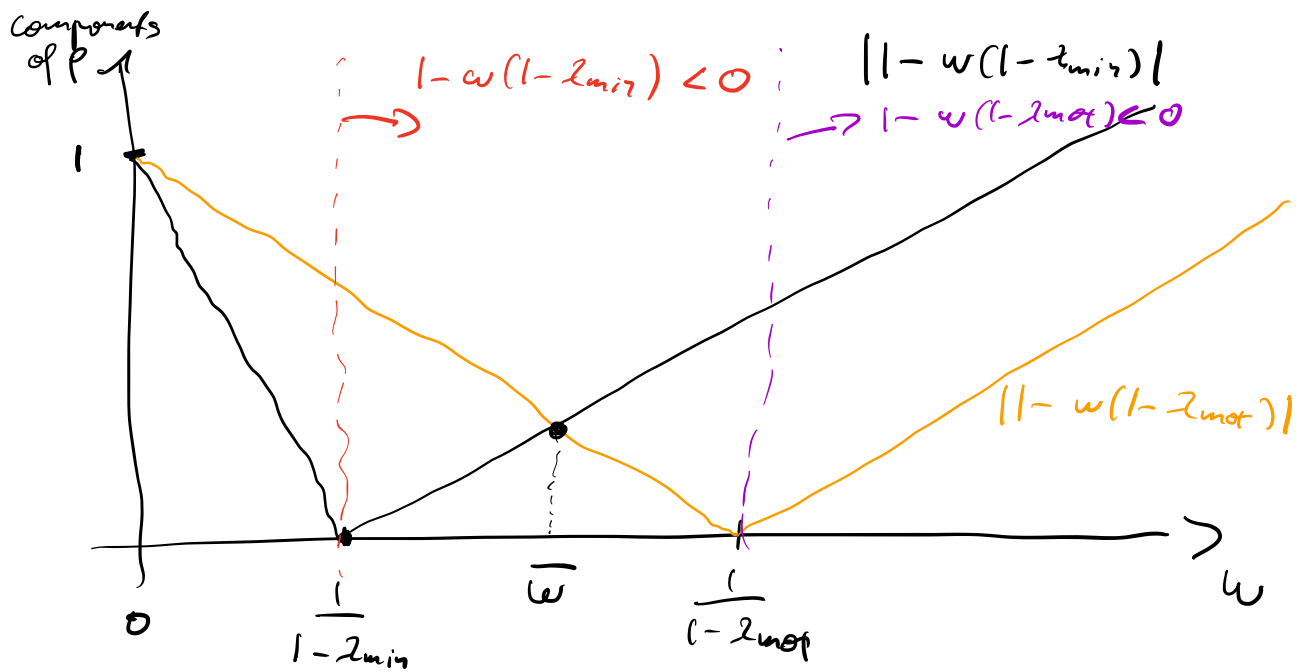$$\lambda_{min} (G) \leq \lambda_{max}(G) < 1$$

$$\lambda_i (G_w) < 1$$

$$\rho (G_w) = max \left\{ | 1 - w (1 - \lambda_{min} (G)) |, \right.$$
$$\left. | 1 - w (1 - \lambda_{max} (G)) | \right\}$$

Let's look at $w_{min} = \dfrac{1}{1 - \lambda_{min} (G)} > 0$

$$| 1 - w_{min} (1 - \lambda_{min} (G)) | = 0$$

Now $w_{max} = \dfrac{1}{1 - \lambda_{max} (G)}$

$$| 1 - w_{max} (1 - \lambda_{max} (G)) | = 0$$

optimal $\bar{\omega}$ satisfies

$$-\left(1 - \bar{\omega}\left(1 - \lambda_{min}(G)\right)\right) = 1 - \bar{\omega}\left(1 - \lambda_{mot}(G)\right)$$

$$\bar{\omega} = \frac{2}{2 - \lambda_{mot}(G) - \lambda_{min}(G)}$$

Richardson : spd matrix $A$ , $G = I - A$

$$\lambda_{min}(G) = 1 - \lambda_{max}(A)$$

$$\lambda_{max}(G) = 1 - \lambda_{min}(A)$$

$$\bar{\omega} = \frac{2}{\lambda_{max}(A) + \lambda_{min}(A)}$$

$$\rho(G_{\bar{\omega}}) = \ldots = \frac{\lambda_{max}(A) - \lambda_{min}(A)}{\lambda_{max}(A) + \lambda_{min}(A)}$$

$$= \ldots = \frac{k_2(A) - 1}{k_2(A) + 1} < 1$$

Finding the optimal damping parameter: $\rightsquigarrow$ board

Finding the optimal damping parameter: $\rightsquigarrow$ board

We obtain that

$$\bar{\omega} = \frac{2}{2 - \lambda_{\mathsf{max}}(G) - \lambda_{\mathsf{min}}(G)}$$

is the optimal damping parameter for symmetrizable iteration methods that minimizes the spectral radius. The spectral radius is

$$\rho(G_{\bar{\omega}}) < 1$$

This means that for a suitable choice $\bar{\omega}$ we can make *any* symmetrizable iteration method convergent for an spd $A$!

Finding the optimal damping parameter: ⇝ board

We obtain that

$$\bar{\omega} = \frac{2}{2 - \lambda_{\max}(G) - \lambda_{\min}(G)}$$

is the optimal damping parameter for symmetrizable iteration methods that minimizes the spectral radius. The spectral radius is

$$\rho(G_{\bar{\omega}}) < 1$$

This means that for a suitable choice $\bar{\omega}$ we can make *any* symmetrizable iteration method convergent for an spd $A$!

Optimal damping parameter for Richardson iteration ⇝ board

# Algorithmic perspective on iterative methods

Richardson iterations

$$\mathbf{x}_{k+1} = (I - A)\mathbf{x}_k + b$$

```
for k = 0, 1, ...
  for i = 0, 1, ..., n-1:
```
$x_{k+1}[i] = x_k[i] + r_k[i]$

where the residual $\mathbf{r}_k$ at iteration $k$ is given by

$$\mathbf{r}_k = \mathbf{b} - A\mathbf{x}_k$$

What would we like to update $\mathbf{x}_k$ with?

# Algorithmic perspective on iterative methods

Richardson iterations

$$\mathbf{x}_{k+1} = (I - A)\mathbf{x}_k + b$$

```
for k = 0, 1, ...
  for i = 0, 1, ..., n-1:   $x_{k+1}[i] = x_k[i] + r_k[i]$
```

where the residual $\mathbf{r}_k$ at iteration $k$ is given by

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$$

What would we like to update $\mathbf{x}_k$ with? We would like to update $\mathbf{x}_k$ in the direction of the error $\mathbf{e}_k = \mathbf{x} - \mathbf{x}_k$ because then

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{e}_k = \mathbf{x}$$

However, we don't have the error $\mathbf{e}_k$ and therefore it is reasonable to use the next best thing which is the residual in many situations

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k = \mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}_k = \mathbf{A}(\mathbf{x} - \mathbf{x}_k) = \mathbf{A}\mathbf{e}_k$$

⤳ Richardson iteration updates $\mathbf{x}_k$ in the direction of the residual $\mathbf{r}_k$

Jacobi iterations

$$\mathbf{x}_{k+1} = (I - D^{-1}A)\mathbf{x}_k + D^{-1}b$$

```
for k = 0, 1, ...
  for i = 0, 1, ..., n-1:    y[i] = 1/a_{ii} r_k[i]
  for i = 0, 1, ..., n-1:    x_{k+1}[i] = x_k[i] + y[i]
```

$$r_k[i] = b_i - a_{i,:} x_k$$

$$a_{i,:} x_{k+1} = b_i$$

- ▶ In every substep $i$ of iteration $k$, an update $y[i]$ is computed and stored
- ▶ Applied immediately, this would lead to the (momentary) disappearance of the $i$-th component of the residual $r_k$
- ▶ Thus, with this current approximation, equation $i$ would be solved exactly—an improvement that would be lost immediately in the following substep for the equation $i + 1$
- ▶ However, the updates of a component are not applied immediately but only at the end of an iteration step (second $i$-loop)

Gauss-Seidel iteration

$$\mathbf{x}_{k+1} = (I - (L+D)^{-1}A)\mathbf{x}_k + (L+D)^{-1}b$$

```
for k = 0, 1, ...
  for i = 0, 1, ..., n-1:
```
$$r_k[i] = b[i] - \sum_{j=1}^{i-1} a_{ij}x_{k+1}[j] - \sum_{j=i}^{n} a_{ij}x_k[j]$$
$$y[i] = 1/a_{ii}r_k[i], \quad x_{k+1}[i] = x_k[i] + y[i]$$

▶ In contrast to Jacobi method, the update is performed immediately
▶ Therefore the new modified values for components $1, ..., i-1$ are already available for updating component $i$

Damping (mostly Jacobi) or over-relaxation (mostly Gauss-Seidel) means to take step lengths different from 1 in the direction of the residual

# The spectral radius of typical iterative matrices

▶ The spectral radius $\rho$ determines convergence and speed; the smaller $\rho$, the faster the error decays. In practice, $\rho$ is often very close to 1 so that even though $\rho < 1$ it takes an unreasonable amount of iterations to get a reasonable answer

▶ An important sample scenario is the discretization of PDEs: It is typical that $\rho$ depends on the dimension $n$ of the matrix $A$, and thus in terms of PDE discretization it depends on the mesh width $h$ of the underlying grid. For example

$$\rho \in \mathcal{O}(1 - h_l^2) = \mathcal{O}(1 - \frac{1}{4^l})$$

with mesh width $h_l = 2^{-l}$ in one dimension ($\mathcal{O}(1 - 16^{-l})$ in two spatial dimensions)

▶ This is a huge disadvantage: Why?

# The spectral radius of typical iterative matrices

▶ The spectral radius $\rho$ determines convergence and speed; the smaller $\rho$, the faster the error decays. In practice, $\rho$ is often very close to 1 so that even though $\rho < 1$ it takes an unreasonable amount of iterations to get a reasonable answer

▶ An important sample scenario is the discretization of PDEs: It is typical that $\rho$ depends on the dimension $n$ of the matrix $A$, and thus in terms of PDE discretization it depends on the mesh width $h$ of the underlying grid. For example

$$\rho \in \mathcal{O}(1 - h_l^2) = \mathcal{O}(1 - \frac{1}{4^l})$$

with mesh width $h_l = 2^{-l}$ in one dimension ($\mathcal{O}(1 - 16^{-l})$ in two spatial dimensions)

▶ This is a huge disadvantage: Why? the finer the grid is (and therefore the more accurate our approximation of the PDE solution should be), the slower these iterative methods get.

Consider the Laplace equation

$$\Delta u(x_1, x_2) = 0\,,$$

in two dimensions and discretize with five-point second-order finite-difference stencil on $N \times N$ grid points ($\rightsquigarrow$ NM2)
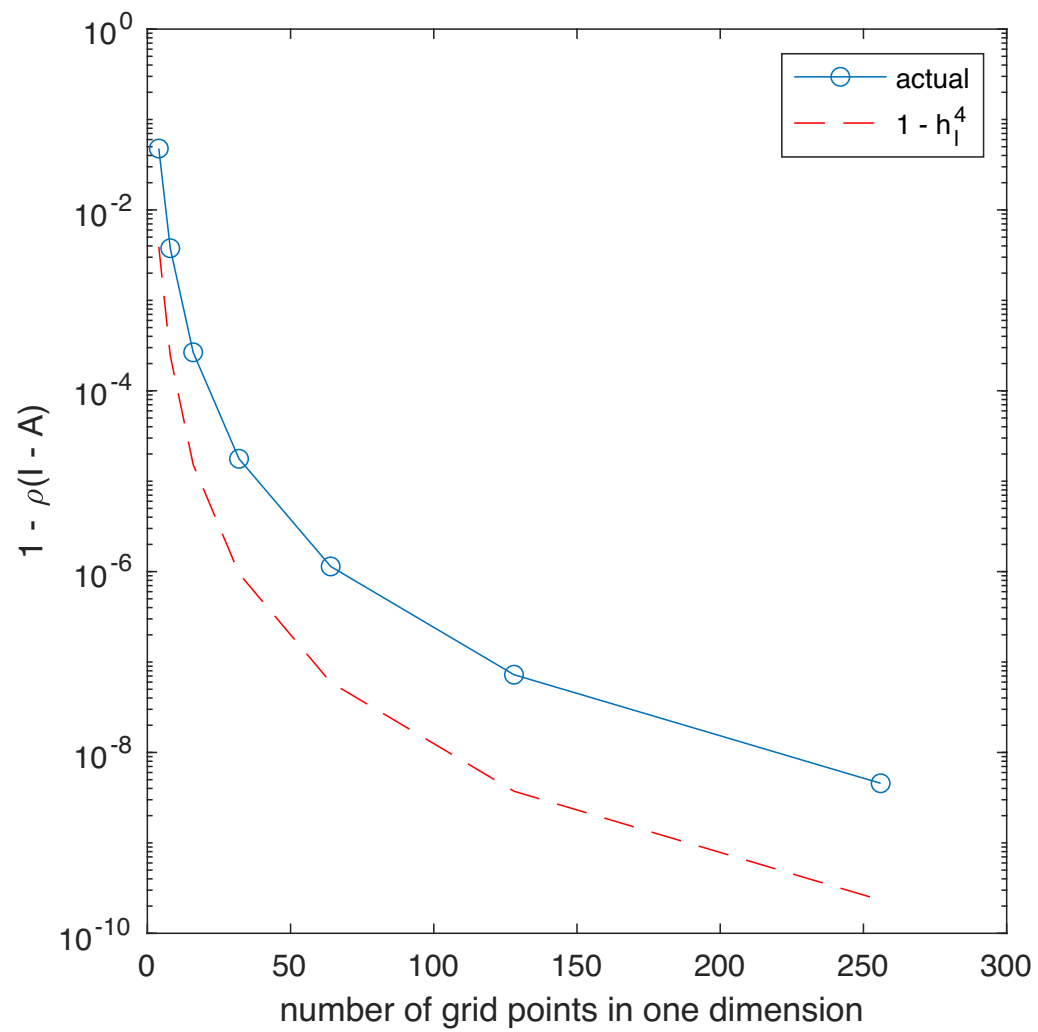
Plot the spectral radius $\rho(I - A)$ (Richardson interpolation) w.r.t. number of grid points

```
1: Nlist = 2.^(2:8);
2: eList = [];
3: for N=Nlist
4:     X = gallery('poisson', N)*(1/N^2);
5:     eList(end + 1) = eigs(speye(N^2) - X, 1);
6: end
```

Running Jacobi on a Poisson matrix with $N = 10$ and $N = 100$ grid points in each dimension:

```
 1: A = gallery('poisson', N)*(1/N^2);
 2: b = randn(N^2, 1);
 3: x = randn(N^2, 1);
 4: xTrue = A\b;
 5: M = speye(N^2) - spdiags(1./diag(A), 0, N^2, N^2)*A;
 6: c = spdiags(1./diag(A), 0, N^2, N^2)*b;
 7:
 8: hist = [];
 9: for iter=1:1000
10:     x = M*x + c;
11:     hist(iter, 1) = norm(x - xTrue)/norm(x);
12:     hist(iter, 2) = norm(A*x - b);
13: end
```
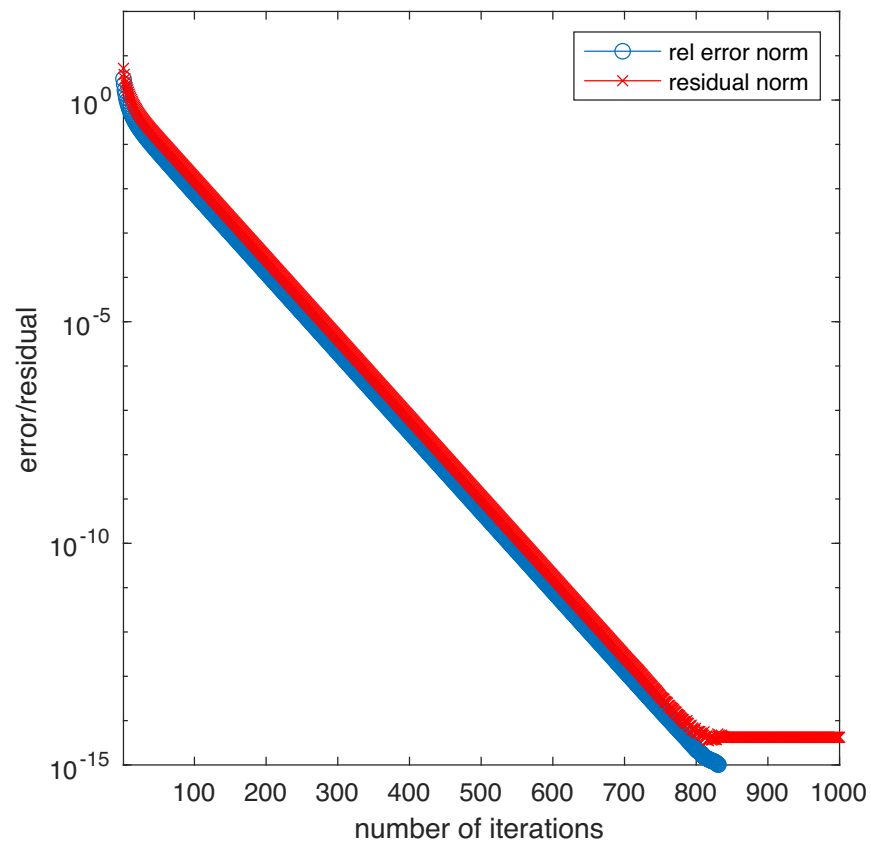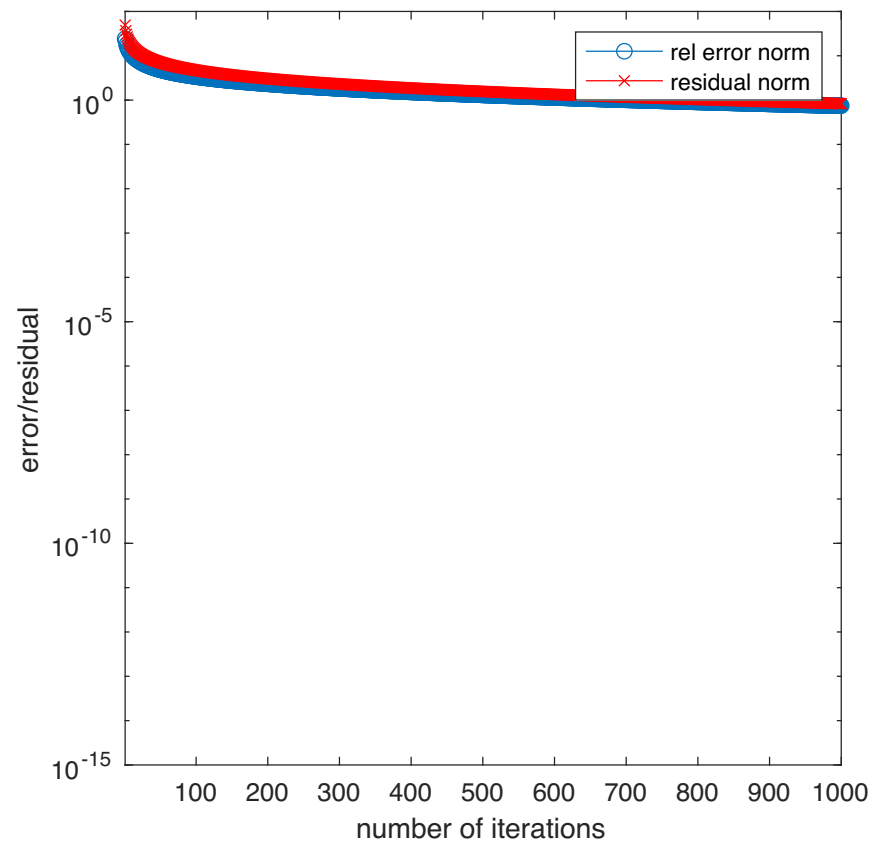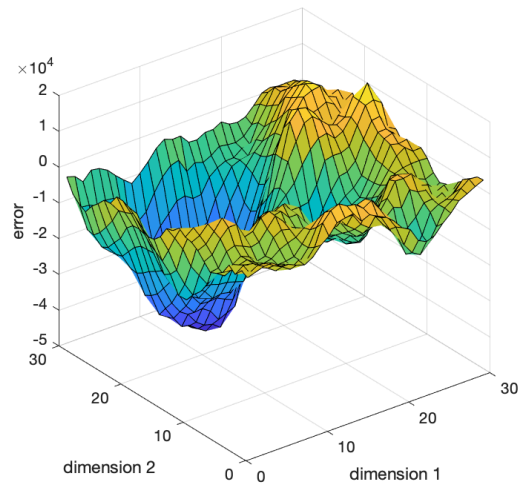
# First few iterations of Jacobi relaxation

```
 1:
 2: N = 30; A = gallery('poisson', N)*(1/N^2);
 3: [X, Y] = meshgrid(linspace(1, N, N), linspace(1, N, N));
 4: b = 10*randn(N^2, 1);
 5: x = randn(N^2, 1);
 6: xTrue = A\b;
 7:
 8: for i=1:5
 9:     x = (speye(N^2) - spdiags(1./diag(A), 0, N^2, N^2)*A)*x + ...
10:                 spdiags(1./diag(A), 0, N^2, N^2)*b;
11: end
```
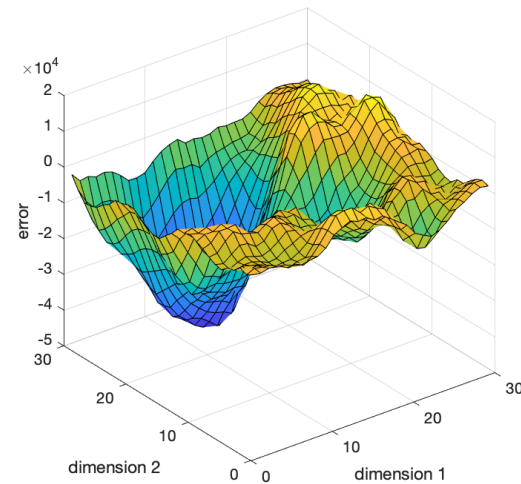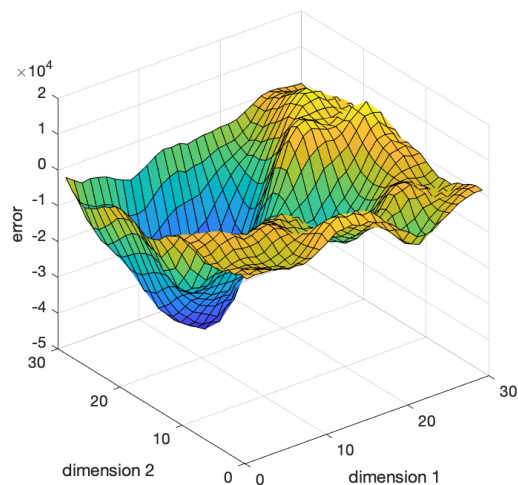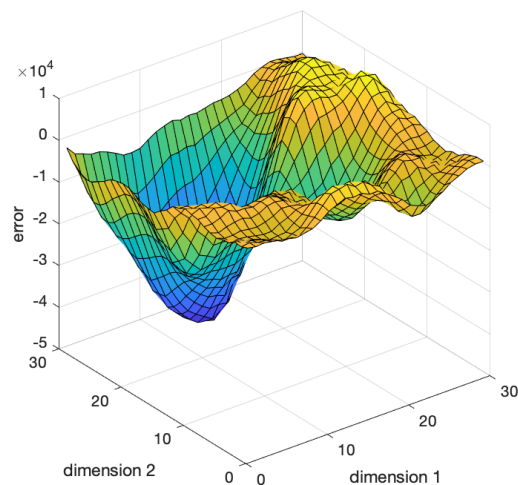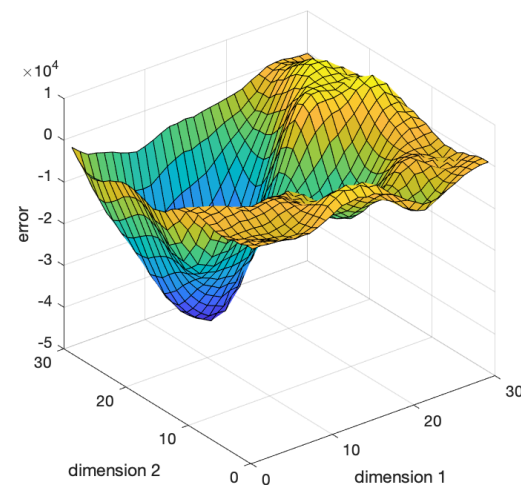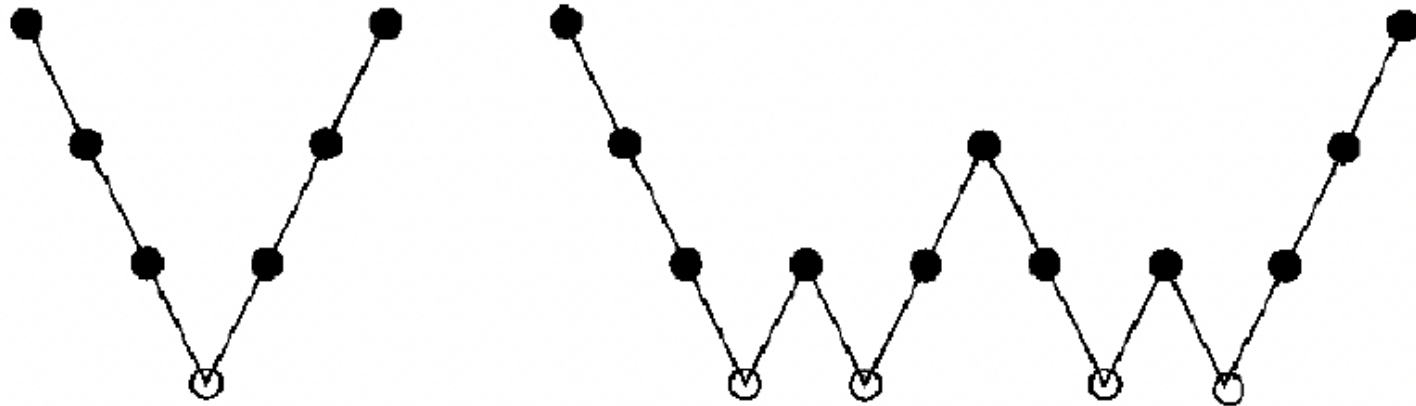
iteration 0

iteration 1

iteration 2

iteration 3

iteration 4

iteration 5

- ▶ Relaxation methods such as Jacobi and Gauss Seidel have a smoothing effect on the error

- ▶ Even if $\rho \approx 1$, only a few iterations are necessary to obtain a smooth error; this means that high-frequency error is reduced very quickly whereas the low-frequency error is reduced slowly

- ▶ Multigrid methods exploit this effect and represent the smoothed error on a coarse grid, where it becomes high-frequency error again, which can be smoothed quickly



⇝ Multigrid methods are among the most efficient solvers for PDE problems ⇝
NM 2 in Spring

Conjugate gradient method

In the following $A$ is symmetric positive definite.

Formulate solving $Ax = b$ as an optimization problem: Define

$$f(x) = \frac{1}{2}x^T A x - b^T x,$$

and minimize

$$\min_{x \in \mathbb{R}^n} f(x)$$

Because $A$ is positive definite, the function $f$ is convex. It is sufficient to look at the gradient

$$\nabla f(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b = Ax - b = -r(x) = 0 \iff Ax = b$$

What is the benefit of this point of view?

In the following $A$ is symmetric positive definite.

Formulate solving $Ax = b$ as an optimization problem: Define

$$f(x) = \frac{1}{2}x^T Ax - b^T x,$$

and minimize

$$\min_{x \in \mathbb{R}^n} f(x)$$

Because $A$ is positive definite, the function $f$ is convex. It is sufficient to look at the gradient

$$\nabla f(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b = Ax - b = -r(x) = 0 \iff Ax = b$$

What is the benefit of this point of view? We now can let loose all what we know about optimization to solve $Ax = b$

Our first try is applying the method of steepest descent in the direction of the negative gradient

$$-\nabla f = r$$

which happens to be the residual

$$r_k = b - Ax_k$$

$$\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k}$$

$$x_{k+1} = x_k + \alpha_k r_k$$

The step length $\alpha_k$ minimizes $f(x_k + \alpha_k r_k)$ as a function of $\alpha_k \rightsquigarrow$ board

$$f(x) = \frac{1}{2}x^T A x - b^T x, \qquad \nabla f(x) = Ax - b$$

$$\min_{\alpha_k} f(x_k + \alpha_k r_k) = \frac{1}{2}(x_k + \alpha_k r_k)^T A (x_k + \alpha_k r_k) -$$
$$b^T (x_k + \alpha_k r_k)$$

$$\frac{\partial}{\partial \alpha_k} f(x_k + \alpha_k r_k) = r_k^T \left[ \nabla f(x_k + \alpha_k r_k) \right]$$

$$= r_k^T \left[ A x_k + A \alpha_k r_k - b \right] \overset{!}{=} 0$$

$$r_k^T A x_k + \alpha_k r_k^T A r_k - r_k^T b \overset{!}{=} 0$$

$$r_k^T \underbrace{(A x_k - b)}_{-r_k} + \alpha_k r_k^T A r_k \overset{!}{=} 0$$

$$\Rightarrow \quad \alpha_k^* = \frac{r_k^T r_k}{r_k^T A r_k}$$

For steepest descent, if $A$ is spd, we obtain

$$\|x^* - x_k\|_A \leq \left( \frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} \right)^k \|x^* - x_0\|_A \,,$$

where $\langle x, y \rangle_A = x^T A y$ and $\| \cdot \|_A = \sqrt{\langle \cdot, \cdot \rangle_A}$.

Proof $\rightsquigarrow$ board

spd $A$:
$$\| x_k - x^* \|_A \leq \left( \frac{k_2(A) - 1}{k_2(A) + 1} \right)^k \| x_k - x_0 \|_A$$

$$x_{k+1}(\alpha) = x_k + \alpha \, r_k \qquad \left( r_k = b - A x_k = A x^* - A x_k \right)$$

$$= x_k + \alpha \, A(x^* - x_k)$$

$$\underbrace{x^* - x_{k+1}(\alpha)}_{\substack{\| \\ e_{k+1}(\alpha)}} = x^* - x_k - \alpha A(x^* - x_k) =$$

$$\left\} \quad = (1 - \alpha A) \underbrace{(x^* - x_k)}_{e_k}$$

$$e_{k+1}(\alpha) = (1 - \alpha A) e_k$$

$$\| e_{k+1}(\alpha) \|_A^2 = e_k^T (1 - \alpha A)^T A (1 - \alpha A) e_k$$

$$\boxed{e_k = \sum_{i=1}^{n} \sigma_i \, z_i}$$