

## Richardson method

$$Q = I$$

$$x^{(k)} = (I - A)x^{(k-1)} + b = x^{(k-1)} + \frac{b - Ax^{(k-1)}}{r^{(k-1)}}$$

converges when  $\|I - A\| < 1$  in some induced matrix norm.

## Jacobi method

$$Q = \begin{bmatrix} a_{11} & & & & \\ & a_{22} & & & \\ & & \dots & & \\ & & & a_{nn} & \end{bmatrix} = \text{diag}(\text{diag}(A)) = \text{diagonal matrix with same entries as diagonal of } A.$$

$$Q^{-1}A = \begin{bmatrix} 1 & a_{12}/a_{11} & a_{13}/a_{11} & \dots & a_{1n}/a_{11} \\ a_{21}/a_{22} & 1 & a_{23}/a_{22} & \dots & a_{2n}/a_{22} \\ \vdots & & & & \\ & & & & a_{n-1n}/a_{n-1n-1} \\ a_{n1}/a_{nn} & \dots & \dots & & 1 \end{bmatrix}$$

Thus:

$$\|I - Q^{-1}A\|_{\infty} = \max_{1 \leq i \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|}$$

Theorem If  $A$  is diagonally dominant, then the Jacobi method's iterates converge to a solution to  $Ax = b$ .

Proof.  $A$  diag dominant  $\Rightarrow |a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, 1 \leq i \leq n$

$$\Rightarrow \|I - Q^{-1}A\|_{\infty} < 1$$

$\Rightarrow$  convergence.

Linear algebra review

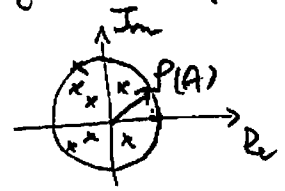
Eigenvalues of  $A$  = roots of characteristic polyn. of  $A$

The characteristic poly. of  $A$  is defined by:

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

Spectral radius:

$$\begin{aligned} \rho(A) &= \max \{ |\lambda| \mid \det(\lambda I - A) = 0 \} \\ &= \text{radius of smallest circle in } \mathbb{C} \text{ containing} \\ &\quad \text{all eigenvalues of } A. \end{aligned}$$



Two matrices  $A$  and  $B$  are similar if there is an invertible matrix  $X$  s.t.

$$AX = XB.$$

note: Two similar matrices have same eigenvalues:

$$\begin{aligned} P_A(\lambda) &= \det(\lambda I - A) = \det(X^{-1}(\lambda I - A)X) \\ &= \det(\lambda I - \underbrace{X^{-1}AX}_B) = P_B(\lambda) \end{aligned}$$

Here we used the facts:

$$\det(AB) = \det(A) \det(B)$$

$$\det(A^{-1}) = \det(A)^{-1}$$

Before we continue studying iterative methods for solving linear systems we need two technical results.

Theorem Every matrix is similar to an upper triangular matrix with arbitrarily small off-diagonal elements.

Proof: Use Schur factorization

$$A = Q T Q^T, \text{ where } Q^T Q = I$$

this is an eigenvalue revealing factorization as the eigenvalues of  $A$  appear on diagonal of  $T$ . we will see more about this factorization when we study eigenvalue algorithms.

$$\text{Let } D = \begin{bmatrix} \varepsilon & & & \\ & \varepsilon^2 & & \\ & & \ddots & \\ & & & \varepsilon^n \end{bmatrix}$$

$$\text{Then } (D^{-1} T D)_{ij} = t_{ij} \varepsilon^{j-i}$$

But  $T = (\nabla)$  thus:

$$t_{ij} = 0 \text{ when } i > j \text{ (elements below diag are zero)}$$

Now for elements above diag: ( $i < j$ ):

$$|t_{ij} \varepsilon^{\overset{?}{j-i}}| \leq \varepsilon |t_{ij}|$$

Thus

$$A = \underbrace{Q D}_X \underbrace{D^{-1} T D}_B \underbrace{D^{-1} Q^T}_{X^{-1}} \quad \text{QED}$$

The second technical result is:

Theorem

$$\rho(A) = \inf_{\|\cdot\|} \|A\|, \text{ where } \inf \text{ is over all induced matrix norms.}$$

This means  $\rho(A) = \|A\|_2 =$  "smallest" induced norm one can choose to measure  $A$ .

proof:

$\rho(A) \leq \inf_{\|\cdot\|} \|A\|$

Let  $\|\cdot\|$  be a vector norm and let  $\lambda, x$  be an eigenvalue/vector of  $A$ , then:

$$\|A\| \|x\| \geq \|Ax\| = \|\lambda x\| = |\lambda| \|x\|$$

$$\Rightarrow |\lambda| \leq \|A\|. \text{ (since eigenvector } x \neq 0)$$

Since this is true for any eigenvalue  $\lambda$  we must have:

$$\rho(A) \leq \|A\|$$

Since this is true for any induced matrix norm  $\|\cdot\|$ , it must be true for the least upper bound as well:

$$\rho(A) \leq \inf_{\|\cdot\|} \|A\|$$

$\rho(A) \geq \inf_{\|\cdot\|} \|A\|$

:  $\forall \epsilon > 0$ , there is a non-singular matrix  $S$  s.t.

$$S^{-1}AS = D + T, \text{ where } \|T\|_\infty \leq \epsilon$$

( $\setminus$ ) ( $\nabla$ )  
(strict)

$$\Rightarrow \|S^{-1}AS\|_\infty = \|D + T\|_\infty \leq \underbrace{\|D\|_\infty}_{=\rho(A)} + \underbrace{\|T\|_\infty}_{\leq \epsilon} \leq \rho(A) + \epsilon$$

To conclude we notice (this is easily shown) that:

$$\|A\|_{\infty}' = \|S^{-1}AS\|_{\infty} \text{ is an induced matrix norm}$$

$$\Rightarrow \inf_{\|\cdot\|} \|A\| \leq \|A\|_{\infty}' \leq \rho(A) + \epsilon$$

$$\Rightarrow \inf_{\|\cdot\|} \|A\| \leq \rho(A) \text{ since inequality above holds } \forall \epsilon > 0. \quad \text{QED}$$

We are now ready for the main result of this section. If you notice all convergence results we have given up to now are simply sufficient conditions for convergence. Here is necessary and sufficient condition for convergence.

First let us write the iterative method in a slightly more general way:

$$x^{(k)} = Gx^{(k-1)} + c \quad (*)$$

The matrix splitting methods (Jacobi etc.) we have seen before fall on this framework:

$$Qx^{(k)} = (Q - A)x^{(k-1)} + b$$
$$x^{(k)} = \underbrace{(I - Q^{-1}A)}_{= G} x^{(k-1)} + \underbrace{Q^{-1}b}_{= c}$$

What is limit of iteration (\*)?

$$x = Gx + c \Rightarrow (I - G)x = c$$

$x = (I - G)^{-1}c$

Theorem The iterative method

$$x^{(k)} = G x^{(k-1)} + c$$

converges to  $(I-G)^{-1}c$  if and only if  $\rho(G) < 1$ .   
 essentially "sharpest" convergence measure

Proof:

- Assume  $\rho(G) < 1$ .

Since  $\rho(G) = \inf_{\|\cdot\|} \|G\|$ , there is an induced matrix norm for which  $\|G\| < 1$ .

$$\begin{aligned} x^{(1)} &= G x^{(0)} + c \\ x^{(2)} &= G^2 x^{(0)} + Gc + c \\ x^{(3)} &= G^3 x^{(0)} + G^2c + Gc + c \\ &\vdots \\ x^{(k)} &= G^k x^{(0)} + \sum_{j=0}^{k-1} G^j c \end{aligned}$$

Note:

$$\|G^k x^{(0)}\| \leq \|G\|^k \|x^{(0)}\| \rightarrow 0 \text{ as } k \rightarrow \infty$$

and

$$\sum_{j=0}^{\infty} G^j c = (I-G)^{-1}c \text{ by Neuman series}$$

Hence  $x^{(k)} \rightarrow (I-G)^{-1}c$ , regardless of initial guess  $x^{(0)}$ .

- Assume  $\rho(G) \geq 1$ : Take  $u \neq 0$  s.t.

$$Gu = \lambda u \text{ and } |\lambda| \geq 1.$$

Then with  $x^{(0)} = 0$  and  $c = u$ :

then: 
$$z^{(k)} = \sum_{i=0}^{k-1} G^i u = \sum_{j=0}^{k-1} \lambda^j u = \begin{cases} \pm k u & \text{if } \lambda = \pm 1 \\ \frac{1 - \lambda^k}{1 - \lambda} & \text{if } |\lambda| \neq 1 \end{cases}$$
 (40)

then  $z^{(k)}$  does not converge as  $k \rightarrow \infty$ . QED.

Corollary The matrix splitting iteration:

$$Q z^{(k)} = (Q - A) z^{(k-1)} + b$$

converges to solution of  $Ax = b$

if and only if

$$\rho(I - Q^{-1}A) < 1.$$

Gauss Seidel method:

$Q$  = lower triangular part of  $A$  (including diag)  
= easy to solve systems (forward substitution)

Theorem If  $A$  is diag. dominant then Gauss-Seidel method converges to sol. of  $Ax = b$ , for any starting vector  $x^{(0)}$ .

Proof: we need to show  $\rho(I - Q^{-1}A) < 1$

Let  $\lambda$  be an eigenvalue of  $I - Q^{-1}A$  and  $x$  a corresp. eigenvector chosen s.t.  $\|x\|_\infty = 1$ .

$$(I - Q^{-1}A)x = \lambda x \Leftrightarrow \overbrace{(Q - A)x}^{\text{strictly upper triangular part}} = \lambda Qx$$

$$\Leftrightarrow - \sum_{j=i+1}^n a_{ij} x_j = \lambda \sum_{j=1}^i a_{ij} x_j$$

for  $i = 1, \dots, n$ .

$$\Rightarrow \lambda a_{ii} x_i = -\lambda \sum_{j=1}^{i-1} a_{ij} x_j - \sum_{j=i+1}^m a_{ij} x_j \quad (i=1, \dots, n) \quad (41)$$

Let  $i$  be s.t.  $|x_i|=1$  (i.e. element realizing  $\max_i |x_i|=1$ )

$$|\lambda| |a_{ii}| \leq |\lambda| \sum_{j=1}^{i-1} |a_{ij}| + \sum_{j=i+1}^m |a_{ij}|$$

$$\Rightarrow |\lambda| \leq \frac{\sum_{j=i+1}^m |a_{ij}|}{|a_{ii}| - \sum_{j=1}^{i-1} |a_{ij}|} < 1 \quad \text{by diagonal dominance.}$$

$\Rightarrow \rho(I - Q^{-1}A) < 1 \Rightarrow$  convergence. QED.

$$\text{Let } A = D - E - F$$

$$\begin{matrix} (\diagdown) & (D) & (\diagup) \\ & \text{strict} & \text{strict} \end{matrix}$$

Successive over relaxation (SOR):

$$Q = \frac{1}{\omega} (D - \omega E), \text{ where } \omega > 0.$$

Iteration becomes:

$$\begin{aligned} \frac{1}{\omega} (D - \omega E) x^{(k)} &= \left( \frac{1}{\omega} (D - \omega E) - D + E + F \right) x^{(k)} + b \\ &= \frac{1}{\omega} \left( (1 - \omega) D + \omega F \right) x^{(k)} + b \end{aligned}$$

$$(D - \omega E) x^{(k)} = ((1 - \omega) D + \omega F) x^{(k)} + \omega b$$

Complexity is similar to Gauss-Seidel, but because of extra parameter  $\omega$  one can get better convergence or even convergence in some systems where G-S does not converge.



Note: When  $\omega = 1$  SOR  $\equiv$  Gauss Seidel.

(42)

Here is another way of looking at SOR:

$$\begin{aligned} (1-\omega)x [ D x^{(k)} &= D x^{(k-1)} ] && \text{identity.} \\ \omega x [ (D-E)x^{(k)} &= F x^{(k-1)} + b ] && \text{G-S.} \end{aligned}$$

$$(D - \omega E) x^{(k)} = ((1-\omega)D + \omega F) x^{(k)} + \omega b$$

Depending on matrix  $A = D - E - F$ , there are ways of choosing  $\omega$  optimally.

Extrapolation methods:

$$\begin{aligned} \gamma x [ x^{(k)} &= G x^{(k-1)} + c ] \\ (1-\gamma)x [ x^{(k)} &= x^{(k-1)} ] \end{aligned}$$

$$\begin{aligned} x^{(k)} &= \gamma (G x^{(k-1)} + c) + (1-\gamma) x^{(k-1)} \\ &= G_\gamma x^{(k-1)} + \gamma c \end{aligned}$$

where  $G_\gamma = \gamma G + (1-\gamma)I$   
= interp between  $G$  and  $I$ .

If iteration converges then:

$$\begin{aligned} x &= \gamma (Gx + c) + (1-\gamma)x \\ \Rightarrow x &= Gx + c \end{aligned}$$

So setting  $G = (I - Q^{-1}A)$  and  $c = Q^{-1}b$  we can solve  $Ax = b$ .

Question How do we choose parameter  $\gamma$  in

$$G_\gamma = \gamma G + (1-\gamma)I$$

s.t.  $P(G_\gamma) = \text{minimum}$

(the smaller  $P(G_\gamma)$  is the faster convergence rate)

If we know

$$a \leq \lambda(G) \leq b$$

why?

then

$\lambda(G_\gamma)$  is between  $\gamma a + (1-\gamma)$   
and  $\gamma b + (1-\gamma)$ .

$$\begin{aligned} \Rightarrow P(G_\gamma) &= \max |\lambda(G_\gamma)| \quad \checkmark \text{ check.} \\ &= \max |\gamma \lambda(G) + 1 - \gamma| \\ &\leq \max_{a \leq \lambda \leq b} |\gamma \lambda + 1 - \gamma|. \end{aligned}$$

This problem has explicit solution:

Theorem If  $\lambda(G) \in [a, b]$ , and  $1 \notin [a, b]$ , then the best choice for  $\gamma$  is:

$$\gamma_* = \frac{2}{2-a-b} \quad \text{and} \quad P(G_\gamma) \leq (1-\gamma) d$$

where  $d = \text{distance from } 1 \text{ to } [a, b]$ .

We shall not prove this result.

However be aware one can think of more complicated functionals of  $G$  and ask the question of which one gives faster convergence

Here is one example that has ties to Chebyshev polynomials which you should have seen in Math 5610, when one is allowed to move the interpolation nodes in order to reduce the Lagrange interpolation error.

Chebyshev acceleration idea

We start with iterative method:

$$x^{(k)} = G x^{(k-1)} + c$$

and ask whether a linear combination of previous iterates

$$u^{(k)} = \sum_{i=0}^k a_i^{(k)} x^{(i)}$$

with coeff s.t.  $\sum_{i=0}^k a_i^{(k)} = 1$  is a better approx

than the  $x^{(0)}, x^{(1)}, \dots, x^{(k)}$ .

$$u^{(k)} = \sum_{i=0}^k a_i^{(k)} x^{(i)}$$

$$\ominus \quad x = \sum_{i=0}^k a_i^{(k)} x$$

---

$$u^{(k)} - x = \sum_{i=0}^k a_i^{(k)} (x^{(i)} - x)$$

$$= \sum_{i=0}^k a_i^{(k)} G^i (x^{(0)} - x)$$

$$= P(G) (x^{(0)} - x) \Rightarrow \|u^{(k)} - x\| \leq \|P(G)\| \|x^{(0)} - x\|$$

where  $P$  is the polynomial:

$$P(z) = \sum_{i=0}^k a_i^{(k)} z^i$$

note:  $P(A)$  means exactly what you would have expected.

$$P(A) = \sum_{i=0}^k a_i^{(k)} A^i, \text{ with } A^0 \equiv I.$$

Now it is not hard to show that

$$\lambda(A) \equiv \text{eigenvalue of } A \Rightarrow P(\lambda(A)) \equiv \text{eigenvalue of } P(A)$$

Thus if we want best approximation we must have:

$$\begin{aligned} P(P(G)) &= \max_{1 \leq i \leq n} |p(\lambda_i(G))| \\ &\leq \max_{z \in S} |p(z)| \end{aligned}$$

where polynomial  $p$  is s.t.

$$p(1) = 1 \quad (\text{since this is equiv. to } \sum_{i=0}^k a_i^{(k)} = 1)$$

$\Rightarrow p$  is a monic polynomial

Morale of story is that we have reduced the problem of finding the best approx in terms of previous iterates to the problem of finding a polynomial with smallest maximum over some set  $S$ .

When  $S = [-1, 1]$  the solution can be given in terms of Chebyshev polynomials:

$$\begin{aligned} T_k(z) &\equiv \text{unique polynomial minimizing } \max_{z \in [-1, 1]} |T_k(z)| \\ &\text{with leading coefficient being } 2^{k-1}. \\ &\equiv \cos(k \cos^{-1} z) \end{aligned}$$

why? Because

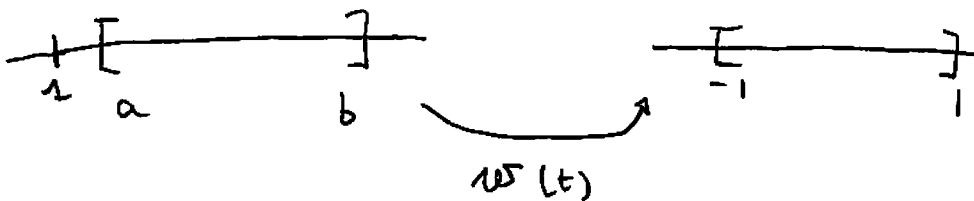
(46)

$$T_k(z) = \min_{\substack{P \in \mathbb{T}_k \\ \text{largest coeff} \\ = 2^{k-1}}} \max_{z \in [-1, 1]} |P(z)|$$

it follows that for  $\beta \in \mathbb{R} \setminus (-1, 1)$ : (to avoid dividing by zero below)  
(all roots of  $T_k$  are in  $(-1, 1)$ )

$$\frac{T_k(z)}{T_k(\beta)} = \min_{\substack{P \in \mathbb{T}_k \\ P(\beta) = 1}} \max_{z \in [-1, 1]} |P(z)|$$

If  $S = [a, b]$  = interval where eigenvalues of  $G$  lie,  
and  $[a, b] \not\subseteq 1$ , then we can find transformation s.t.:



where

$$w(t) = (-1) \frac{t-b}{a-b} + (1) \frac{a-t}{a-b} \quad (\text{Lagrange interp})$$
$$= \frac{a+b-2t}{a-b}$$

then:

$$\frac{T_k(w(t))}{T_k(w(1))} = \min_{\substack{P \in \mathbb{T}_k \\ P(1) = 1}} \max_{z \in [a, b]} |P(z)|$$

↳ can be computed on the fly with a three term recurrence.  
So we only need to remember 2 past iterates.

# CONJUGATE GRADIENT

(47)

Objective: Solve  $Ax = b$  when  $A \in \mathbb{R}^{n \times n}$ , symm. pos. def.

We shall use the notation  $\langle x, y \rangle = x^T y = \sum_{i=1}^n x_i y_i$

- inner product.

Recall properties of an inner product (these hold in general not only in  $\mathbb{R}^n$ )

inner product  $\equiv$  positive symmetric bilinear form

i)  $\langle x, y \rangle = \langle y, x \rangle$  (symm)

ii)  $\langle x, x \rangle \geq 0$  for all  $x \in \mathbb{R}^n$  } (pos.)  
and  $\langle x, x \rangle = 0 \iff x = 0$

iii)  $\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$  (bilin)

Also we have the following property.

iv)

$$\langle x, Ay \rangle = x^T (Ay) = (A^T x)^T y = \langle A^T x, y \rangle$$

Conjugate gradient exploits the following equivalence between linear system and an optimization problem:

(\*)

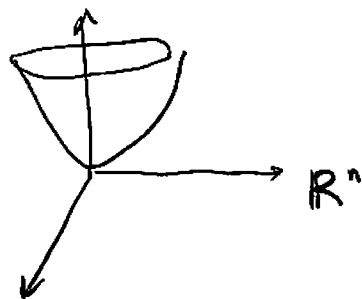
Find  $x$  s.t.  
 $Ax = b$

$\Leftrightarrow$

min  $\frac{1}{2} x^T A x - x^T b = q(x)$   
 $x \in \mathbb{R}^n$

Here  $q(x) \equiv$  quadratic function. When  $x \in \mathbb{R}^n$ :

(48)



Finding solutions of  $Ax=b$  ( $\Leftrightarrow$ ) finding bottom of bowl.

note: the fact that we have an (upward) bowl and not a hyperbowl is due to  $A$  being s.p.d.

( $x^T Ax \geq 0 \forall x \in \mathbb{R}^n$ , so we cannot have branches going to  $-\infty$ )

If you know some vector calculus and first and second order sufficient conditions for having a minimizer then showing (\*) ... is simply:

$$\nabla q(x) = Ax - b$$

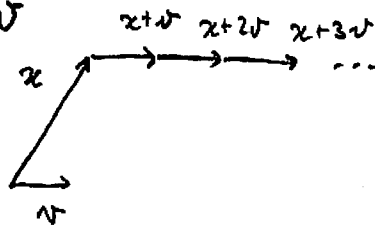
$$\nabla^2 q(x) = A = \text{s.p.d.}$$

Think of  $f: \mathbb{R} \rightarrow \mathbb{R}$ ;  $f \in C^2$   
 $f'(x_*) = 0$  and  $f''(x_*) > 0$   
 $\Rightarrow x_*$  is a minimizer of  $f$

However here is a direct proof of (\*)

proof of (\*):

Look at ray  $x + tv$



$$\begin{aligned}
 q(x+tv) &= \frac{1}{2} (x+tv)^T A (x+tv) - (x+tv)^T b \\
 &= q(x) + t v^T A x - t v^T b + \frac{1}{2} t^2 v^T A v
 \end{aligned}$$

$$\frac{d}{dt} q(x+tv) = v^T (Ax - b) + t v^T A v$$

q as minimum when  $t = t_* = - \frac{v^T (Ax - b)}{v^T A v}$

and:

$$\begin{aligned}
 q(x+t_*v) &= q(x) + t_* (v^T (Ax - b) + \frac{1}{2} t_* v^T A v) \\
 &= q(x) + t_* (v^T (Ax - b) - \frac{1}{2} v^T (Ax - b)) \\
 &= q(x) - \frac{(v^T (Ax - b))^2}{2 v^T A v}
 \end{aligned}$$

What does this mean?

If  $x$  is minimizer of  $q(x)$ , then if we move from  $x$  in any direction  $v$ , then we should not be able to get a smaller value for  $q$ .

$$\Rightarrow v^T (Ax - b) = 0 \text{ for all directions } v$$

$$\Rightarrow Ax - b = 0$$

Conversely: If  $Ax = b$  then

$$q(x+tv) = q(x) + \frac{1}{2} t^2 v^T A v > q(x) \text{ when } v \neq 0.$$