

## Convergence in $\mathbb{R}^n$

24  
1

We say a sequence of vectors  $v^{(k)} \in \mathbb{R}^n$  converges to  $v$  as  $k \rightarrow \infty$  if

$$\lim_{k \rightarrow \infty} \|v^{(k)} - v\| = 0$$

In  $\mathbb{R}^n$  the choice of norm does not matter as all norms are equivalent: (this is true for any finite dimensional space)

Two norms  $\|\cdot\|$  and  $\|\cdot\|'$  are equivalent if there are two constants  $c_1$  and  $c_2 > 0$  s.t.

$$c_1 \|x\| < \|x\|' < c_2 \|x\|$$

Also  $\mathbb{R}^n$  is complete meaning any sequence satisfying the Cauchy criterion converges.

$$\forall \epsilon > 0 \exists N \text{ s.t. } \forall i, j \geq N \|v^{(i)} - v^{(j)}\| < \epsilon$$

## Neumann Series

Let  $A \in \mathbb{R}^{n \times n}$  be such that  $\|A\| < 1$ , then

i)  $I - A$  is invertible

$$\text{ii) } (I - A)^{-1} = \sum_{k=0}^{\infty} A^k$$

Proof: Assume for contradiction that  $I - A$  is singular

$\Rightarrow \exists z \in \mathbb{R}^n \text{ s.t. } (I - A)z = 0$ , take w.l.o.g.  $z$  s.t.  $\|z\| = 1$

then:

$\therefore$

$$1 = \|z\| = \|Az\| < \underbrace{\|A\|}_{\sim} \|z\| < \|z\| = 1 \text{ contradiction!}$$

Now we need to show that

$$\sum_{k=0}^m A^k \rightarrow (I - A)^{-1} \quad (\text{i.e. partial sums converge})$$

↙ telescoping series.

Notice:

$$(I - A) \sum_{k=0}^m A^k = \sum_{k=0}^m A^k - A^{k+1} = A^0 - A^{m+1}$$

$$= I - A^{m+1}$$

Now using properties of induced matrix norm: ( $\|AB\| \leq \|A\| \|B\|$ )

$$\|A^{m+1}\| \leq \|A\|^m \rightarrow 0$$

Thus:

$$\left\| (I - A) \sum_{k=0}^m A^k - I \right\| \leq \|A\|^{m+1} \rightarrow 0 \text{ as } m \rightarrow \infty.$$

QED.

Note: This theorem is a very intuitive generalization of what you already know about geometric series.

$$\frac{1}{1-x} = 1 + x + x^2 + \dots, \text{ when } |x| < 1.$$

Here is the kind of result we can show with Neumann Series:

$$\|(I - A)^{-1}\| \leq \sum_{k=0}^{\infty} \|A^k\| \leq \sum_{k=0}^{\infty} \|A\|^k = \frac{1}{1 - \|A\|}.$$

Note: This theorem is very useful for the theory, however it is seldom used in practice. (there are faster iterative methods for solving  $Ax = b$ ).

Here is an easy but handy generalization of  
Neumann Series:

(29)

3

Theorem

If  $A$  and  $B$  are matrices s.t.  $\|I - AB\| < 1$  then  
 $A, B$  are invertible and

$$A^{-1} = B \sum_{k=0}^{\infty} (I - AB)^k$$

$$B^{-1} = \left[ \sum_{k=0}^{\infty} (I - AB)^k \right] A$$

Proof: By Neumann series:

$I - (I - AB) = AB$  is invertible and

$$(AB)^{-1} = \sum_{k=0}^{\infty} (I - AB)^k$$

$$\Rightarrow A^{-1} = B B^{-1} A^{-1} = B \sum_{k=0}^{\infty} (I - AB)^k$$

$$B^{-1} = B^{-1} A^{-1} A = \left[ \sum_{k=0}^{\infty} (I - AB)^k \right] A$$

## Iterative refinement

Let  $x^{(0)}$  be an approx. sol to

$$Ax = b.$$

then:

$$x = x^{(0)} + A^{-1}(b - Ax^{(0)})$$

We often call

$$r^{(0)} = b - Ax^{(0)} = \text{residual vector}$$

$$e^{(0)} = A^{-1}r^{(0)} = \text{error vector.}$$

## Algorithm

$x^{(0)}$  given

for  $k = 0, 1, \dots$

$$r^{(k)} = b - Ax^{(k)}$$

$$Ae^{(k)} = r^{(k)}$$

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

usually inexact solve

Solution can be improved even if the solves are inexact for example:

→ LU factorization is stopped prematurely ( $\equiv$  Incomplete LU factorization)

→  $Ae^{(k)} = r^{(k)}$  is solved with an iterative method

→  $A$  is not known exactly.

To be more precise about what we mean by "Solving approximately"

Let  $B \equiv$  "approximate" inverse of  $A$

$\approx A^{-1}$  (in some sense we shall see)

The iterates are :

$$x^{(0)} = Bb$$

$$x^{(k+1)} = x^{(k)} + B(b - Ax^{(k)}), k \geq 0.$$

Looking at a few iterates we see a pattern emerge.

$$x^{(0)} = Bb$$

$$x^{(1)} = x^{(0)} + B(b - Ax^{(0)})$$

$$= Bb + B(b - A B b)$$

$$= Bb + B(I - AB)b$$

$$x^{(2)} = x^{(1)} + B(b - Ax^{(1)})$$

$$= Bb + B(I - AB)b + B(b - A(Bb + B(I - AB)b))$$

$$= Bb + B(I - AB)b + B((I - AB)b - AB(I - AB)b)$$

$$= Bb + B(I - AB)b + B(I - AB)^2 b$$

$$\boxed{x^{(m)} = B \sum_{k=0}^m (I - AB)^k b}$$

Theorem (Iterative refinement)

Iterates from iterative refinement are:

$$x^{(m)} = B \sum_{k=0}^m (I - AB)^k b \quad (m \geq 0)$$

and because of the generalized Neumann series:

$$x^{(m)} \rightarrow x \text{ as } m \rightarrow \infty \text{ when } \|I - AB\| < 1.$$

Proof : Can be done by induction on  $m$ .

## Solving equations with iterative methods

Idea: Instead of having a direct method which finds solution in finitely many steps, use an iterative method to find successive approx that converge to sol.

### Motivational example

$$\begin{bmatrix} 7 & -6 \\ -8 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ -4 \end{bmatrix}$$

Jacobi method: solve i-th equation for i-th unknown:

$$\begin{cases} x_1^{(k)} = \frac{3}{7} + \frac{6}{7} x_2^{(k-1)} \\ x_2^{(k)} = -\frac{4}{9} + \frac{8}{9} x_1^{(k-1)} \end{cases}$$

Gauss Seidel method: use improved value of  $x_1^{(k)}$  in second equation.

$$\begin{cases} x_1^{(k)} = \frac{3}{7} + \frac{6}{7} x_2^{(k-1)} \\ x_2^{(k)} = -\frac{4}{9} + \frac{8}{9} x_1^{(k)} \end{cases} \quad \text{converges faster!}$$

Note: convergence for these methods depends on initial guess  $(x_1^{(0)}, x_2^{(0)})$

## Matrix splitting

Consider the system  $Ax = b$ .

Introduce a non-singular matrix  $Q$  and rewrite system as:

$$Qx = (Q - A)x + b$$

Get iteration:

$$Qx^{(k)} = (Q - A)x^{(k-1)} + b, \quad k \geq 1, \quad x^{(0)} \text{ given}$$

$Q$  should have the desirable properties:

- ①  $x^{(k)} \rightarrow x$  where  $x$  is a sol. of  $Ax = b$
- ② Solving systems with  $Q$  is cheap.

If  $x^{(k)} \rightarrow x$  then the limit  $x$  must satisfy:

$$Qx = (Q - A)x + b \Rightarrow Ax = b$$

Assuming  $A$  &  $Q$  are non-singular:

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

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

$$x^{(k)} - x = (I - Q^{-1}A)(x^{(k-1)} - x)$$

$$\Rightarrow \|x^{(k)} - x\| \leq \|I - Q^{-1}A\| \|x^{(k-1)} - x\|$$

⋮

$$\leq \|I - Q^{-1}A\|^k \|x^{(0)} - x\|$$

Thus when  $\|I - Q^{-1}A\| < 1$  we have  $x^{(k)} \rightarrow x$ , where  $Ax = b$ .

Theorem If  $\|I - Q^{-1}A\| < 1$  for some induced matrix norm

then seq.

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

converges to sol. of  $Ax = b$  for any starting vector  $x^{(0)}$

## Richardson method

(34)

$$Q = I$$

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

8

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

## Jacobi method

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

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

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.

$$A \text{ diag dominant} \Rightarrow |a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad 1 \leq i \leq n$$

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

$\Rightarrow$  convergence.

## Linear algebra review

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

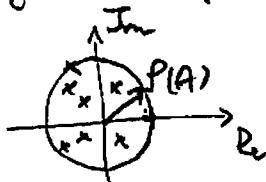
The characteristic poly. of  $A$  is defined by:

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

Spectral radius:

$$r(A) = \max \{ |\lambda| \mid \det(\lambda I - A) = 0 \}$$

= radius of smallest circle in  $\mathbb{C}$  containing all eigenvalues of  $A$ .



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.

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

Then

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

But  $T = (\nabla)$  thus:

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

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

$$|t_{ij} \varepsilon^{\frac{j-i}{2}}| \leq \varepsilon |t_{ij}|$$

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

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

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:

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

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

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:

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

$\rho(A) \geq \inf_{\|\cdot\|} \|A\|$  :  $\forall \epsilon > 0$ , there is a non-singular matrix S s.t.

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

(\wedge) ( $\nabla$ )

(strict)

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

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

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

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

$$\Rightarrow \inf_{\|x\|} \|Ax\| \leq \rho(A) \text{ since the above holds } \forall \epsilon > 0.$$

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 in this framework:

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

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

What is limit of iteration (\*)?

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

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

Theorem The iterative method

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

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

Proof:

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

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

$$x^{(1)} = Gx^{(0)} + c$$

$$x^{(2)} = G^2x^{(0)} + Gc + c$$

$$x^{(3)} = G^3x^{(0)} + G^2c + Gc + c$$

⋮

$$x^{(k)} = G^kx^{(0)} + \sum_{j=0}^{k-1} G^j c$$

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  $\mu \neq 0$  s.t.

$$G\mu = \lambda\mu \text{ and } |\lambda| \geq 1.$$

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

then:

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

thus  $z^{(k)}$  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 correspond. eigenvector chosen s.t.  $\|x\|_\infty = 1$ .

$$(I - Q^{-1}A)x = \lambda x \Leftrightarrow \underbrace{(Q - A)x}_{\substack{\text{strictly upper} \\ \text{triangular}} \atop \text{def.}} = \lambda Qx$$

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

$i = 1, \dots, n$

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

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

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

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

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

let  $A = D - E - F$

$$( \backslash ) \quad (D) \quad (\nabla)$$

strict      strict

### Successive Over relaxation (SOR):

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

Iteration becomes:

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

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

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)

16

Here is another way of looking at SOR:

$$(1-\omega)x \left[ D\mathbf{x}^{(k)} = D\mathbf{x}^{(k-1)} \right] \quad \text{identity.}$$

$$\omega x \left[ (D-E)\mathbf{x}^{(k)} = F\mathbf{x}^{(k-1)} + b \right] \quad G-S.$$

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

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

Extrapolation methods:

$$\gamma x \left[ \mathbf{x}^{(k)} = G\mathbf{x}^{(k-1)} + c \right]$$

$$(1-\gamma)x \left[ \mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} \right]$$

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

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

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

If iteration converges then:

$$\mathbf{x} = \gamma (G\mathbf{x} + c) + (1-\gamma)\mathbf{x}$$

$$\Rightarrow \mathbf{x} = G\mathbf{x} + c$$

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

Question How do we choose parameter  $\gamma$  in

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

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

(the smaller  $S(G_\gamma)$  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  $\gamma b + (1-\gamma)$ .

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

This problem has explicit solution:

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

$$\gamma_* = \frac{2}{2-a-b} \quad \text{and} \quad S(G_{\gamma_*}) \leq 1 - 18d$$

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.  $\sum_{i=0}^k a_i^{(k)} = 1$  do a better approx

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

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

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


---

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

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

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

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.

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

Now it is not hard to show that

$$\lambda(A) = \text{eigenvalue of } A \Rightarrow p(\lambda(A)) = \text{eigenvalue of } p(A)$$

Thus if we want best approximation we must have :

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

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

Moral 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 :

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

$$= \cos(k \cos^{-1} z)$$

why? Because

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

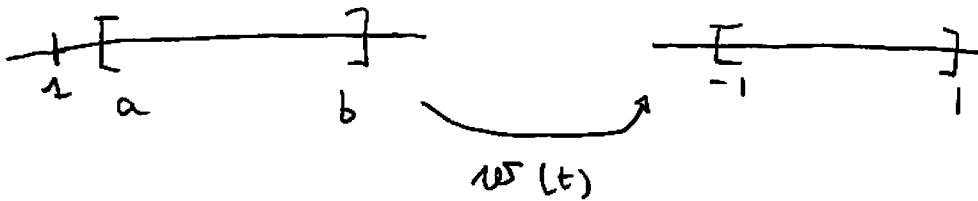
largest coeff  
 $= 2^{k-1}$

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

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

$P(\beta) = 1$

If  $S = [a, b] =$  interval where eigenvalues of  $G$  lie,  
and  $[a, b] \not\subset [-1, 1]$ , 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:

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

$P(1) = 1$

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

## CONJUGATE GRADIENT

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  $\langle x, x \rangle = 0 \Leftrightarrow x = 0$
- iii)  $\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$  (bilinear)

Also we have the following property.

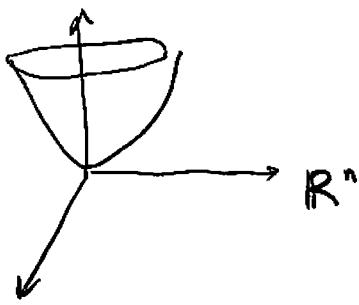
(iv)

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

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

$$(*) \quad \boxed{\begin{array}{|l|} \hline \text{Find } x \text{ s.t.} \\ Ax = b \\ \hline \end{array}} \Leftrightarrow \boxed{\begin{array}{|l|} \hline \min_{x \in \mathbb{R}^n} \frac{1}{2} x^T A x - x^T b = q(x) \\ \hline \end{array}}$$

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



Finding solutions of  $Ax = b \Rightarrow$  finding bottom of bowl.

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

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

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

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

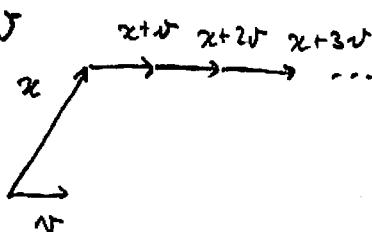
$$\nabla^2 q(x) = A = \text{s.p.d.}$$

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

However here is a direct proof of (\*)

Proof of (\*):

Look at ray  $x + t v$



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

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

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

and:

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

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  $q$ .

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

$$\Rightarrow Ax - b = 0$$

Conversely: If  $Ax = b$  then

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

This gives idea on how to find solution to  $Ax = b$ .

$$x^{(k+1)} = x^{(k)} + t_k v^{(k)}$$

↑ "search direction"  
step size

If  $v^{(k)}$  is given then our calculation reveals that the best possible step size we can take (i.e. the one that reduces the most value of objective function  $q(x)$ ) is:

$$t_k = \frac{\langle v^{(k)}, b - Ax^{(k)} \rangle}{\langle v^{(k)}, Av^{(k)} \rangle} = \frac{\langle v^{(k)}, r^{(k)} \rangle}{\langle v^{(k)}, Av^{(k)} \rangle}$$

where  $r^{(k)} = b - Ax^{(k)}$  = residual at step  $k$ .

One example of such method is:

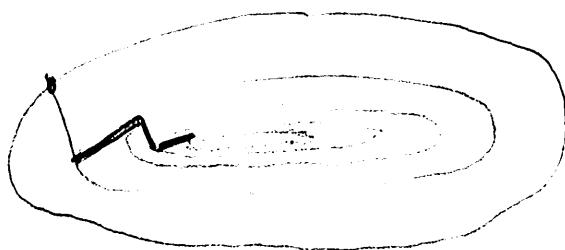
Steepest descent where  $v^{(k)} = r^{(k)} = -\nabla q(x^{(k)})$

Algorithm.

for  $k = 1, 2, \dots$

$$\left| \begin{array}{l} v = b - Ax \\ t = \frac{(v, v)}{(v, Av)} \\ x = x + t v \end{array} \right.$$

This is not a very efficient method. It has slow, typically "zig-zagging" convergence to solution:



Level sets of  $q(x)$  are ellipses when  $x \in \mathbb{R}^2$ .

Note:  $r^{(k)}$  is  $\perp$  to level sets of  $q(x)$ .

## Conjugate Gradient method (Hestenes & Stiefel, 1952)

$x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k, p_k$  is determined s.t.  
 $x_{k+1}$  solves

$$(1) \quad \min \frac{1}{2} x^T A x - b^T x$$

$x \in \mathcal{X}_0 + \text{span}\{p_0, \dots, p_{k-1}, r_k\}$

$$x = \hat{x} + x_0$$

$$\Leftrightarrow (2) \quad \min \frac{1}{2} \hat{x}^T A \hat{x} - \hat{x}^T b$$

$\hat{x} \in \text{span}\{p_0, \dots, p_{k-1}, r_k\}$

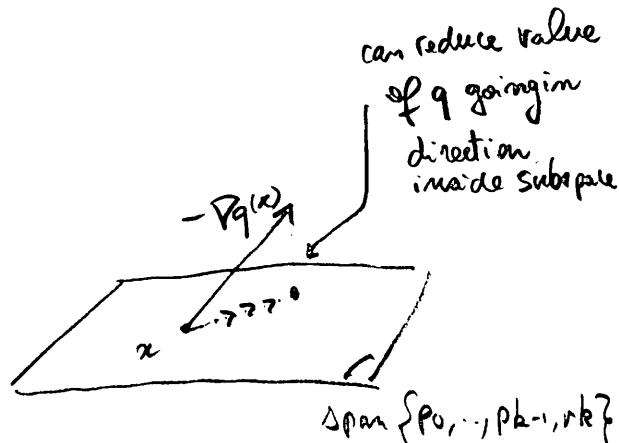
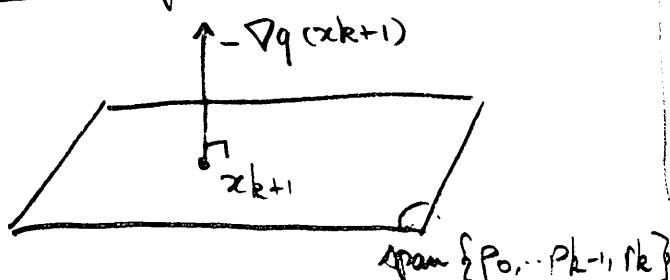
Recall:  $\text{span}\{u_1, u_2, \dots, u_n\} = \left\{ \sum_{i=1}^n d_i u_i \mid d_i \in \mathbb{R} \text{ for } i=1, \dots, n \right\}$   
= set of all possible linear combinations  
of vectors  $u_1, \dots, u_n$   
= linear span of family  $\{u_1, \dots, u_n\}$ .

It can be shown that  $\hat{x}_{k+1}$  solves (2) iff:

$$(A \hat{x}_{k+1} - b)^T v = 0 \quad \forall v \in \text{span}\{p_0, \dots, p_{k-1}, r_k\}$$

$$\Leftrightarrow \underbrace{(A x_{k+1} - b)^T v}_{= -\nabla q(x_{k+1})} = 0 \quad " \quad " \quad "$$

Intuitively:



can reduce value  
of  $q$  going in  
direction  
inside sub-space

OPTIMAL: no direction of descent in  $\text{span}\{\dots\}$

SUB-OPTIMAL

Note:  $v$  is called a descent direction for  $q$  if

$$v^T \nabla q(x) < 0.$$

Looking at Taylor expansion around  $x$ :

$$q(x + t v) = q(x) + \underbrace{t \nabla q(x)^T v}_{< 0 \text{ when } t > 0} + o(t)$$

and  $\nabla q(x)^T v < 0$ .

$\Rightarrow$  means there is a step length  $t > 0$  for which  $q(x + t v) < q(x)$ .

Now look at previous step  $\hat{x}_k$ , which must solve:

$$\min_{\hat{x}} \frac{1}{2} \hat{x}^T A \hat{x} - \hat{x}^T r_0 \Leftrightarrow \begin{aligned} v^T (b - Ax_k) &= 0, \\ \forall v \in \text{span}\{p_0, \dots, p_{k-1}\} \end{aligned}$$

$$\Leftrightarrow P_j^T (b - Ax_k) = 0$$

$$\text{for } j = 0, \dots, k-1.$$

We also have:

$$x_{k+1} = x_k + \alpha_k p_k \quad = 0$$

$$0 = (Ax_{k+1} - b)^T p_j = (\cancel{Ax_k - b})^T p_j + \alpha_k p_k^T A p_j$$

$$\text{for } j = 0, \dots, k-1$$

$$\Rightarrow \boxed{p_k \text{ is } A\text{-orthogonal to previous } k \text{ directions}}$$

A-orthogonal means orthogonality w.r.t. inner product:

$$\langle u, v \rangle_A = \langle u, Av \rangle = u^T A v.$$

It is not hard to check  $\langle \cdot, \cdot \rangle_A$  is indeed an inner product when A is s.p.d.

So we can obtain  $p_k$  be A-orthogonalizing family of vectors  $\{r_0, r_1, \dots, r_{k-1}, r_k\}$ .

Use Gram-Schmidt orthogonalization.

$$P_0 = r_0$$

$$P_k = r_k - \sum_{j=0}^{k-1} \underbrace{\frac{r_k^T A P_j}{P_j^T A P_j} P_j}_{= \perp \text{ proj w.r.t. } \langle \cdot, \cdot \rangle_A \text{ inner prod of } r_k \text{ along } P_j}$$

Now take the optimal step size along this direction:

$$\alpha_k = \frac{P_k^T (b - Ax_k)}{P_k^T A P_k} = \frac{P_k^T r_k}{P_k^T A P_k}$$

Gram Schmidt + this choice of step size is essentially a preliminary version of CG.

However in this version:

- cost of iteration grows linearly with iteration number.
- storage needed " " " " " .

Fortunately the Gram-Schmidt orthogonalization takes a simpler form if we look closely at subspaces we use.

The subspace that CG uses to look for new direction is: (54)

28

$$\begin{aligned}\text{Span} \{ p_0, r, \dots, p_{k-1}, r_k \} &= \text{Span} \{ p_0, \dots, p_{k-1}, p_k \} \\ &= \text{Span} \{ r_0, \dots, r_{k-1}, r_k \} \\ &= K_{k+1}(A, r_0) \\ &= \underline{\text{Krylov subspace}}\end{aligned}$$

why?

$$\begin{aligned}r_0 &= b - Ax_0 \in K_1(A, r_0) \\ \text{If } r_{k-1} &= b - Ax_{k-1} \in K_k(A, r_0) \text{ then} \\ r_k &= b - Ax_k = b - A(x_{k-1} + \alpha p_{k-1}) \\ &= \underbrace{b - Ax_{k-1}}_{= r_{k-1} \in K_k} - \alpha \underbrace{Ap_{k-1}}_{\in K_k} \\ &\quad \underbrace{\in}_{\in K_{k+1}}\end{aligned}$$

$\Rightarrow r_k \in K_{k+1}$

Conjugate Gradient forms part of a large family of iterative solvers called Krylov Subspace methods.

Why did we introduce these subspaces?

Recall optimality conditions, now written with Krylov subspaces.

$$r_{k+1}^T v = 0, \quad \forall v \in K_{k+1}(A, r_0)$$

$$r_k^T v = 0, \quad \forall v \in K_k(A, r_0)$$

Take  $p_i \in K_{k-1}(A, r_0) \Rightarrow Ap_i \in K_k(A, r_0)$

$$\Rightarrow r_k^T Ap_i = 0, \quad \text{for } i = 0, \dots, k-2$$

This orthogonality greatly simplifies Gram Schmidt and reduces sum to only one term:

$$p_k = r_k - \boxed{\frac{r_k^T A p_{k-1}}{p_{k-1}^T A p_{k-1}}} p_{k-1}$$

"  $\beta_k$

- Storage and cost of iteration remains the same regardless of  $k$ .
- One needs roughly only 5 vectors of length  $n$ .
- One requires only knowledge of action of  $A$  on a vector.  
 $\Rightarrow$  no need to know entries of  $A$  as in direct methods.  
 $A$  could be even specified via a "black-box" code.

In the classical formulation of CG  $\alpha_k$  and  $\beta_k$  take different forms, which can be obtained by applying optimality conditions:

$$\alpha_k = \frac{\|r_k\|^2}{p_k^T A p_k}$$

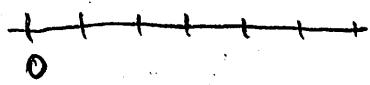
$$\beta_k = \frac{\|r_{k+1}\|^2}{\|r_k\|^2}$$

### Convergence of CG:

- If  $A$  is s.p.d. CG converges in at most  $n$  steps, where  $n$  = dimension of  $A$ .
- In general convergence of CG depends on how many "clusters" of eigenvalues does  $A$  have.

# Spectrum of A (all eigenvalues)

(56)  
30



Slow convergence

(S1)



fast convergence in roughly as many iterations as there are clusters (here n iter)

(S2)

## Preconditioning

To speed up convergence of CG (i.e. going from (S1) to (S2)) one solves system:

$$M^{-1}A\bar{x} = M^{-1}\bar{b}, \text{ where } M \text{ is invertible}$$

$$\Leftrightarrow A\bar{x} = \bar{b}$$

$M$  = (left) preconditioner

= approximation of  $A$  that is easy to invert.

One can also think of right preconditioners:

$$\underbrace{AM^{-1}y}_{=x} = b \quad \Leftrightarrow A\bar{x} = \bar{b}$$

For symmetric systems; we don't want to lose symmetry,  
so letting

$$M = CC^T, M \text{ invertible:}$$

$$Ax = b \quad \Leftrightarrow \underbrace{C^{-1}AC^{-T}C^T x}_{\text{s.p.d. if } A \text{ s.p.d.}} = C^{-1}b$$

Note: there is no single way of choosing a preconditioner.  
A preconditioner that works for a particular problem may not work as well for another.

Here are some preconditioning techniques that work  
( see Trefethen and Bau for a broader overview )

- incomplete LU or Cholesky factorization: drop elements that are below a threshold to get sparse LU (L) factor that give cheap systems to solve
- Jacobi : use diagonal or block-diagonal of A
- Discretization based : use solution to a coarse grid problem, a constant coeff problem, a periodic problem or solve problem in alternating directions (we will see this more closely when we get to FDI methods)

The final algorithm is as follows:

### CONJUGATE GRADIENT (unpreconditioned)

$$x_0 = \text{given}$$

$$r_0 = b - Ax_0$$

$$p_0 = r_0$$

for  $n = 1, 2, 3, \dots$

$$\alpha_n = \frac{r_{n-1}^T r_{n-1}}{p_{n-1}^T A p_{n-1}}$$

$$x_n = x_{n-1} + \alpha_n p_{n-1}$$

$$r_n = r_{n-1} - \alpha_n A p_{n-1} \quad \leftarrow \text{check convergence here by looking at } \|r_n\|$$

$$\beta_n = \frac{r_n^T r_n}{r_{n-1}^T r_{n-1}}$$

$$p_n = r_n + \beta_n p_{n-1}$$

and the preconditioned version is:

## PRECONDITIONED CONJUGATE GRADIENT

$x_0 = \text{given}$

$M = \text{preconditioner} = \text{given}$

$$r_0 = b - Ax_0$$

Solve  $My_0 = r_0$  for  $y_0$

$$p_0 = y_0$$

for  $n = 1, 2, 3, \dots$

$$x_n = \frac{r_{n-1}^T y_{n-1}}{p_{n-1}^T A p_{n-1}}$$

$$x_n = x_{n-1} + \alpha_n p_{n-1}$$

$$r_n = r_{n-1} - \alpha_n A p_{n-1} \quad \leftarrow \text{check convergence here.}$$

Solve  $My_n = r_n$  for  $r_n$ , should be  $y_n$

$$\beta_n = \frac{y_n^T r_n}{y_{n-1}^T r_{n-1}}$$

$$p_n = y_n + \beta_n p_{n-1}$$

Note: This is equivalent to applying CG to system

$$C^T A C^{-T} C^T x = C^{-1} b$$

where  $M = C C^T = \text{preconditioner.}$

We will not see equivalence but notice:

$$y_n = M^{-1} r_n \Rightarrow r_n^T y_n = r_n^T M^{-1} r_n = \underbrace{(C^{-1} r_n)^T (C^{-1} r_n)}_{\text{prec. residual}}$$

$$= (C^{-1} r_n)^T (C^{-1} r_n)$$