

Operator (matrix) norms (induced by the corresponding vector norms)

$$\|A\|_{p,q} = \max_{\|x\|_q \leq 1} \|Ax\|_p = \max_x \frac{\|Ax\|_p}{\|x\|_q}$$

$$\|A\|_p \cong \|A\|_{p,p}$$

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| \quad (\text{max column sum})$$

$$\|A\|_2 = \sqrt{\max_{1 \leq i \leq n} \lambda_i(A^T A)} \quad (\text{Spectral norm; used most commonly})$$

$\lambda_i(M)$ - eigenvalue of M

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (\text{max row sum})$$

$$\|A\|_F = \sqrt{\sum_{i,j=1}^n |a_{ij}|^2} \quad (\text{Frobenius "norm"; not really a norm...})$$

* Ill-Conditioning and Amplification of Errors

Denote

x - the actual (unknown) solution

\bar{x} - the computed solution

$r = b - A\bar{x}$ - the residual

$e = x - \bar{x}$ - the solution error

We have: $Ax = b \Rightarrow x = A^{-1}b$

Using norms:

$$\left. \begin{aligned} \|b\| = \|Ax\| &\leq \|A\| \|x\| \\ \|x\| = \|A^{-1} b\| &\leq \|A^{-1}\| \|b\| \end{aligned} \right\} \Rightarrow \frac{\|b\|}{\|A\|} \leq \|x\| \leq \|A^{-1}\| \|b\| \quad (*)$$

Likewise: $r = b - Ax = Ax - Ax = A(x - \bar{x}) = Ae$

$$\frac{\|r\|}{\|A\|} \leq \|e\| \leq \|A^{-1}\| \|r\| \quad (**)$$

Combining (*) and (**)

$$\frac{\|r\|}{\|A\|} \leq \|e\| / \|x\|$$

$$\left. \begin{aligned} \frac{1}{\|A\|} \frac{\|r\|}{\|x\|} &\leq \frac{\|e\|}{\|x\|} \\ \|x\| &\leq \|A^{-1}\| \|b\| \end{aligned} \right\} \Rightarrow \frac{1}{\|A\| \|A^{-1}\|} \frac{\|r\|}{\|b\|} \leq \frac{\|e\|}{\|x\|}$$

by the same argument

Thus,

$$\frac{1}{\|A\| \|A^{-1}\|} \frac{\|r\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|r\|}{\|b\|}$$

condition number $\triangleq \|A\| \cdot \|A^{-1}\| = \kappa$

so that

$$\frac{1}{\kappa} \frac{\|r\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \kappa \frac{\|r\|}{\|b\|}$$

Thus, this gives a relative "error bar" on the solution in terms of the condition number κ and the residual r .

$\|r\| \approx$ machine precision ϵ
 $\|b\| \approx O(1)$

error e can be ~~arbitrarily~~ ^{arbitrarily} large if the system matrix A has a large condition number.

ITERATIVE METHODS

Iterative methods start with an initial guess $x^{(0)}$ which is then ~~used~~ repeatedly refined. They are a viable alternative for sparse problems.

Jacobi Method

Split the matrix as follows

$$\square = \begin{array}{|c|} \hline \square \\ \hline \end{array} = \begin{array}{|c|} \hline \square \\ \hline \end{array} + \begin{array}{|c|} \hline \square \\ \hline \end{array} + \begin{array}{|c|} \hline \square \\ \hline \end{array}$$

$$A = L + D + U$$

so that $Ax = (L + D + U)x = b$

Rewrite as

$$Dx^{(n+1)} = -(L+U)x^{(n)} - b$$

$$x^{(n+1)} = -D^{-1}[(L+U)x^{(n)} - b], \quad n=1, \dots, \quad x^{(0)} \text{ - given initial guess}$$

$$x_i^{(n+1)} = -\sum_{\substack{j=1 \\ j \neq i}}^N \frac{a_{ij}}{a_{ii}} x_j^{(n)} + \frac{b_i}{a_{ii}}, \quad i=1, \dots, N \\ n=1, \dots$$

This is in fact a fixed-point iteration in \mathbb{R}^N with

$$x^{(n+1)} = G(x^{(n)})$$

Contractivity condition
on G necessary
for
convergence

\Rightarrow diagonal dominance
of matrix A

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^N |a_{ij}| \quad i=1, \dots, N$$

Gauss-Seidel Method

The Jacobi method can be accelerated by using at a given iteration the elements which have already been updated

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

In the matrix notation

$$(L+D)x^{(n+1)} = -Ux^{(n)} + b$$

$$x^{(n+1)} = -(L+D)^{-1}Ux^{(n)} + (L+D)^{-1}b$$

This is also a fixed-point approach. The necessary conditions for convergence are similar to the Jacobi method, but the Gauss-Seidel method is faster.
* Computational cost of ~~many~~ iterative methods: $O(mn^2)$
m - # of iterations

MULTIDIMENSIONAL ROOT-FINDING

(Section 8.2 of Grasselli & Pelinovsky)

Essentially, a generalization of the 1D problem but with ~~important~~ important technical modifications

Problem

Given a continuous function $\bar{F}: \mathbb{R}^n \rightarrow \mathbb{R}^n$, find all the solutions of $\bar{F}(x) = 0$ for x in the region $[a_1, b_1] \times \dots \times [a_n, b_n]$.