Lecture 8 Conditioning, Banded Systems

David Semeraro

University of Illinois at Urbana-Champaign

September 19, 2013

- How do we know if GE if will be accurate?
 - norms, condition number, theory
- Can we reduce cost for special systems
 - tridiagonals, banded, etc

Geometric Interpretation of Singularity

Consider a 2×2 system describing two lines that intersect

$$y = -2x + 6$$
$$y = \frac{1}{2}x + 1$$

The matrix form of this equation is

$$\begin{bmatrix} 2 & 1 \\ -1/2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 1 \end{bmatrix}$$

The equations for two parallel but not intersecting lines are

$$\begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 5 \end{bmatrix}$$

Here the coefficient matrix is singular (rank(A) = 1), and the system is inconsistent

The equations for two parallel and coincident lines are

$$\begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 6 \end{bmatrix}$$

The equations for two nearly parallel lines are

$$\begin{bmatrix} 2 & 1 \\ 2+\delta & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 6+\delta \end{bmatrix}$$

Geometric Interpretation of Singularity



Consider the solution of a 2×2 system where

$$b = \begin{bmatrix} 1\\ 2/3 \end{bmatrix}$$

One expects that the exact solutions to

$$Ax = \begin{bmatrix} 1\\ 2/3 \end{bmatrix}$$
 and $Ax = \begin{bmatrix} 1\\ 0.6667 \end{bmatrix}$

will be different. Should these solutions be a lot different or a little different?

Norms

Vectors:

$$||x||_{p} = (|x_{1}|^{p} + |x_{2}|^{p} + \dots + |x_{n}|^{p})^{1/p}$$
$$||x||_{1} = |x_{1}| + |x_{2}| + \dots + |x_{n}| = \sum_{i=1}^{n} |x_{i}|$$
$$||x||_{\infty} = \max(|x_{1}|, |x_{2}|, \dots, |x_{n}|) = \max_{i}(|x_{i}|)$$

Matrices:

$$||A|| = \max_{x \neq 0} \frac{||Ax||}{||x||}$$
$$||A||_{p} = \max_{x \neq 0} \frac{||Ax||_{p}}{||x||_{p}}$$
$$||A||_{1} = \max_{1 \leq i \leq n} \sum_{i=1}^{m} |a_{ij}|$$
$$||A||_{\infty} = \max_{1 \leq i \leq m} \sum_{i=1}^{n} |a_{ij}|$$

< ≣⇒

I

 $\|\alpha x\| = |\alpha| \|x\|$ $\|Ax\| \le \|A\| \|x\|$ $\|x + y\| \le \|x\| + \|y\|$

Challenge: Make sure that you can prove these properties.

Effect of Perturbations to b

Perturb b with δb such that

$$\frac{\|\delta b\|}{\|b\|} \ll 1,$$

The perturbed system is

$$A(x+\delta x_b)=b+\delta b$$

The perturbations satisfy

 $A\delta x_b = \delta b$

Analysis shows (see next two slides for proof) that

$$\frac{\|\delta x_b\|}{\|x\|} \leqslant \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}$$

Thus, the effect of the perturbation is small only if $||A|| ||A^{-1}||$ is small.

$$\frac{\|\delta x_b\|}{\|x\|} \ll 1 \quad \text{only if} \quad \|A\| \|A^{-1}\| \sim 1$$

Let $x + \delta x_b$ be the *exact* solution to the perturbed system

$$A(x + \delta x_b) = b + \delta b \tag{1}$$

Expand

 $Ax + A\delta x_b = b + \delta b$

Subtract *Ax* from left side and *b* from right side since Ax = b

 $A\delta x_b = \delta b$

Left multiply by A^{-1}

$$\delta x_b = A^{-1} \delta b \tag{2}$$

Effect of Perturbations to *b* (Proof, p. 2)

Take norm of equation (2)

$$\left\|\delta x_{b}\right\| = \left\|A^{-1}\,\delta b\right\|$$

Applying consistency requirement of matrix norms

$$\|\delta x_b\| \leqslant \|A^{-1}\| \|\delta b\| \tag{3}$$

Similarly, Ax = b gives ||b|| = ||Ax||, and

$$\|b\|\leqslant \|A\|\|x\|$$

Rearrangement of equation (4) yields

$$\frac{1}{\|x\|} \leqslant \frac{\|A\|}{\|b\|}$$

(4)

(5)

Multiply Equation (5) by Equation (3) to get

$$\frac{\|\delta x_b\|}{\|x\|} \le \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}$$
(6)

• □ ▶ • □ ▶ • □ ▶

Summary:

If $x + \delta x_b$ is the *exact* solution to the perturbed system

$$A(x+\delta x_b)=b+\delta b$$

then

$$\frac{\|\delta x_b\|}{\|x\|} \leqslant \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}$$

Effect of Perturbations to A

Perturb A with δA such that

$$\frac{\|\delta A\|}{\|A\|} \ll 1,$$

The perturbed system is

$$(A + \delta A)(x + \delta x_A) = b$$

Analysis shows that

$$\frac{\|\delta x_A\|}{\|x+\delta x_A\|} \leqslant \|A\| \|A^{-1}\| \frac{\|\delta A\|}{\|A\|}$$

Thus, the effect of the perturbation is small *only if* $||A|| ||A^{-1}||$ is small.

$$\frac{\|\delta x_A\|}{\|x+\delta x_A\|} \ll 1 \quad \text{only if} \quad \|A\| \|A^{-1}\| \sim 1$$

Effect of Perturbations to both *A* and *b*

Perturb both A with δA and b with δb such that

$$\frac{\|\delta A\|}{\|A\|} \ll 1 \quad \text{and} \quad \frac{\|\delta b\|}{\|b\|} \ll 1$$

The perturbation satisfies

$$(A + \delta A)(x + \delta x) = b + \delta b$$

Analysis shows that

$$\frac{\|\delta x\|}{\|x+\delta x\|} \ \leqslant \ \frac{\|A\|\|A^{-1}\|}{1-\|A\|\|A^{-1}\|\frac{\|\delta A\|}{\|A\|}} \left[\frac{\|\delta A\|}{\|A\|}+\frac{\|\delta b\|}{\|b\|}\right]$$

Thus, the effect of the perturbation is small *only if* $||A|| ||A^{-1}||$ is small.

$$\frac{\|\delta x\|}{\|x+\delta x\|} \ll 1 \quad \text{only if} \quad \|A\| \|A^{-1}\| \sim 1$$

The condition number

 $\kappa(A) \equiv \|A\| \|A^{-1}\|$

indicates the sensitivity of the solution to perturbations in A and b. The condition number can be measured with any p-norm. The condition number is always in the range

 $1\leqslant \kappa(A)\leqslant \infty$

- $\kappa(A)$ is a mathematical property of A
- Any algorithm will produce a solution that is sensitive to perturbations in A and b if κ(A) is large.
- In exact math a matrix is either singular or non-singular. $\kappa(A) = \infty$ for a singular matrix
- $\kappa(A)$ indicates how close A is to being numerically singular.
- A matrix with large κ is said to be ill-conditioned

In Practice, applying Gaussian elimination with partial pivoting and back substitution to Ax = b gives the **exact solution**, \hat{x} , to the **nearby problem**

 $(A+E)\hat{x} = b$ where $||E||_{\infty} \leq \varepsilon_m ||A||_{\infty}$

Gaussian elimination with partial pivoting and back substitution "gives exactly the right answer to nearly the right question." — Trefethen and Bau An algorithm that gives the exact answer to a problem that is near to the original problem is said to be **backward stable**. Algorithms that are not backward stable will tend to amplify roundoff errors present in the original data. As a result, the solution produced by an algorithm that is not backward stable will not necessarily be the solution to a problem that is close to the original problem.

Gaussian elimination without partial pivoting is *not* backward stable for arbitrary *A*.

If *A* is symmetric and positive definite, then Gaussian elimination without pivoting is backward stable.

Let \hat{x} be the numerical solution to Ax = b. $\hat{x} \neq x$ (x is the exact solution) because of roundoff.

The **residual** measures how close \hat{x} is to satisfying the original equation

$$r = b - A\hat{x}$$

It is not hard to show that

$$\frac{\hat{x} - x\|}{\|x\|} \leqslant \kappa(A) \frac{\|r\|}{\|b\|}$$

Small ||r|| does not guarantee a small $||\hat{x} - x||$. If $\kappa(A)$ is large the \hat{x} returned by Gaussian elimination and back substitution (or any other solution method) is not guaranteed to be anywhere near the true solution to Ax = b.

- Applying Gaussian elimination with partial pivoting and back substitution to Ax = b yields a numerical solution \hat{x} such that the residual vector $r = b A\hat{x}$ is small *even if* the $\kappa(A)$ is large.
- If A and b are stored to machine precision ε_m, the numerical solution to Ax = b by any (good) variant of Gaussian elimination is correct to d decimal digits where

$$d = |\log_{10}(\varepsilon_m)| - \log_{10}(\kappa(A))$$

$$d = |\log_{10}(\varepsilon_m)| - \log_{10}(\kappa(A))$$

Example:

MATLAB computations have $\varepsilon_m \approx 2.2 \times 10^{-16}$. For a system with $\kappa(A) \sim 10^{10}$ the elements of the solution vector will have

$$\begin{aligned} d &= |\log_{10}(2.2 \times 10^{-16})| - \log_{10} (10^{10}) \\ &\approx 16 - 10 \\ &= 6 \end{aligned}$$

correct digits

Summary of Limits to Numerical Solution of Ax = b

- $\kappa(A)$ indicates how close *A* is to being numerically singular
- **2** If $\kappa(A)$ is "large", *A* is **ill-conditioned** and *even the best* numerical algorithms will produce a solution, \hat{x} that cannot be guaranteed to be close to the true solution, *x*
- In practice, Gaussian elimination with partial pivoting and back substitution produces a solution with a small residual

$$r = b - A\hat{x}$$

even if $\kappa(A)$ is large.

- tridiagonal systems
- banded systems
- LU decomposition
- Cholesky factorization

3



A tridiagonal matrix A

- storage is saved by not saving zeros
- only n + 2(n 1) = 3n 2 places are needed to store the matrix versus n^2 for the whole system
- can operations be saved? yes!

-

_

Tridiagonal

Start forward elimination (without any special pivoting)

- subtract a_1/d_1 times row 1 from row 2
- 2 this eliminates a_1 , changes d_2 and does not touch c_2
- continuing:

$$d_i = d_i - \left(\frac{a_{i-1}}{d_{i-1}}c_{i-1}\right)$$

for $i = 2 \dots n$

$$\begin{bmatrix} \tilde{d}_1 & c_1 & & & \\ & \tilde{d}_2 & c_2 & & & \\ & & \tilde{d}_3 & c_3 & & & \\ & & & \ddots & \ddots & & \\ & & & & \tilde{d}_i & c_i & & \\ & & & & \ddots & \ddots & \\ & & & & & & \tilde{d}_n \end{bmatrix}$$

This leaves an upper triangular (2-band). With back substitution:

1
$$x_n = \tilde{b}_n / \tilde{d}_n$$
2 $x_{n-1} = (1 / \tilde{d}_{n-1}) (\tilde{b}_{n-1} - c_{n-1} x_n)$
3 $x_i = (1 / \tilde{d}_i) (\tilde{b}_i - c_i x_{i+1})$

E

Tridiagonal Algorithm

```
input: n, a, d, c, b
1
      for i=2 to n
2
         xmult = a_{i-1}/d_{i-1}
3
         d_i = d_i - xmult \cdot c_{i-1}
4
         b_i = b_i - xmult \cdot b_{i-1}
5
      end
6
      x_n = b_n/d_n
7
      for i = n - 1 down to 1
8
         x_i = (b_i - c_i x_{i+1})/d_i
9
      end
10
```

Challenge: Will this algorithm make good use of the processor cores in a multicore processor? Why or why not?



Image: Image:

m-band



- the *m* correspond to the total width of the non-zeros
- after a few passes of GE fill-in with occur within the band
- so an empty band costs (about) the same an a non-empty band
- one fix: reordering (e.g. Cuthill-McKee)
- generally GE will cost $O(m^2n)$ for *m*-band systems