

## Numerical singularity (vs. exact singularity)

23 Oct '23

During numerical factorization in a floating-point system,

we usually cannot expect exact zeros.

But a matrix is "near singular" if all elements below the diagonal in column  $k$  (and  $\hat{a}_{kk}$ )

at the  $k^{\text{th}}$  stage have magnitude  $\lesssim \epsilon_{\text{ps}} \max_{i,j} |a_{ij}|$ .

$i \geq j < k$

We say the matrix has numerical singularity.

↑ commonly used heuristic

## Complexity of Gaussian Elimination

We will count multiplication/addition pairs; i.e.,  $mx + b$ .

Each is a flop (floating-point operation).  $\uparrow$  mult  $\uparrow$  add

Comparisons and divisions are also there, but they are relatively cheap ( $O(n)$  per stage);

so we just count flops.

### Computing the LU-factorization

1<sup>st</sup> stage

$$\begin{bmatrix} \boxed{\phantom{0}} & & & & \\ 0 & * & * & \dots & * \\ 0 & * & * & \dots & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & * & * & \dots & * \end{bmatrix} \quad \begin{matrix} n-1 \\ *s \text{ are elements inducing flops} \end{matrix}$$

Adding a multiple of row ① to rows ②, ..., ⑤ costs  $(n-1)^2$  flops.

2<sup>nd</sup> stage  $(n-2)^2$  flops

$\vdots$

$(n-1)$  stage 1 flop

$$\begin{aligned} \text{Total: } (n-1)^2 + (n-2)^2 + \dots + 1 &= \sum_{i=1}^{n-1} (n-i)^2 \\ &= \frac{(n-1)n(2(n-1)+1)}{6} \\ &= \frac{(n-1)n(2n-1)}{6} \\ &= \frac{n^3}{3} + O(n^2) \text{ flops.} \end{aligned}$$

$\uparrow$  this constant matters in practice

### Computing forward solve ( $Ld = b$ )

Structure:

$$\begin{bmatrix} 1 & & & 0 \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \begin{bmatrix} d_1 \\ \vdots \\ d_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \Rightarrow \begin{cases} d_1 = b_1 & 0 \text{ flops} \\ d_2 = b_2 - l_{21}d_1 & 1 \text{ flop} \\ d_3 = \underbrace{b_3 - l_{32}d_2 - l_{31}d_1}_{\text{2 flops}} & 2 \text{ flops} \\ \vdots & \\ d_n = b_n - l_{n(n-1)}d_{n-1} - \dots - l_{n2}d_2 & n-1 \text{ flops} \end{cases}$$

$$\begin{aligned} \text{Total: } 0 + 1 + 2 + \dots + (n-1) &= \sum_{i=0}^{n-1} i = \frac{(n-1)n}{2} \\ &= \frac{n^2}{2} + O(n) \text{ flops.} \end{aligned}$$

### Computing backward solve ( $Ux = d$ )

Similar:  $\frac{n^2}{2} + O(n)$  flops.

Total forward/backward:  $n^2 + O(n)$ .

Compare to LU-factorization:  $\frac{n^3}{3} + O(n^2)$ .  $\left. \begin{matrix} \uparrow \\ \frac{n}{3} \text{ times more expensive (in flops)!} \end{matrix} \right\}$

This is why we factorize once if there are several systems to solve, only differing in  $b$ .

If we applied Gaussian elim for each, it would take  $O(n^3)$  each time.

## Roundoff error analysis of GE

Factorization:  $PA = LU$

Due to rounding error (initial and propagated), actually get  $\hat{L}$ ,  $\hat{U}$  and  $\hat{P}$  st

$$\hat{P}(A+E) = \hat{L}\hat{U}.$$

Hopefully,  $\|E\|$  is small compared to  $\|A\|$  (if we pivot during the factorization).

( $\hat{P}$  is potentially a diff't pivoting strategy from  $P$ , but it is still just  $I_n$  with rows permuted.)

Actually, solving (forward and backward) can introduce more roundoff error,

but it can all still be represented as

$$(A+\tilde{E})\hat{x} = b$$

where  $\tilde{E}$  is slightly diff't from  $E$  (but is essentially  $E$ ).

Dropping the distinction:  $(A+E)\hat{x} = b$ .

Equivalently, let  $E\hat{x} = r$ ; then

$$(A+E)\hat{x} = b \Leftrightarrow r = b - A\hat{x}.$$

↑ residual  
(what's left over)

Is  $\|E\|$  small compared to  $\|A\|$ ?

Result: If we use row partial pivoting during the factorization,  
can show that

$$(*) \quad \|E\| \lesssim k \cdot \text{eps} \|A\|,$$

where  $k$  is not too large, grows with  $n$ , depends on pivoting.

Similarly, for computed sol<sup>n</sup>  $\hat{x}$  and  $r = b - A\hat{x}$ ,

$$\|r\| \lesssim k \cdot \text{eps} \|b\|$$

$$(*) \quad \Leftrightarrow \frac{\|r\|}{\|b\|} \lesssim k \cdot \text{eps}.$$

↑ relative residual

Can prove this, but very technical.

Does this mean that  $\frac{\|\hat{x} - x\|}{\|x\|}$  (relative error) is small?  
↑ true sol<sup>n</sup>

Remember: We don't have  $x$  (true sol<sup>n</sup>), so we can't directly calculate relative error.

e.g.:

$$\begin{bmatrix} .700 & .563 \\ .913 & .659 \end{bmatrix} x = \begin{bmatrix} .227 \\ .254 \end{bmatrix}. \quad (\text{True sol}^n: x = \begin{bmatrix} 1 \\ -1 \end{bmatrix}; \text{ in general not known})$$

Consider the computed sol<sup>n</sup>s

$$\hat{x}_\alpha = \begin{bmatrix} .999 \\ -1.001 \end{bmatrix}, \quad \hat{x}_\beta = \begin{bmatrix} .341 \\ -.087 \end{bmatrix}.$$

Residuals:

$$\begin{aligned} r_\alpha &= b - A\hat{x}_\alpha & r_\beta &= b - A\hat{x}_\beta \\ &= \begin{bmatrix} -.001243 \\ -.001572 \end{bmatrix}, & &= \begin{bmatrix} -.000001 \\ 0 \end{bmatrix}. \end{aligned}$$

↑ smaller residual

$$\text{So } \frac{\|r_\beta\|}{\|b\|} \ll \frac{\|r_\alpha\|}{\|b\|}.$$

Then why is  $\frac{\|\hat{x}_\alpha - x\|}{\|x\|}$  so much smaller than  $\frac{\|\hat{x}_\beta - x\|}{\|x\|}$ ??

Need the relationship b/w relative error and relative residual.

In particular, when does a small relative residual (which is guaranteed if we use row

partial pivoting) guarantee a small relative error — how do we check sol<sup>n</sup> accuracy?

$$\frac{\|r\|}{\|b\|} \xleftrightarrow{\text{relationship?}} \frac{\|\hat{x} - x\|}{\|x\|}$$

relative residual                      relative error