

ex) find  $y = \cos x = f(x)$ .

For this case,  $\delta y \approx -\sin x \delta x$ . Hence

$$\left| \frac{\delta y}{y} \right| \approx k \left| \frac{\delta x}{x} \right| \quad k = \frac{x f'(x)}{f} = -x \tan x$$

$$\therefore \left| \frac{\delta y}{y} \right| = |x| |\tan x| \left| \frac{\delta x}{x} \right| = k \left| \frac{\delta x}{x} \right|$$

take  $x \approx \frac{\pi}{2}$  for example. We see  
that  $f(x)$  is highly sensitive to perturbations in  $x$ . //

## CONDITION NUMBER OF MATRIX-VECTOR MULTIPLICATION

Take:  $x \xrightarrow{Ax} y$  i.e.  $y = Ax$

$$A \in \mathbb{C}^{m \times n}$$

$$k_2 = \sup_{\delta x} \left( \frac{\| \delta(Ax + \delta x) - Ax \| \| Ax \|}{\| Ax \| \| \delta x \|} \right)$$

$$= \sup_{\delta x} \frac{\| A \delta x \|}{\| \delta x \|} \frac{\| Ax \|}{\| Ax \|}$$

$$\therefore k = \|A\| \frac{\|x\|}{\|\Delta x\|}$$

If  $A$  is square & nonsingular

$$\frac{\|x\|}{\|\Delta x\|} < \|A^{-1}\|$$

$$\therefore k = \alpha \|A\| \|A^{-1}\|$$

$$\alpha = \frac{\|x\|}{\|\Delta x\|} / \|A^{-1}\| \approx 1$$

If  $\Delta \in \mathbb{C}^{m \times n}$  use

$$k = \alpha \|A\| \|A^+\| \quad //$$

Thm: Let  $A \in \mathbb{C}^{m \times n}$  be nonsingular &  $Ax = b$ .

The problem of computing  $b$ , given  $x$  has  
condition number

$$k = \|A\| \frac{\|x\|}{\|b\|} \leq \|A\| \|A^+\|$$

with respect to perturbations in  $x$ .

The problem of finding  $x$ , given  $b$ :

$$k = \|A\| \frac{\|b\|}{\|x\|} \leq \|A\| \|A^{-1}\| \quad //$$

### CONDITION NUMBER OF A MATRIX

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

if  $k$  is small  $\approx$  "well-conditioned"

if  $k$  is very large "ill-conditioned"

( $k$  can be  $\infty$ )

Spec  $\|\cdot\|_2$  then

$$\|A\|_2 = G_1$$

$$\|A^{-1}\| = \frac{1}{G_m}$$

$$\therefore k(A) = \frac{G_1}{G_m}$$

it "measures" the ellipticity of the hyperellipse.

What happens when  $A \rightarrow A + \delta A$ ?

How do we estimate the sensitivity

of outcomes to "errors" in  $A$ ?

let's fix  $b$ .

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

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

$$\delta Ax + A\delta x + \underbrace{\delta A\delta x}_\text{too small} = 0$$

$$\therefore \delta Ax = -A\delta x$$

$$\delta x = -A^{-1}(\delta A)x$$

Bound  $\delta x$ :

$$\|\delta x\| \leq \|A^{-1}\| \|\delta A\| \|x\|$$

Want  $\frac{\delta x}{\delta A}$  :

$$\frac{\|\delta x\|}{\|x\|} \leq \|A^{-1}\| \|\delta A\| = k(A)$$

$\therefore k(A)$  captures sensitivity to  
errors in  $A$  //

# FLOATING POINT NUMBERS

## Lecture 13 Finite-Precision Computing

Finite Precision Machine (double precision binary machine)

max positive number  $\sim 1.79 \cdot 10^{308}$

min positive number  $\sim 2.23 \cdot 10^{-308}$

resolution (largest gap) =  $2^{-52} \sim 2.22 \cdot 10^{-16}$

representation (binary, octal, hex, decimal, etc.)

The IEEE standard is binary (memory representation)

Floating Point Representation: means position of decimal (or binary point) "moves" (stored separately). The gaps between the numbers will scale in proportion to size of the numbers.

$fl(x)$  is the "floating point representable" of the number  $x$

Representation: consists of 4 fields:

//

$\beta$ base	{
$t$ precision	
$e \in [L, U]$ exponent	
Sign indicator	

$$x = \pm \left( \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots + \frac{d_t}{\beta^t} \right) \beta^e$$

$d_1, d_2, \dots, d_t$        $\beta^e$

Normalized notation

The normalized representation in decimal is 0..... The normalized representation in binary is 1....

$$0 \leq d_i \leq \beta - 1 \quad i=0,1,\dots,t$$

$$e \in [l, u]$$

ex) Normalized representation

$$100 \text{ in } \beta = (0 \ 0.1 \cdot 10^3)$$

Hexadecimal

$$\beta = 16 \quad 0 \leq d_i \leq 15$$

$$d_i = 0, 1, 2, \dots, 9, A, B, C, D, E, F$$

$$(23A.D)_{16} = 2 \cdot 16^2 + 3 \cdot 16^1 + 10 \cdot 16^0 + 13 \cdot 16^{-1}$$

$$= 512 + 48 + 10 + 0.8125 = 570.8125$$

Binary  $\beta = 2 \quad 0 \leq d_i \leq 1$

$$\text{ex) } (1101.0011)_2 = 10.11010011 \cdot 2^4$$

$$= 1 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 + 0 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} + 1 \cdot 2^{-4}$$

$$= 13.1875$$

ex) Convert  $\frac{2}{3}$  to binary:

$$\frac{2}{3} = (0.a_1a_2\dots)_2$$

x2

$$\frac{4}{3} = (a_1.a_2\dots)_2 \rightarrow a_1=1$$

subtract 1 from b.s.

$$\frac{4}{3}-1=\frac{1}{3}=(0.a_2\dots)_2$$

x2

$$\frac{2}{3}=(a_2.a_3\dots)_2 \quad a_2=0$$

x2

$$\frac{4}{3}=(a_3.a_4\dots)_2 \quad a_3=1$$

etc

$$\frac{2}{3}=(0.1010\dots)_2$$

## IEEE 754-1985 Standard

	$\beta$	$t$	$L$	$U$	wordlength
HP (Half)	2	11	-14	15	16
SP (single)	2	24	-126	127	32
DP (double)	2	53	-1022	1023	64
QP (quad)	2	113	-16383	16384	128

Example SP 32 bit



UFL (underflow)  $\beta^{-t}$

OFL (overflow)  $\beta^{t+1} (1 - \beta^{-t})$

In binary  $x = (-1)^s q \cdot 2^m$

$s \geq 0$  positive     $s=1$  negative

$q = (1.f)$   
↳ implied (Normalized binary notation)

$m = e - b$

$m = e - 2^7$  (for 32-bit)  
8-bit biased

Range  $0 < e < 1111111 = 2^8 - 1 = 255$

The number 0 is reserved for ± 0  
" .. 255 " " " " for NaN

$m = e - 127 \therefore m \in [-126, 127]$

The smallest positive #  $2^{-126} \approx 1.2 \cdot 10^{-38}$

The largest positive #  $(2 - 2^{-23}) \times 2^{127} \approx 3.4 \cdot 10^{38}$

The least significant bit  $2^{-23}$

$$\epsilon_{PS} = \frac{1}{2} \beta^{1-t} \quad (\text{rounding})$$

Machine epsilon

$$\epsilon_{PS_{SP}} = \frac{1}{2} 2^{-23} : 2^{-24} \approx 6 \cdot 10^{-8}$$

$$\epsilon_{PS_{DP}} = \frac{1}{2} 2^{-53} \approx 1 \cdot 10^{-16}$$

The  $\frac{1}{2}$  results from using rounding approximation scheme.