

LECTURE

LAST TIME

□ SUMMARY OF GRAM SCHMIDT ALGORITHM.

□ STABILITY OF A FUNCTION $f: \mathbb{C}^n \rightarrow \mathbb{C}^m$

- ERRORS $\begin{cases} \text{NOISE} \\ \text{TRUNCATION} \\ \text{ROUNDING} \end{cases}$

- CONDITION NUMBER

$$k(x) = \lim_{\delta x \rightarrow 0} \sup_{\delta x} \frac{\|\delta f\| / \|f\|}{\|\delta x\| / \|x\|}$$

$$\delta f = f(x + \delta x) - f(x)$$
$$f: \mathbb{C}^n \rightarrow \mathbb{C}^m$$

\Rightarrow

RELATIVE OUTPUT ERROR (at point x)

$$= k(x) \text{ RELATIVE INPUT ERROR}$$

② IF f IS DIFFERENTIABLE $k(x) = \frac{\|J(x)\| \|x\|}{\|f(x)\|}$

③ IF f IS LINEAR ($f = Ax$) $k(x) = \frac{\|A\| \|x\|}{\|Ax\|} \leq \|A\| \|A\|^{-1}$

④ CONDITION NUMBER OF A SQUARE INVERTIBLE MATRIX:

$$k(A) = \|A\| \|A\|^{-1}. \quad \text{FOR THE 2-NORM } k_2(A) = \frac{\sigma_1}{\sigma_m}$$

$$k(A) \geq 1. \quad (\text{SHOULD BE AN EASY EXERCISE}) \quad \text{FOR ANY NORM}$$

where $\sigma_{\max} = \sigma_1$; $\sigma_{\min} = \sigma_m$

THE LARGEST AND SMALLEST SINGULAR VECTORS.

□ COROLLARY : $K_2(A) \geq 1$.

TODAY

□ FLOATING POINT ARITHMETIC.

□ STABILITY TO ROUNDING [Chapters 13, 14, 15].

□ BACKWARD ERROR

A NOTE OF THE SHARPNESS OF THE CONDITION NUMBER

• CONSIDER $A = \begin{bmatrix} 1 & 0 \\ 0 & \epsilon \end{bmatrix}$, $0 < \epsilon < 1$.

THEN BY DEFINITION, $K = \frac{1}{\epsilon}$.

Let $x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$. Let $\delta x = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$

Then $\delta f = \begin{bmatrix} 3 \\ \epsilon \end{bmatrix} - \begin{bmatrix} 0 \\ \epsilon \end{bmatrix} = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$; $f = \begin{bmatrix} 0 \\ \epsilon \end{bmatrix}$

$\|\delta f\|_2 = 3$; $\|f\|_2 = \epsilon$; $\|x\| = 1$; $\delta x = 3$.

$\frac{\|\delta f\|}{\|f\|} = \frac{3}{\epsilon}$; $\frac{\|\delta x\|}{\|x\|} = 3$; $\frac{\|\delta f\|/\|f\|}{\|\delta x\|/\|x\|} = \frac{1}{\epsilon}$.

FLOATING POINT NUMBERS

- Normalized Decimal REPRESENTATION

$$x = \pm S \times 10^E$$

$$S = J.d_1 \dots d_{p-1}$$

$$d_i: 0 \dots 9$$

$$J: 1, \dots, 9 \quad (J \neq 0)$$

e.g. $43.501 = 4.3501 \times 10^1$
 $-0.0134 = -1.3400 \times 10^{-2}$

- NORMALIZED BINARY REPRESENTATION

$$x = \pm m' \times 2^{\textcircled{E}} \rightarrow \text{Exponent}$$

$$\textcircled{m'} = 1.b_1 b_2 \dots b_{p-1}$$

$$b_i = \{0, 1\}$$

↑ mantissa or fraction or significant

p: precision.

Example: $x = -1.01 \times 2^{0101}$
 $= - (1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2}) \times 2^{(0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0)}$

$$(\text{in base 10}) \quad = -\left(1 + \frac{1}{2^2}\right) 2^5 = -40$$

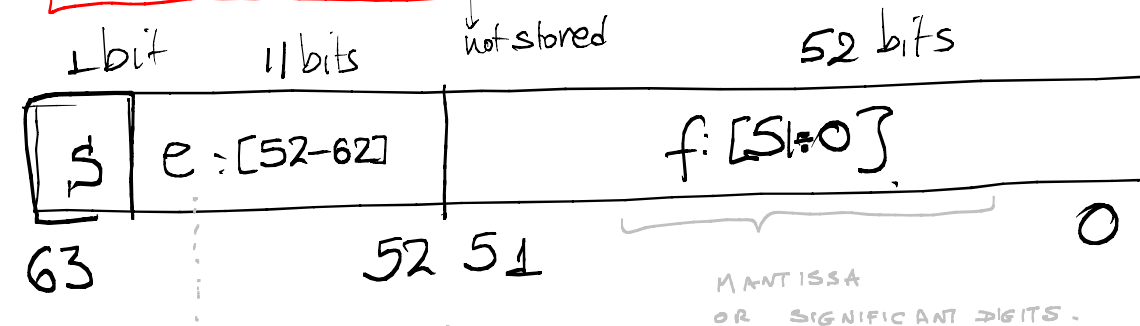
IEEE-754 FLOATING POINT STANDARD.

- 64-bits to represent a number.

- WE CAN EXACTLY REPRESENT 2^{64} NUMBERS ONLY
($\sim 1.84 \times 10^{19}$)

- $$x = (-1)^s \times \underbrace{1.f}_{\text{not stored}} \times 2^{e-1023}$$

FORMAT
OF AN
IEEE-754 NUMBER



NOTE THAT "e" IS NOT THE EXPONENT. THE EXPONENT "E" IS $E = e - 1023$. "e" IS POSITIVE BUT E CAN BE POSITIVE & NEGATIVE

$$0 < e < 2047$$

$$\downarrow$$

$$-102 \leq E \leq 1023$$

Notice that $2^{11} - 1 = 2047$
So not all bits of e-
ARE ALLOWED TO BE ONES.

($\sum_{i=0}^{10} 2^i = 2047$)
THE "e" = 1...1 IS USED TO CODE SPECIAL CASES. SEE TABLE BELOW.

BIT PATTERN	VALUE
$0 < e < 2047$	$(-1)^s \times 2^{e-1023} \times 1.f$
$e=0; f \neq 0$	$(-1)^s \times 2^{-1022} \times 0.f$ "Subnormal numbers"
$e=0; f=0$	$(-1)^s \times 0.0$ (signed zero)
$s=0; e=2047; f=0$	INF
$s=-1; e=2047; f=0$	-INF
$e=2047; f \neq 0$	NAN : not a number

powers of two : represented exactly (if in range)

$$\pm (1.000 \dots) 2^{e-1023}$$

max number : $\sim 1.8e+308$

min number > 0 : $\sim 2.2e-308$

min subnormal > 0 : $\sim 5e-324$

(larger numbers \rightarrow overflow)
(smallest numbers \rightarrow underflow)

SPACING OF FLOATING POINT NUMBERS

• Choosing $f=0$ and $e=1023$ we get

$$1 \rightarrow + (1.0000 \dots 0) \cdot 2^{1023-1023} = + (1.000) \times 2^0 = 1$$

The next largest number is

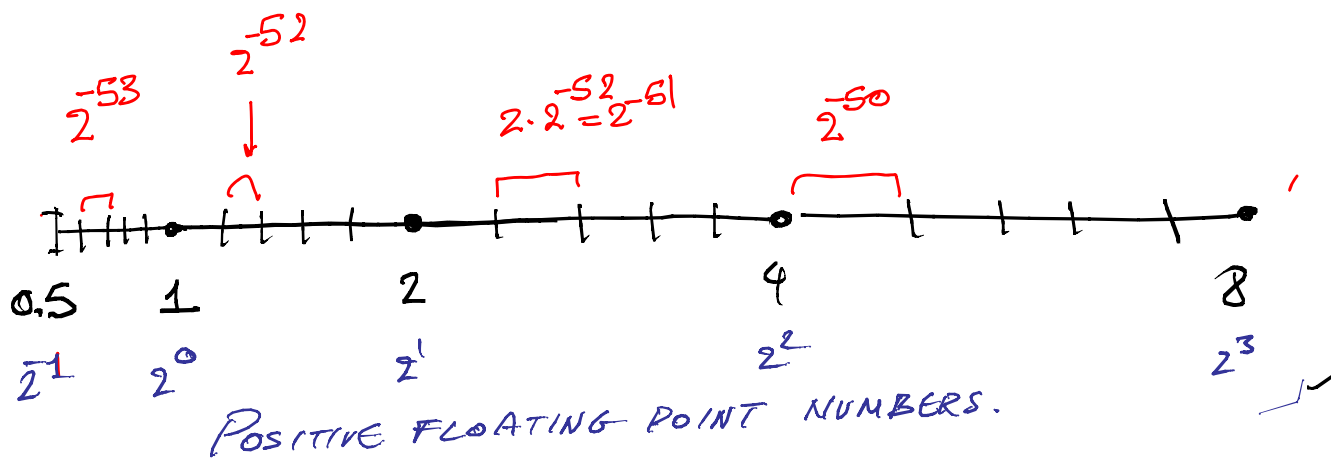
$$\begin{aligned} &+ (1.000 \dots 001) 2^0 = 1 + 2^{-52} \\ &\quad \downarrow \\ &\quad 1 \times 2^{-52} \\ &= 1 + 2.20446049e-16 \end{aligned}$$

$$2 \rightarrow (1.000 \dots) \times 2^{1024-1023}$$

The next largest number is

$$(1.00 \dots 001) 2^1 = 2 + \underline{2} \cdot 2^{-52}$$

Therefore the gap between floating point numbers is non-constant



SIDE-NOTE (HISTORICAL REFERENCE).

WM - KAHAN (UC BERKELEY) MAIN

DESIGNER OF IEEE-754.

RECEIVED TURING AWARD IN 1988.

DEFINITIONS.

Precision: number of bits for the mantissa

DOUBLE PRECISION: 53 bits \approx 16 decimal digits.

SINGLE PRECISION: 24 bits.

Accuracy: THE NUMBER OF CORRECT DIGITS.

Machine epsilon: ϵ_M : $1 + \epsilon_M$ is the next floating point number after 1.

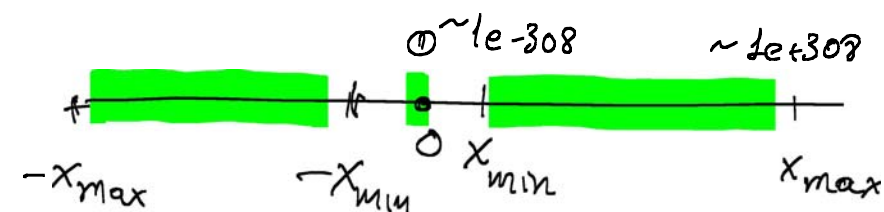
In IEEE-754 $\epsilon_M = 2^{-52} \sim 2.2e-16$

NOTATION

$F_0 :=$ SET OF ALL NORMAL FLOATING POINT NUMBERS

$F :=$ SET OF ALL FLOATING POINT NUMBERS (NORMAL + SUBNORMAL + $\pm INF$)

• Range (F_0):

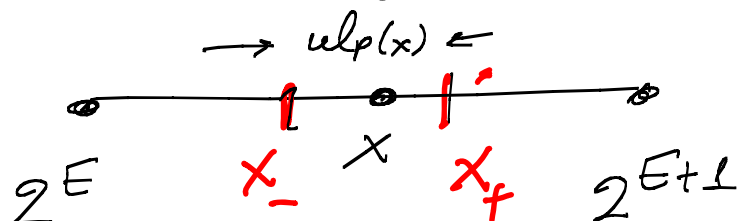


• $ulp(x)$: gap in the power-2 interval that $|x|$ belongs to.

$$\Rightarrow ulp(x) = (0.000...001) \times 2^E \Rightarrow$$

$$ulp(x) = \epsilon_M 2^E$$

• FOR ANY $x \in \mathbb{R}$ that belongs in the range of normal numbers,



$$x_+ = x_- + ulp(x).$$

x_- : NEAREST $y \in F_0 : y < x$, x_+ : NEAREST $y \in F_0 : y > x$

\tilde{x} : rounding of $x \in \text{Range}(F_0)$

Depending on the hardware settings

↓
we will also use $[x]$ to indicate rounding.

$$\tilde{x} = x_+$$

$$\tilde{x} = x_-$$

$$\tilde{x} = \underset{x_+, x_-}{\operatorname{argmin}} (|\tilde{x} - x_+|, |\tilde{x} - x_-|)$$

EXERCISE: SHOW THAT
 $\forall x \in \text{Range}(F_0)$

$$\frac{|x - \tilde{x}|}{|x|} \leq \epsilon_M$$

\Rightarrow relative rounding error is bounded by the machine epsilon.