# CSE330: Numerical Methods

Topic: Representation of Numbers in Computer System and Human System

Prepared by:

Saad Bin Sohan

BRAC University

Email: sohan.academics@gmail.com
GitHub: https://github.com/saad-bin-sohan
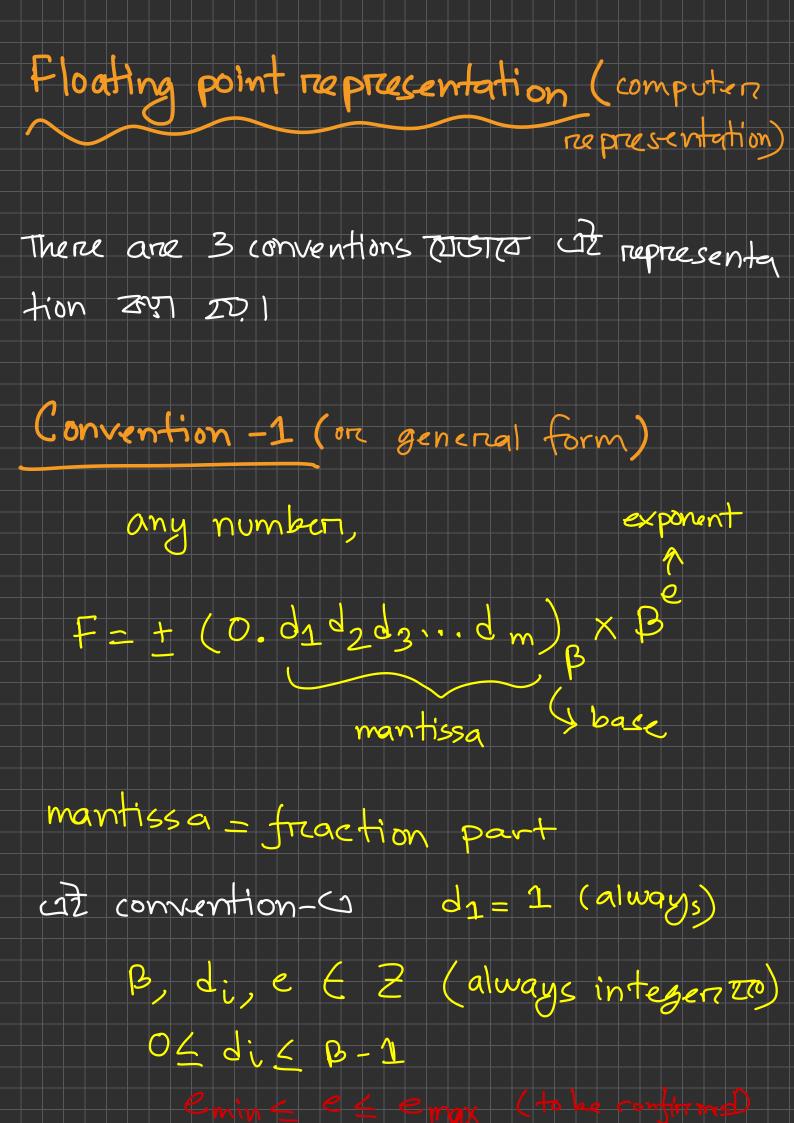
# CSE330

Numerical methods

Faculty name: Md Aquib Azmain
Initial: AQU
email : aquib.azmain@bracu.ac.bd
Room : 80615

## For Course Curriculum

Quiz- 10% (Best 4 out of 6)

Assignment- 15% (Best 4 out of 6)

Mid - 20%

Final - 30%

Lab - 20%

Attendence - 5%


All resource you need :

Numerical Analysis - II
(Lecture Notes) - Anthony Yeates

Topic: Representation of numbers in computer system (and human system)

# Floating Point Arithmetic

Fixed point representation (human represen-tation of numbers)

→ floating point

any number, $X = \pm ( d_1 d_2 \ldots d_{k-1} \cdot d_k \ldots d_n )$

$\beta$

→ base

where $d_1, d_2, d_3 \ldots, d_n \in \{0, 1, \ldots, \beta-1\}$

Ques:
$$X = + \ (10.1)_2$$

# Floating point representation (computer representation)

There are 3 conventions যেভাবে এই representation করা হয়।

## Convention -1 (or general form)

any number,

$$F = \pm \left(0.d_1 d_2 d_3 \ldots d_m\right)_\beta \times \beta^e$$

exponent → $e$

$\beta$ → base

mantissa = $\underbrace{0.d_1 d_2 d_3 \ldots d_m}$ = fraction part

এই convention-এ $\qquad d_1 = 1$ (always)

$$\beta, d_i, e \in \mathbb{Z} \text{ (always integer হও)}$$

$$0 \le d_i \le \beta - 1$$

$$e_{min} \le e \le e_{max} \quad \text{(to be confirmed)}$$

(fixed point থেকে difference টা আসলো $\beta^e$ এর জন্য)

/// convention-1 এ $d_1 = 1$ হওয়ার কারণ যেন

there is a unique representation for each number.

যেমন: $(1001.11)_2 \times 2^2$ কে

$(0.100111)_2 \times 2^6$ এ এবং

$(0.00100111)_2 \times 2^8$ এই both কে

represent logical হলেও computer এসব different representation করতে চায় না।

এজন্য $(1001.11)_2 \times 2^2$ এর convention 1

অনুযায়ী representation হলো

$$(0.100111)_2 \times 2^6$$

# Que-1

$\beta = 2$, $e_{min} = -1$

$m = 3$, $e_{max} = 2$

(i) highest possible positive number কত?

(ii) lowest possible non-negative number কত?

(iii) lowest possible negative number কত?

# Solution:

(i) $+ (0.111)_2 \times 2^2$

(ii) $+ (0.100)_2 \times 2^{-1}$

(iii) $- (0.111)_2 \times 2^2$

[* so far zero represnt possible না convention-1 ↵]

NB: decimal এর জন্য convention এ $\triangleleft$

$d_1 \geq 0$ $\qquad$ $1 \leq d_1 \leq 9$

# Convention-2 (or normalized form)

any number,

$$F = \pm (1.d_1 d_2 \ldots d_m)_\beta \times \beta^e$$

যেখানে

$\beta, d_i, e \in \mathbb{Z}$ (always integer ≥ 0)

$0 \leq d_i \leq \beta - 1$

$d_i$ হল single digit

$$e_{min} \leq e \leq e_{max}$$

**Que-2** $\beta = 2$, $e_{min} = -1$

$m = 3$, $e_{max} = 2$

convention-2 তে,

(i) highest possible positive number কত?

(ii) lowest possible non-negative number কত?

(iii) lowest possible negative number / lowest possible number considering signed bit কত?

**solution**:

(i) $+(1.111)_2 \times 2^2$

(ii) $+(1.000)_2 \times 2^{-1}$

(iii) $-(1.111)_2 \times 2^2$

# Convention-3 (de-normalized form)

any number,

$$F = \pm (0.1 d_1 d_2 \ldots d_m) \times \beta^e$$

যেখানে

$\beta, d_i, e \in \mathbb{Z}$ (always integer ≥0)

$0 \le d_i \le \beta - 1$

$d_i$ হল single digit

$\boxed{\text{Que-3}}$  $\beta = 2$ ,    $e_{min} = -1$

$m = 3$ ,    $e_{max} = 2$

convention-3 তে,

(i) highest possible positive number কত?

(ii) lowest possible non-negative number কত?

(iii) lowest possible negative number / lowest possible number considering signed bit কত?

$\boxed{\text{solution}}$ :    $'$

(i)  $+ (0.1\underbrace{1111})\times 2^2$
$\qquad\qquad\quad\; m=3$

(ii) $+(0.1\underbrace{0\,0\,0})\times 2^{-1}$
$\qquad\qquad\;\; m=3$

(iii) $-(0.11111)_2 \times 2^2$

# Que-4

$$\beta = 2 \quad , \quad e_{min} = -1$$

$$m = 3 \quad , \quad e_{max} = 2$$

Total কতগুলো নাম্বার represent করা possible এই কম্বিনেশন দিয়ে and number গুলো কী কী?

## solution :

### convention- 1

এই চারটা number for $e = -1$

$$(0.100)_2 \times 2^{-1} =$$

$$(0.101)_2 \times 2^{-1}$$

$$(0.110)_2 \times 2^{-1}$$

$$(0.111)_2 \times 2^{-1}$$

এরকম total 16টা positive number হবে ২৯ই (neg গুলা $x$ m এ ধরতে না)

$e = 0, 1, 2$ এর জন্য $4 \times 3 = 12$ টা নাম্বার পাওয়া যাবে

· · ·

· · ·

· · ·

নাম্বারগুলো যের কোন ভাগ সংখ্যালোকে decimal
এ represent করতে হবে।

**Topic:** Decimal numberগুলোকে number line এ represent করা



$\frac{1}{4}$ $\frac{5}{16}$ $\frac{3}{8}$

NB: zero এর ১০ numberর গুলো decimal এ convert করা হয় পদ্ধতি, তাকে লেখা যায় না।

# IEEE Standards

(How computer operates is in the denormalised form of numbers. But before that, computer converts every number to normalized form, then to denormalized form. And then the operations are done in de-norm form.)

the IEEE standards are

$$\beta = 2$$

52 bits for mantissa/fraction part

11 bits for exponents

1 bit for sign

(+)
———————————
64 bits in total (thats the standard how computers are designed)

## The normalized form (part in IEEE Standards)

$$\pm (1. d_1 d_2 \ldots d_{52}) \times 2^e$$

$$e_{min} = 0 \qquad e_{max} = 2^{11} - 1$$

$$= 2047$$

largest non-negative number,

$$= + (1.\underbrace{111\ldots 111}_{52 \text{ till } 1}) \times 2^{2047}$$

smallest non negative number,

$$= + (1.\underbrace{000\ldots 00}_{52 \text{ till zero}}) \times 2^0$$

# topic: exponent biasing

$2^{11}$ য 2048 টা bit এর zero একটা
(সেটা non negative এর আগে), তাই negative
-এ আছে-এর অর্ধেক যেহে 1 কম বা
$((2048/2)-1) = 1023$ টা.

positive এ $(2048/2) = 1024$ টা

zero এ কটা

negative এ $(2048/2)-1 = 1023$ টা

exponent biasing -এ,

$\rightarrow 0 \leq e \leq 2047$

$$\Rightarrow (1.\, d_1 d_2 d_3 \cdots d_{52}) \times 2^{e-1023}$$