

Chapter 2

Floating point numbers. Convergence. Bisection method. [Section 1.2]

2.0.1 Round-off errors and Computer arithmetic

Remark: Computer don't operate with real numbers!

Representations of decimal numbers(rational).

1. **Fixed Point** Let $d_k = \{0, 1, \dots, 9\}$ Then we can represent a decimal number as

$$x = (d_5 d_4 d_3 d_2 d_1 d_0 . d_{-1} d_{-2} d_{-3} d_{-4} d_{-5}) = \sum_{k=-5}^5 d_k 10^k \quad (2.1)$$

, where $m < 0 \leq M$. Similarly, for binary numbers:

$$x = b_M b_{M-1} \dots b_1 b_0 . b_{-1} b_{-2} \dots b_m = \sum_{k=m}^M b_k 2^k, \quad b_k \in \{0, 1\} \quad (2.2)$$

This is called the fixed-point representation. However, it's inconvenient for large numbers, or ones with many decimals, so it is rarely used(nowadays).

2. **Floating Point** Idea: represent the number as an integer scaled by an exponent of a fixed base.

$$12.345 = (1)12345 \times 10^{-3} \quad (2.3)$$

where (1) is the sign bit, 12345 is the significand(mantissa), 10 is the base, and -3 is the exponent.

Remark Base 10 corresponds to the "scientific notation used in calculators. Computers(usually) use base 2.

Remark The common type of floating point representation follows the IEEE 754 standard.

- (a) **Single precision** (binary 32, called float in C) 32 bits, mantisa 24 bits(approximately 7 decimal digits).
- (b) **Double precision** (binary 64, called double in C) 64 bits, mantisa 53 bits(approximately 16 decimal digits).

Example 2.0.1. The first 33 bits of π in binary is

$$\pi = 110010010000111111011010101000100 \quad (2.4)$$

To write this as a single precision floating point number, take the 24 bit rounding approximation

$$110010010000111111011011 = (1 \times 2^0 + 1 \times 2^{-1} + 0 \times 2^{-2} + \dots + 1 \times 2^{-23}) \times 2^1 \approx 3.1415928 \quad (2.5)$$

where the last bit is rounded up and the exponent is 1.

Machine epsilon(or machine precision) is an best bound for relative approximation error due to the floating point arithmetic.

The standard values of machine epsilon are

- binary 32(single precision): $\epsilon \approx 1.19 \times 10^{-7}$
- binary 64(double precision): $\epsilon \approx 2.22 \times 10^{-16}$

This value is contained in the module numpy of python (import numpy as np)

```

1 import numpy as np
2 print(np.finfo(float).eps)
3 # 2.220446049250313e-16
4 print(np.finfo(np.float32).eps)
5 # 1.1920929e-07

```

Remark

1. "Machine epsilon accuracy" is the ultimate standard for numerical algorithms. (no better accuracy can be expected)
2. Accuracy within 3 or 4 decimals from the machine epsilon is already very sufficient and more than that cannot be expected in practice.

Definition 2.0.1. Suppose p^* is an approximation of $p \in \mathbb{R}$. Then

1. The **actual error** is $p - p^*$;
2. The **absolute error** is $|p - p^*|$;
3. The **relative error** is $\frac{|p - p^*|}{|p|}$, provided $p \neq 0$.

The number p^* is said to approximate $p \neq 0$ to t significant digits if t is the largest non-negative integer such that

$$\frac{|p - p^*|}{|p|} \leq 5 \times 10^{-t} \quad (2.6)$$

Remark It seems that there is a typo in the lecture note where it mistakes \leq for $<$ in the definition above. We have modified it.

Definition 2.0.2 (Floating Point Arithmetic). Denote by $fl(x)$ the floating point approximation of x . Then

$$x \oplus y = fl(fl(x) + fl(y)) \quad (2.7)$$

$$x \otimes y = fl(fl(x) \cdot fl(y)) \quad (2.8)$$

$$x \ominus y = fl(fl(x) - fl(y)) \quad (2.9)$$

$$x \oslash y = fl\left(\frac{fl(x)}{fl(y)}\right), \quad y \neq 0 \quad (2.10)$$

2.0.2 Convergence

Definition 2.0.3 (Rates of Convergence). Suppose $(\alpha_n), n = 0, 1, 2, \dots$ is a sequence of real numbers such that $\lim_{n \rightarrow \infty} \alpha_n = \alpha \in \mathbb{R}$ and (β_n) is a sequence such that $\lim_{n \rightarrow \infty} \beta_n = 0$. If there exist $K > 0, N \geq 0$ such that $|\alpha_n - \alpha| \leq K|\beta_n|$ for all $n \geq N$, then we say that (α_n) converges to α with rate $O(\beta_n)$.

Remark Usually $\beta_n = \frac{1}{n^p}$ for some $p > 0$.

2.0.3 Bisection method

Recall: Intermediate Value Theorem(Bolzano's theorem): Suppose $f : [a, b] \rightarrow \mathbb{R}$ is continuous, and there exists K such that either $K \in (f(a), f(b))$ if $f(a) < f(b)$ or $K \in (f(b), f(a))$ if $f(a) > f(b)$. Then there exists $c \in (a, b)$ such that $f(c) = K$. As a corollary, if $f(a)f(b) < 0$, then there exists $c \in (a, b)$ such that $f(c) = 0$.

Application Finding an interval $[a, b]$ that contains a solution of $x - 2^{-x} = 0$. For $x = 0, f(0) = -1 < 0$. For $x = 1, f(1) = 0.5 > 0$. So by the corollary there exists $c \in (0, 1)$ such that $f(c) = 0$.

The following is an more advanced application.

Algorithm 2.0.1 (Bisection Method). Let $f : [a, b] \rightarrow \mathbb{R}$ be continuous and suppose that $f(a)f(b) < 0$. Then

1. Set $n = 1, a_1 = a, b_1 = b$.
2. Compute $c_n = \frac{a_n + b_n}{2}$. If $f(c_n) = 0$, stop. Otherwise, go to step 3.

3. If $f(a_n)f(c_n) < 0$, set $a_{n+1} = a_n, b_{n+1} = c_n$. If $f(b_n)f(c_n) < 0$, set $a_{n+1} = c_n, b_{n+1} = b_n$. Increment n by 1 and go to step 2.

Example 2.0.2. Let $f(x) = \cos(x) - 2x, [a, b] = [-8, 10]$. Then $f(-8) > 0, f(10) < 0$. So we can apply the bisection method.

Theorem 2.0.1. Suppose $f \in C([a, b])$ and $f(a)f(b) < 0$. Then the sequence (c_n) generated by the bisection method approximates the zero $p \in (a, b)$ of f with

$$|c_n - p| \leq \frac{b - a}{2^n}, n \geq 1 \quad (2.11)$$

Proof Let $n \geq 1$. Then $b_n - a_n = \frac{b - a}{2^{n-1}}$ and $p \in (a_n, b_n)$ where a_n, b_n are the endpoints of the interval at the n -th step. Because $c_n = \frac{a_n + b_n}{2}$, it follows that

$$|c_n - p| \leq \frac{b_n - a_n}{2} = \frac{b - a}{2^n} \rightarrow 0 \quad (2.12)$$

as $n \rightarrow \infty$. ■

Remark As $|p_n - p| \leq (b - a)2^{-n}$. So the convergence converges with the rate $O(2^{-n})$. Note: Here we can take $K = b - a$.

Remark To avoid taking off points by graders, we're expected to show(at lest state) the conditions of the theorems we use. For example we need to say $K = b - a$ as in the previous proof of the theorem.