

1. Floating-point arithmetic

- **Fixed point representation:**

$$x = \pm (d_1 d_2 \dots d_{k-1} \cdot d_k \dots d_n)_\beta$$

- **Floating-point representation:**

$$x = (0. d_1 \dots d_{k-1})_\beta \beta^{d_k \dots d_n - B}$$

where B is an **exponent bias**.

- If $d_1 \neq 0$ then the floating point system is **normalised** and each float has a unique representation.
- **binary64**: stored as

$$s e_{10} \dots e_0 d_1 \dots d_{52}$$

where s is the **sign** (0 for positive, 1 for negative), $e_{10} \dots e_0$ is the **exponent**, and $d_1 \dots d_{52}$ is the **mantissa**. The bias is 1023. The number represented is

$$\begin{cases} (-1)^s (1. d_1 \dots d_{52})_2 2^e & \text{if } e \neq 0 \text{ or } 2047 \\ (-1)^s (0. d_1 \dots d_{52})_2 2^{-1022} & \text{if } e = 0 \end{cases}$$

where $e = (e_{10} \dots e_0)_2$ $e = 2047$ is used to store NaN, $\pm\infty$. The first case $e \neq 0$ is a **normal** representation, the $e = 0$ case is a **subnormal representation**.

- Floating-point numbers have finite range and precision.
- **Underflow**: where floating point calculation result is smaller than smallest representable float. Result is set to zero.
- **Overflow**: where floating point calculation result is larger than largest representable float. **Floating-point exception** is raised.
- **Machine epsilon** ε_M : difference between smallest representable number greater than 1 and 1. $\varepsilon_M = \beta^{-k+1}$.
- $\text{fl}(x)$ maps real numbers to floats.
- **Chopping**: rounds towards zero. Given $x = (0. d_1 \dots d_k d_{k+1} \dots)_\beta \cdot \beta^e$, if the float has k mantissa digits, then

$$\text{fl}_{\text{chop}}(x) = (0. d_1 \dots d_k)_\beta \cdot \beta^e$$

- **Rounding**: rounds to nearest. Given $x = (0. d_1 \dots d_k d_{k+1} \dots)_\beta \cdot \beta^e$, if the float has k mantissa digits, then

$$\tilde{\text{fl}}_{\text{round}}(x) = \begin{cases} (0. d_1 \dots d_k)_\beta \cdot \beta^e & \text{if } \rho < \frac{1}{2} \\ \left((0. d_1 \dots d_k)_\beta + \beta^{-k} \right) \cdot \beta^e & \text{if } \rho \geq \frac{1}{2} \end{cases}$$

where $\rho = (0. d_{k+1} \dots)$.

- **Relative rounding error:**

$$\varepsilon_x = \frac{\text{fl}(x) - x}{x} \iff \text{fl}(x) = x(1 + \varepsilon_x)$$

- $$\left| \frac{\text{fl}_{\text{chop}} - x}{x} \right| \leq \beta^{-k+1}, \quad \left| \frac{\tilde{\text{fl}}_{\text{round}}(x) - x}{x} \right| \leq \frac{1}{2} \beta^{-k+1}$$

- **Round-to-nearest half-to-even:** fairer rounding than regular rounding for discrete values. In the case of a tie, round to nearest even integer:

$$\text{fl}_{\text{round}}(x) = \begin{cases} (0.d_1 \dots d_k)_\beta \cdot \beta^e & \text{if } \rho < \frac{1}{2} \text{ or } (\rho = \frac{1}{2} \text{ and } d_k \text{ is even}) \\ \left((0.d_1 \dots d_k)_\beta + \beta^{-k} \right) \cdot \beta^e & \text{if } \rho > \frac{1}{2} \text{ or } (\rho = \frac{1}{2} \text{ and } d_k \text{ is odd}) \end{cases}$$

- $x \oplus y = \text{fl}(\text{fl}(x) + \text{fl}(y))$ and similarly for \otimes, \ominus, \odot .
- Relative error in $x \pm y$ can be large:

$$\text{fl}(x) \pm \text{fl}(y) - (x \pm y) = x(1 + \varepsilon_x) \pm y(1 + \varepsilon_y) - (x \pm y) = x\varepsilon_x \pm y\varepsilon_y$$

so relative error is

$$\frac{x\varepsilon_x \pm y\varepsilon_y}{x \pm y}$$

- In general, $x \oplus (y \oplus z) \neq (x \oplus y) \oplus z$
- For some computations, can avoid round-off errors (usually caused by subtraction of numbers close in value) e.g. instead of

$$x = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

compute

$$x = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \cdot \frac{-b - \sqrt{b^2 - 4ac}}{-b - \sqrt{b^2 - 4ac}} = \frac{-2c}{b + \sqrt{b^2 - 4ac}}$$

2. Polynomial Interpolation

- \mathcal{P}_n is set of polynomials of degree $\leq n$.
- $\text{conv}\{x_0, \dots, x_n\}$ is smallest closed interval containing $\{x_0, \dots, x_n\}$.
- **Taylor's theorem:** for function f , if for $t \in \mathcal{P}_n$, $t^{(j)}(x_0) = f^{(j)}(x_0)$ for $j \in \{0, \dots, n\}$ then

$$f(x) - t(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} (x - x_0)^{n+1}$$

for some $\xi \in \text{conv}\{x_0, x\}$ (**Lagrange form of remainder**).

- **Polynomial interpolation:** given nodes $\{x_j\}_{j=0}^n$ and function f , there exists unique $p \in \mathcal{P}_n$ such that p interpolates f : $p(x_j) = f(x_j)$ for $j \in \{0, \dots, n\}$.
- **Cauchy's theorem:** let $p \in \mathcal{P}_n$ interpolate f at $\{x_j\}_{j=0}^n$, then

$$\forall x \in \text{conv}\{x_j\}, f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} (x - x_0) \cdots (x - x_n) \quad \text{for some } \xi \in \text{conv}\{x_j\}$$

- **Chebyshev polynomials:**

$$T_n(x) = \cos(n \cos^{-1}(x)), \quad x \in [-1, 1]$$

- $T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x)$.
- Roots of $T_n(x)$ are $x_j = \cos(\pi(j + \frac{1}{2}) / n)$ for $j \in \{0, \dots, n-1\}$. Local extrema at $y_j = \cos(j\pi / n)$ for $j \in \{0, \dots, n-1\}$.
- Let $\omega_n(x) = (x - x_0) \cdots (x - x_n)$, $\{x_j\}_{j=0}^n \subset [-1, 1]$ (if $\{x_j\} \not\subset [-1, 1]$ so interval is $[a, b]$, then we can map $x_j \rightarrow a + \frac{1}{2}(x_j + 1)(b - a)$). Then $\sup_{x \in [-1, 1]} |\omega_n(x)|$ attains its min value iff $\{x_j\}$ are zeros of $T_{n+1}(x)$. Also,

$$2^{-n} \leq \sup_{x \in [-1, 1]} |\omega_n(x)| < 2^{n+1}$$

- **Convergence theorem:** let $f \in C^2([-1, 1])$, $\{x_j\}_{j=0}^n$ be zeros of Chebyshev polynomial $T_{n+1}(x)$ and $p_n \in \mathcal{P}_n$ interpolate f at $\{x_j\}$. Then

$$\sup_{x \in (-1, 1)} |f(x) - p_n(x)| \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

- **Weierstrass' theorem:** let $f \in C^0([a, b])$. $\forall \varepsilon > 0$, exists polynomial p such that

$$\sup_{x \in (a, b)} |f(x) - p(x)| < \varepsilon$$

- **Lagrange construction:** basis polynomials given by

$$L_k(x) = \prod_{j \neq k} \frac{x - x_j}{x_k - x_j}$$

satisfy $L_k(x_j) = \delta_{jk}$. Then

$$p(x) = \sum_{k=0}^n L_k(x) f(x_k)$$

interpolates f at $\{x_j\}$.

- **Note:** Lagrange construction not often used due to computational cost and as we have to recompute from scratch if $\{x_j\}$ is extended.
- **Divided difference operator:**

$$[x_j]f := f(x_j)$$

$$[x_j, x_k]f := \frac{[x_j]f - [x_k]f}{x_j - x_k}, \quad [x_k, x_k]f := \lim_{y \rightarrow x_k} [x_k, y] = f'(x_k)$$

$$[x_j, \dots, x_k, y, z]f := \frac{[x_j, \dots, x_k, y]f - [x_j, \dots, x_k, z]f}{y - z}$$

These can be computed incrementally as new nodes are added.

- **Newton construction:** Interpolating polynomial p is

$$p(x) = [x_0]f + (x - x_0)[x_0, x_1]f + (x - x_0)(x - x_1)[x_0, x_1, x_2]f \\ + \cdots + (x - x_0) \cdots (x - x_{n-1})[x_0, \dots, x_n]f$$

- **Hermite construction:** for nodes $\{x_j\}_{j=0}^n$, exists unique $p_{2n+1} \in \mathcal{P}_{2n+1}$ that interpolates f and f' at $\{x_j\}$. Can be found using Newton construction, using nodes $(x_0, x_0, x_1, x_1, \dots, x_n, x_n)$. Generally, if $p'(x_k) = f'(x_k)$ is needed, include x_k twice. If $p^{(n)}(x_k) = f^{(n)}(x_k)$ is needed, include x_k $n + 1$ times.
- If y_0, \dots, y_k is permutation of x_0, \dots, x_k then $[y_0, \dots, y_k]f = [x_0, \dots, x_k]f$.
- Interpolating error is

$$f(x) - p(x) = (x - x_0) \cdots (x - x_n)[x_0, \dots, x_n, x]f$$

which gives

$$[x_0, \dots, x_{n-1}, x]f = \frac{f^{(n+1)}(\xi)}{(n+1)!}$$

- **Range reduction:** when computing a function e.g. $f(x) = \arctan(x)$, $f(-x) = -f(x)$ and $f(1/x) = \frac{\pi}{2} - f(x)$ so only need to compute for $x \in [0, 1]$.

3. Root finding

- **Intermediate value theorem:** if f continuous on $[a, b]$ and $f(a) < c < f(b)$ then exists $x \in (a, b)$ such that $f(x) = c$.
- **Bisection:** let $f \in C^0([a_n, b_n])$, $f(a_n)f(b_n) < 0$. Then set $m_n = (a_n + b_n) / 2$ and

$$(a_{n+1}, b_{n+1}) = \begin{cases} (m_n, b_n) & \text{if } f(a_n)f(m_n) > 0 \\ (a_n, m_n) & \text{if } f(b_n)f(m_n) > 0 \end{cases}$$

Then:

- $b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n)$.
- By intermediate value theorem, exists $p_n \in (a_n, b_n)$ with $f(p_n) = 0$.
- $|p_n - m_n| \leq 2^{-(n+1)}(b_0 - a_0)$.
- **False position:** same as bisection except set m_n as x intercept of line from $(a_n, f(a_n))$ to $(b_n, f(b_n))$:

$$m_n = b_n - \frac{f(b_n)}{f(b_n) - f(a_n)}(b_n - a_n)$$

- Bisection and false position are **bracketing methods**. Always work but slow.
- **Fixed-point iteration:** rearrange $f(x_*) = 0$ to $x_* = g(x_*)$ then iterate $x_{n+1} = g(x_n)$.
- f is **Lipschitz continuous** if for some L ,

$$|f(x) - f(y)| \leq L|x - y|$$

- Space of Lipschitz functions on X is $C^{0,1}(X)$.
- Smallest such L is **Lipschitz constant**.
- Every Lipschitz function is continuous.
- Lipschitz constant is bounded by derivative:

$$\sup_{x \neq y} \frac{|f(x) - f(y)|}{|x - y|} \leq \sup_x |f'(x)|$$

- f is **contraction** if Lipschitz constant $L < 1$.

- **Contraction mapping or Banach fixed point theorem:** if g is a contraction and $g(X) \subset X$ (g maps X to itself) then:
 - Exists unique solution $x_* \in X$ to $g(x) = x$ and
 - The fixed point iteration method converges $x_n \rightarrow x_*$.
- **Local convergence theorem:** Let $g \in C^1([a, b])$ have fixed point $x_* \in (a, b)$ with $|g'(x_*)| < 1$. Then with x_0 sufficiently close to x_* , fixed point iteration method converges to x_* .
 - If $g'(x_*) > 0$, $x_n \rightarrow x_*$ monotonically.
 - If $g'(x_*) < 0$, $x_n - x_*$ alternates in sign.
 - If $|g'(x_*)| > 1$, iteration method almost always diverges.
- $x_n \rightarrow x_*$ with **order at least $\alpha > 1$** if

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x_*|}{|x_n - x_*|^\alpha} = \lambda < \infty$$

If $\alpha = 1$, then $\lambda < 1$ is required.

- **Exact order of convergence** of $x_n \rightarrow x_*$:

$$\alpha := \sup \left\{ \beta : \lim_{n \rightarrow \infty} \frac{|x_{n+1} - x_*|}{|x_n - x_*|^\beta} < \infty \right\}$$

Limit must be < 1 for $\alpha = 1$.

- Convergence is **superlinear** if $\alpha > 1$, **linear** if $\alpha = 1$ and $\lambda < 1$, **sublinear** otherwise.
- If $g \in C^2$, then with fixed point iteration,

$$\frac{|x_{n+1} - x_*|}{|x_n - x_*|} \rightarrow |f'(x_*)| \text{ as } n \rightarrow \infty$$

so $x_n \rightarrow x_*$ superlinearly if $g'(x_*) = 0$ and linearly otherwise.

- If $g \in C^N$, fixed point iteration converges with order $N > 1$ iff

$$g'(x_*) = \dots = g^{(N-1)}(x_*) = 0, \quad g^{(N)}(x_*) \neq 0$$

- **Newton-Raphson:** fixed point iteration with $g(x) = x - f(x) / f'(x)$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

- For Newton-Raphson, $g'(x_*) = 0$ so quadratic convergence.
- Can use Newton-Raphson to solve $1 / x - b = 0$:

$$x_{n+1} = x_n - \frac{1 / x_n - b}{-1 / x_n^2} = x_n(2 - bx_n)$$

- **Newton-Raphson in d dimensions:**

$$\underline{x}_{n+1} = \underline{x}_n - (Df)^{-1}(\underline{x}_n) f(\underline{x}_n)$$

where Df is **Jacobian**.

- **Secant method:** approximate $f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$ with Newton-Raphson:

$$x_{n+1} = x_n - \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n)$$

4. Numerical differentiation

- **Taylor expansion:**

$$f(x \pm h) = f(x) \pm hf'(x) + \frac{h^2}{2!}f''(x) \pm \frac{h^3}{3!}f'''(x) + \dots$$

- **Forward difference approximation:**

$$f'(x) = \frac{f(x+h) - f(x)}{h} - \frac{h}{2}f''(\xi), \quad \xi \in \text{conv}\{x, x+h\}$$

with $h > 0$.

- **Backward difference approximation:** forward difference but with $h < 0$.
- **Centred difference approximation:**

$$f'(x) = \frac{f(x+h) - f(x-h)}{2h} - \frac{h^2}{12}(f'''(\xi_-) + f'''(\xi_+)), \quad \xi_{\pm} \in [x-h, x+h]$$

- **Richardson extrapolation:** for approximation of $R(x; 0)$ of the form

$$R(x; h) = R^{(1)}(x; h) = R(x; 0) + a_1(x)h + a_2(x)h^2 + a_3(x)h^3 + \dots$$

we have

$$R^{(1)}(x; h/2) = R(x; 0) + a_1(x)\frac{h}{2} + a_2(x)\frac{h^2}{4} + a_3(x)\frac{h^3}{8} + \dots$$

This gives **second order approximation**:

$$R^{(2)}(x; h) = 2R^{(1)}(x; h/2) - R^{(1)}(x; h) = R(x; 0) - a_2(x)\frac{h^2}{2} + \dots$$

Similarly,

$$R^{(3)}(x; h) = \frac{4R^{(2)}(x; h/2) - R^{(2)}(x; h)}{3} = R(x; 0) + \tilde{a}_3(x)h^3 + \dots$$

is **third order approximation**. Generally,

$$R^{(n+1)}(x; h) = \frac{2^n R^{(n)}(x; h/2) - R^{(n)}(x; h)}{2^n - 1} = R(x; 0) + O(h^{n+1})$$