

$$\text{flp}(x) = x(1 \pm \varepsilon)$$

Consequence :

ARITHMETICS:  $(x + y) + z = x + (y + z)$

F.P. ARITHMETICS:  $\text{flp}(x + y) + \text{flp}(z) \neq \text{flp}(x) + \text{flp}(y + z)$

EXAMPLE:  $x = 1 \quad y = 10^{-16} \quad z = 10^{-16}$

in F.P.:  $(x + y) + z = 1 \mid x + (y + z) =$

$\underbrace{1}_{1} \mid \underbrace{1}_{1} \quad \underbrace{2 \cdot 10^{-16}}_{2 \cdot 10^{-16}}$

$\downarrow$

$1.000\,000\,000\,000\,000\,2$

## IEEE 754 STANDARD

→ FLOATING point implementation

GOAL: reproducibility of results on different machines

FEATURES

① Binary format + HIDDEN BIT TECHNIQUE  
[ $\beta = 2$ ]

Ex:  $F(\beta = 2, t = 3, L = -1, U = 2)$

$[-1, 0, 1, 2]$

$F = \{0\} \cup \left\{ \left( \sum_{k=0}^{t-1} d_k \beta^{-k} \right) \beta^p \right\}$

MANTISSA:  $d_0 \neq 0$

	$d_0$	$d_1$	$d_2$
$\beta = 2$	1	0	0
	1	0	1
	1	1	0
	1	1	1

$\beta = 10$

$\Rightarrow$

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

$1 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} = 1,25$

$1 \cdot 2^0 + 1 \cdot 2^{-1} + 0 = 1,5$

$1 \cdot 2^0 + 1 \cdot 2^{-1} + 1 \cdot 2^{-2} = 1,75$

VALUE of LAST BIT  $\rightarrow -2$

SPACING  $\Delta = 0,25$



Hidden bit technique: since  $d_0 \neq 0$ , with  $\beta = 2$

$$d_0 < \cancel{1} \Rightarrow d_0 \text{ ALWAYS} = 1$$

### EXPONENT PART

$$\begin{matrix} L = -1 \\ U = 2 \end{matrix} \rightarrow \beta^p = \{ 2^{-1}; 2^0; 2^1; 2^2 \}$$

$$F(\beta = 2, t = 3, L = -1, U = 2)$$

$$\begin{aligned} p = -1 &\rightarrow \beta^{-1} = 2^{-1} = 0,5 \\ \vdots \\ p = 1 &\rightarrow \beta^1 = 2 \end{aligned}$$

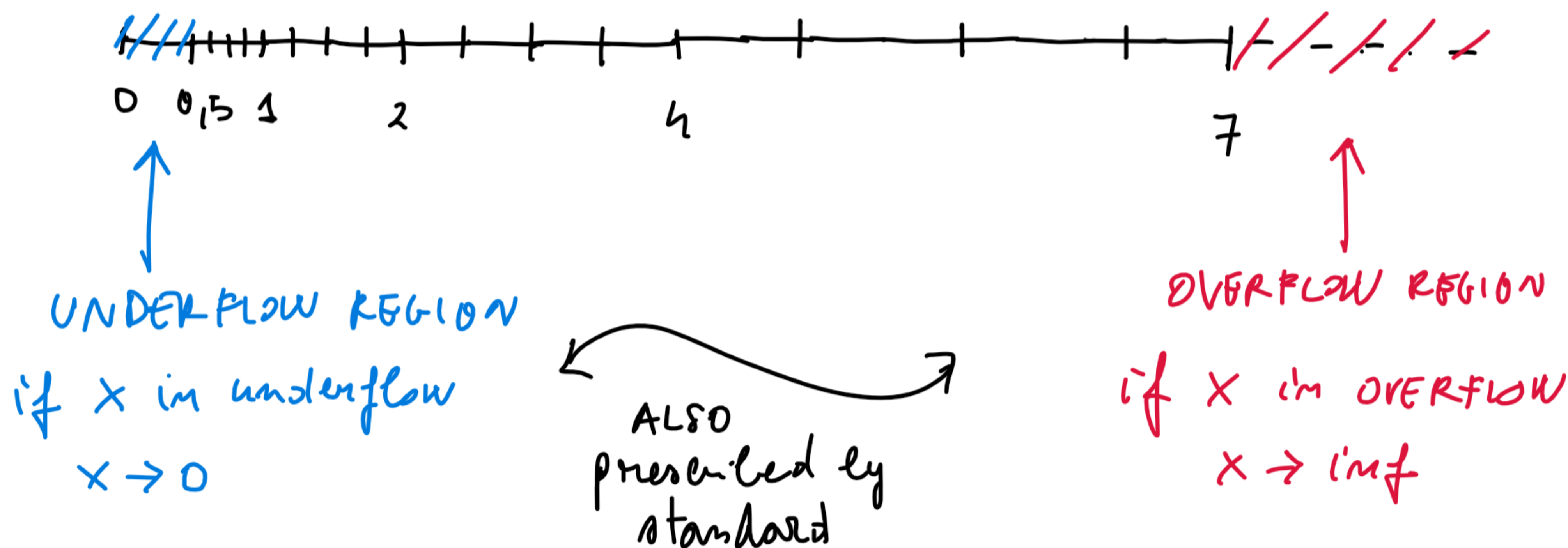
$p = -1$	$p = 0$	$p = 1$	$p = 2$
0,5	1	2	4
0,625	1,25	2,5	5
0,75	1,5	3	6
0,875	1,75	3,5	7

$$\underbrace{\quad}_{\Delta = 0,125} \quad \underbrace{\quad}_{\Delta = 0,25} \quad \underbrace{\quad}_{\Delta = 0,5} \quad \underbrace{\quad}_{\Delta = 1}$$

$$\begin{aligned} \varepsilon &= \frac{1}{2} \beta^{1-t} = \frac{1}{2} \beta^{-2} = 0,125 \\ &\quad \downarrow \quad \quad \quad \uparrow \\ &\quad \quad \quad = 2 \end{aligned}$$

$X \Rightarrow X(1 \pm \varepsilon)$   
12,5% UNCERTAINTY over numbers

REL. ERR



## ② PRECISION LEVELS

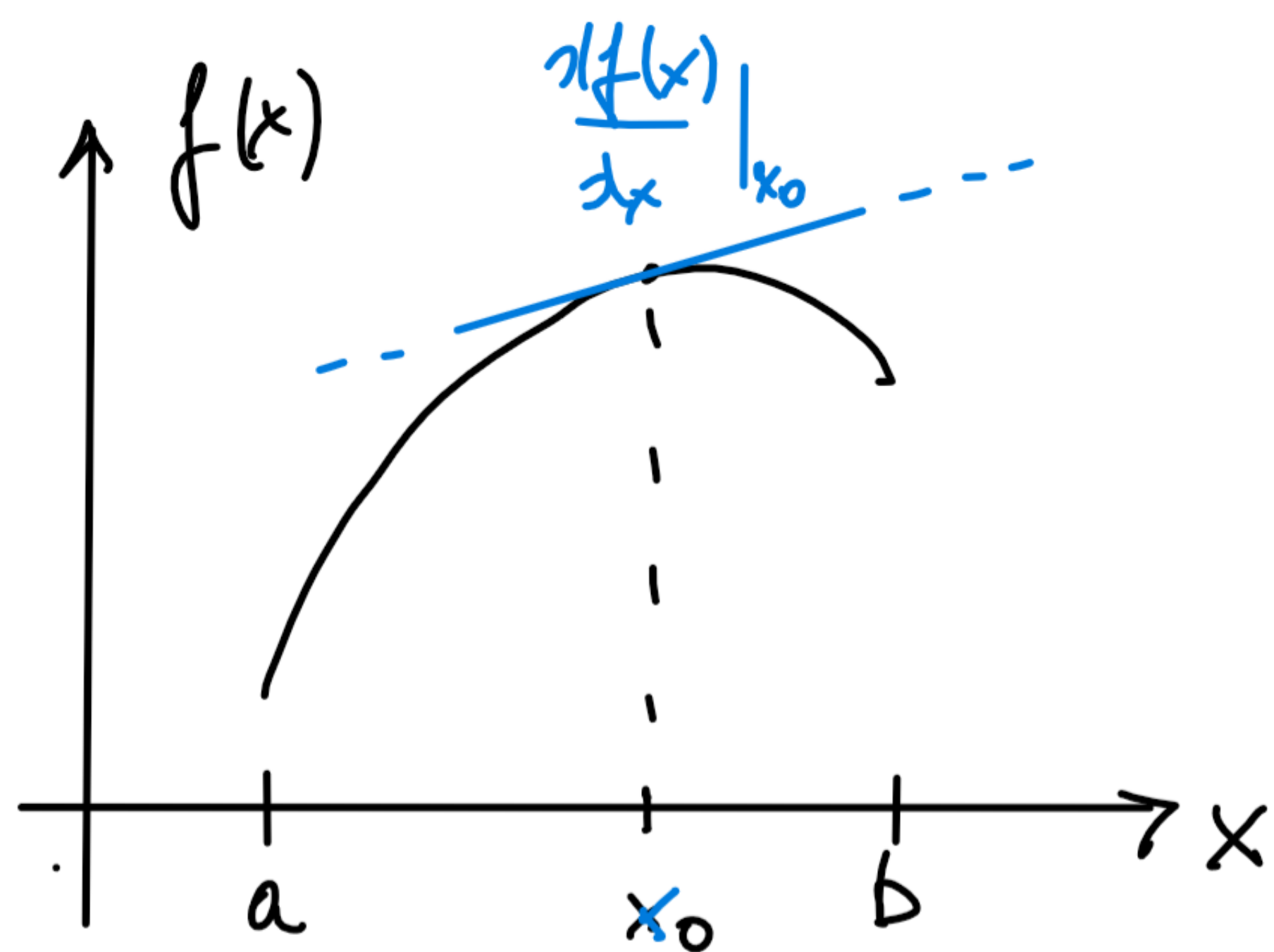
	SINGLE PRECISION 32 - BIT	DOUBLE PRECISION 64 - BIT
SIGN	1 BIT	1 BIT
EXPONENT	8 BIT $(2^8 = 256 \text{ } p \in [-126; 127])$	11 BIT $(2^{11} = 2048 \text{ } p \in [-1022; 1023])$
MANTISSA	23 BIT	52 BIT
$\epsilon$	$\beta^{1-(23+1)} = \beta^{-23} \sim 10^{-7}$	$\beta^{1-(52+1)} = \beta^{-52} \sim 10^{-16}$
$\beta^{1-t}$		
$\beta^{\uparrow}$		
# of digits for mantissa		



# NUMERICAL DIFFERENTIATION

$$f(x) \in [a, b]$$

$$\left. \frac{df(x)}{dx} \right|_{x_0} = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$$



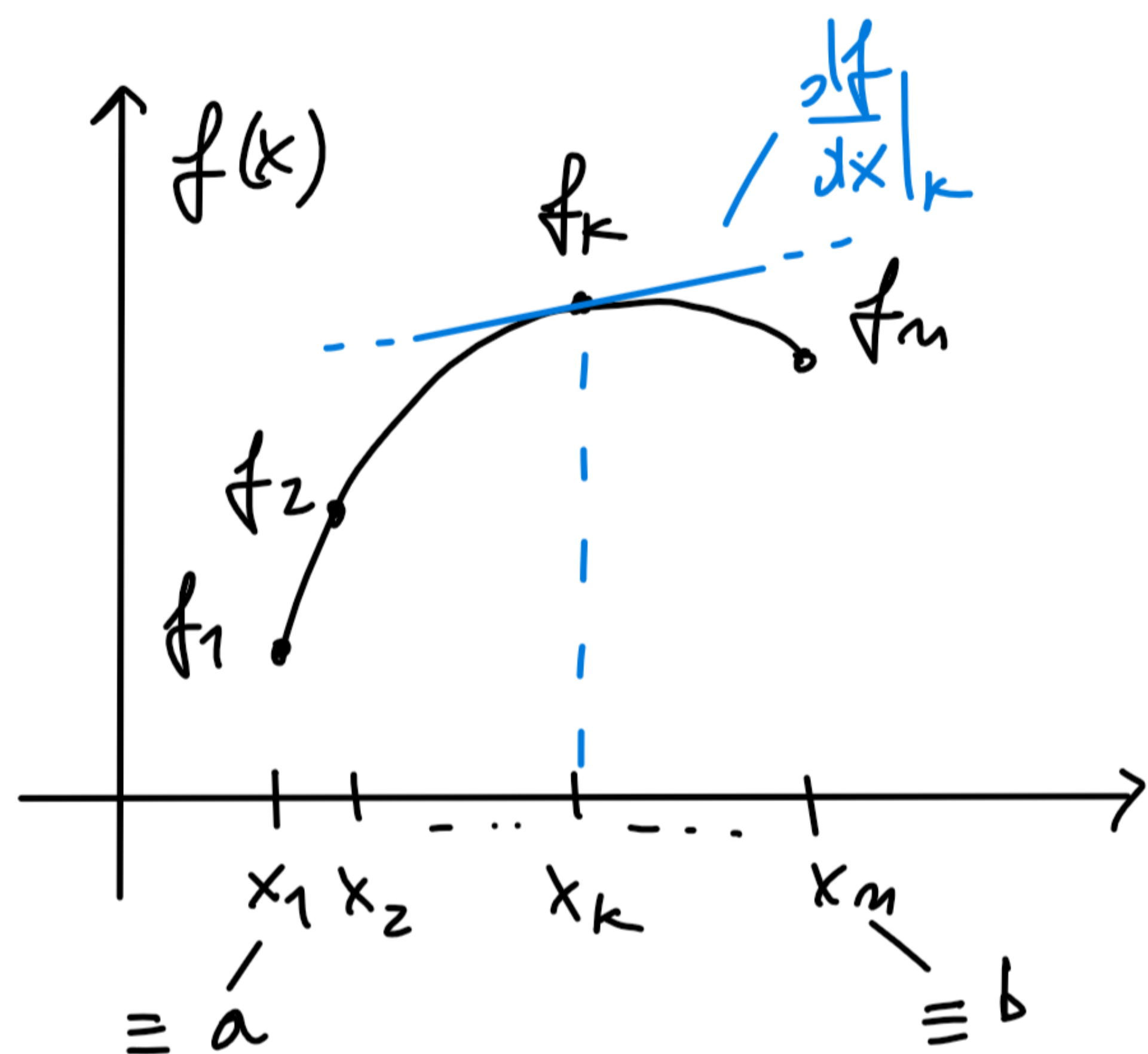
We need ALGEBRAIC APPROXIMATION of derivative

- Divide  $[a, b]$  into set of  $n$ -nodes

$$[x_1, x_2, x_3, \dots, x_k, \dots, x_n]$$

$$x_1 \equiv a \quad x_n \equiv b$$

$$f(x_1) = f_1; \quad f(x_k) = f_k; \quad f(x_n) = f_n$$



Taylor series around  $x_k$

$$f(x) = f_k + \left. \frac{df(x)}{dx} \right|_k (x - x_k) + \frac{1}{2} \left. \frac{d^2 f(x)}{dx^2} \right|_k (x - x_k)^2 + \dots + \underbrace{\dots}_{O(x - x_k)^3}$$

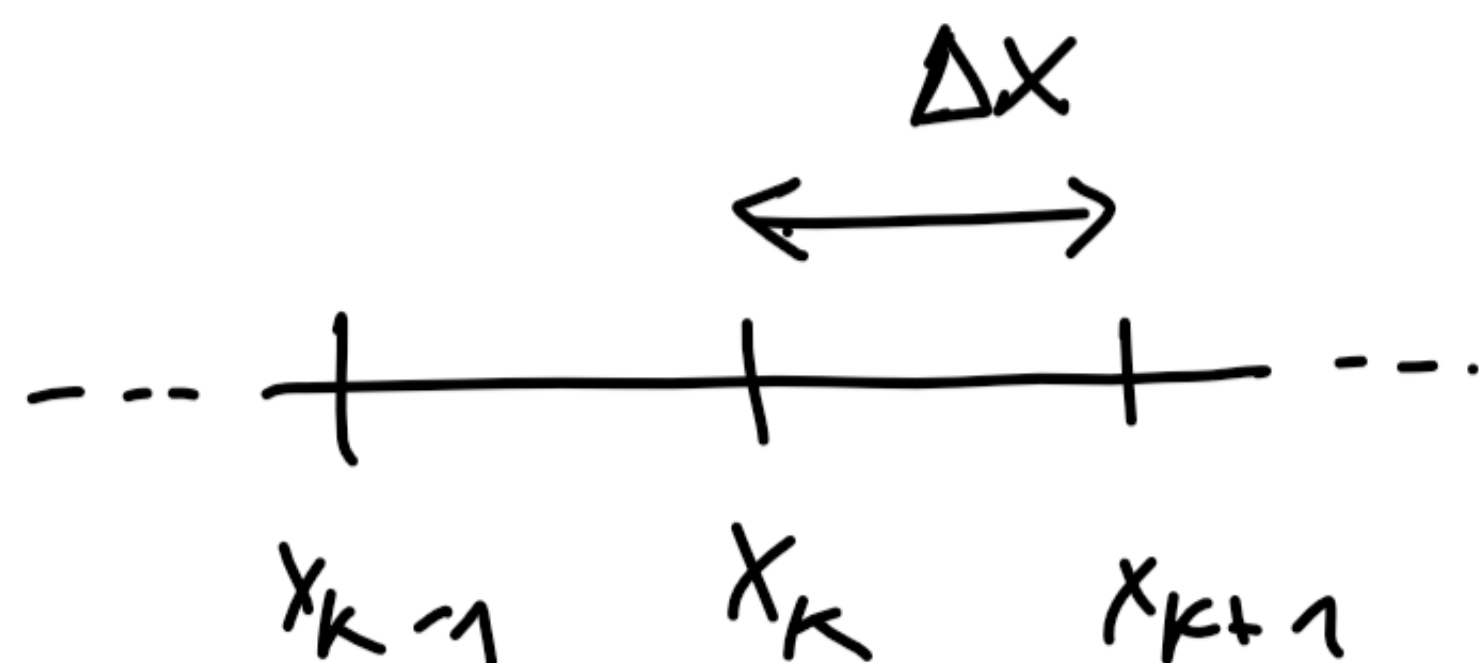
$$\exists M \in \mathbb{R}^+, |err| \leq M(x - x_k)^3, \text{ when } (x - x_k) \rightarrow 0$$

$$|err| = \left| f(x) - \left[ f_k + \left. \frac{df}{dx} \right|_k (x - x_k) + \frac{1}{2} \left. \frac{d^2 f}{dx^2} \right|_k (x - x_k)^2 \right] \right|$$

"the error SCALES AS  $(x - x_k)^3$  when  $(x - x_k)$  goes to zero

Assume equal spacing:

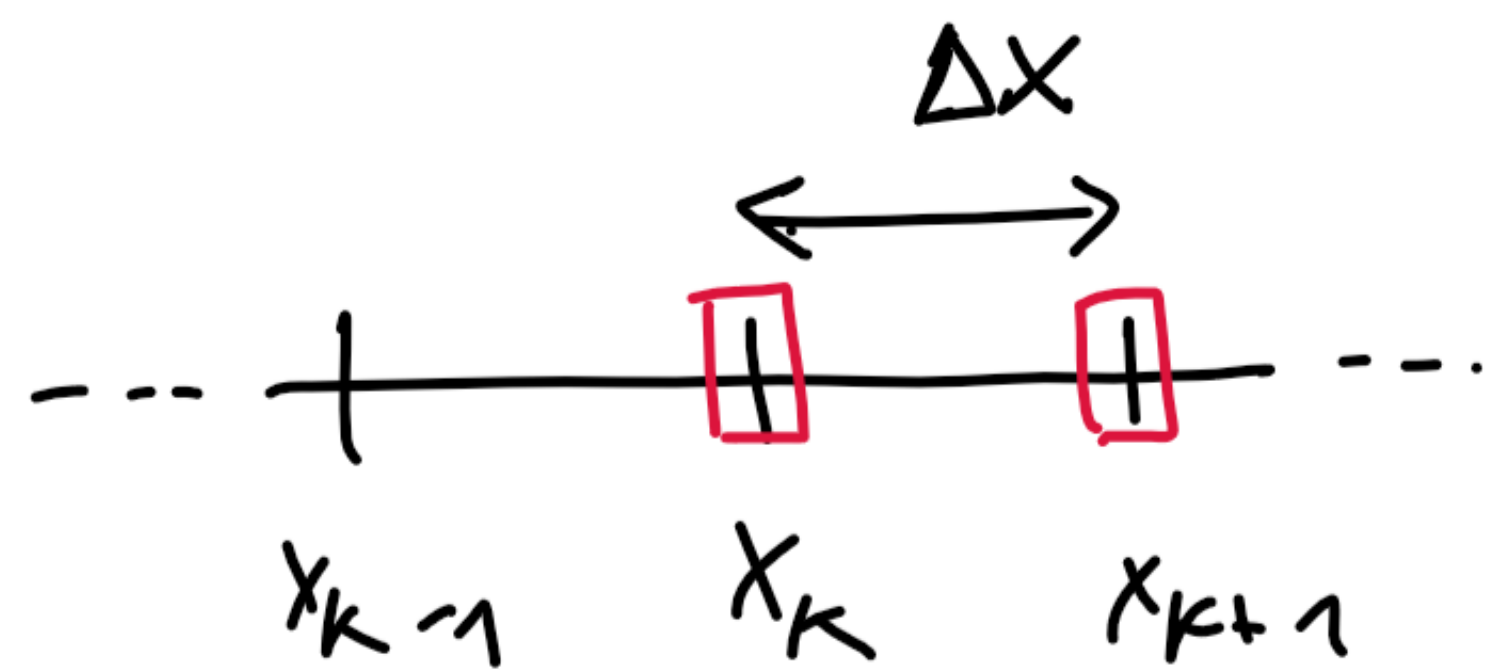
$$\Delta x = x_{k+1} - x_k = \frac{b - a}{n - 1}$$





from  $x_k$  to  $x_{k+1}$

$$f_{k+1} = f_k + \frac{df(x)}{dx} \bigg|_k (\overbrace{x_{k+1} - x_k}^{\Delta x}) + \frac{1}{2} \frac{d^2 f(x)}{dx^2} \Delta x^2 + O(\Delta x^3)$$



TRUNCATION TO SECOND-ORDER

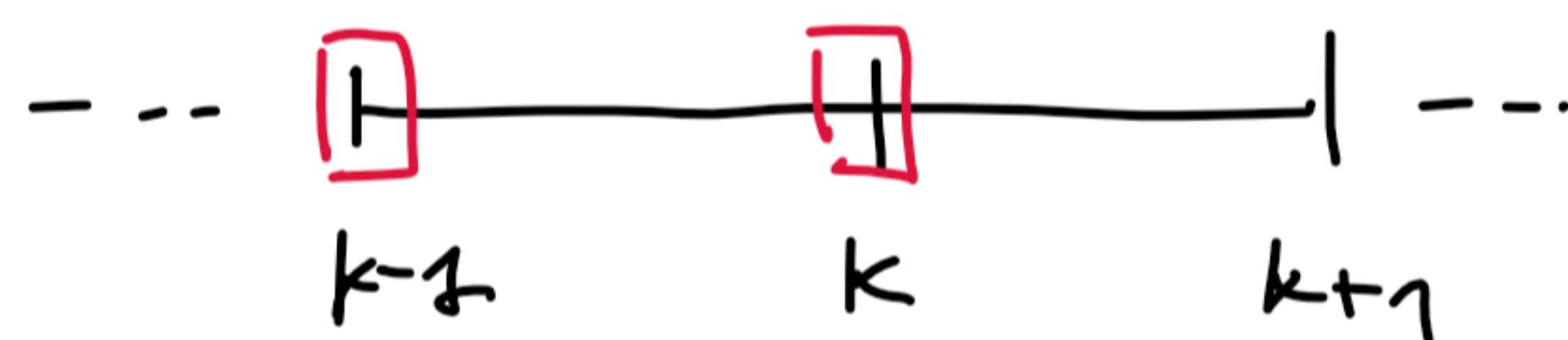
$$f_{k+1} = f_k + \frac{df(x)}{dx} \bigg|_k \Delta x + O(\Delta x^2) \quad (1)$$

$$\frac{df(x)}{dx} \bigg|_k = \frac{f_{k+1} - f_k}{\Delta x} + O(\Delta x) \approx \frac{f_{k+1} - f_k}{\Delta x}$$

FORWARD finite  
difference  
formula

FIRST-ORDER ACCURATE  $\rightarrow$  if  $\Delta x' = \frac{\Delta x}{2} \rightarrow |err'| = \left| \frac{err}{2} \right|$   
 $\downarrow$   
 exponent of  $O(\Delta x^1)$

TAYLOR for  $x_k$  to  $x_{k-1}$



$$f_{k-1} = f_k - \frac{df(x)}{dx} \bigg|_k \Delta x + \frac{1}{2} \frac{d^2 f(x)}{dx^2} \Delta x^2 + O(\Delta x^3)$$

$\uparrow$   $(x_{k-1} - x_k)$        $\uparrow$   $(x_{k-1} - x_k)^2$   
 $\leq 0$

2<sup>nd</sup> order.

$$f_{k-1} = f_k - \frac{df(x)}{dx} \bigg|_k \Delta x + O(\Delta x^2) \quad (2)$$

$$\frac{df(x)}{dx} \bigg|_k = \frac{f_k - f_{k-1}}{\Delta x} + O(\Delta x) \approx \frac{f_k - f_{k-1}}{\Delta x}$$

BACKWARD finite  
difference formula

[ACCURACY: first-order]

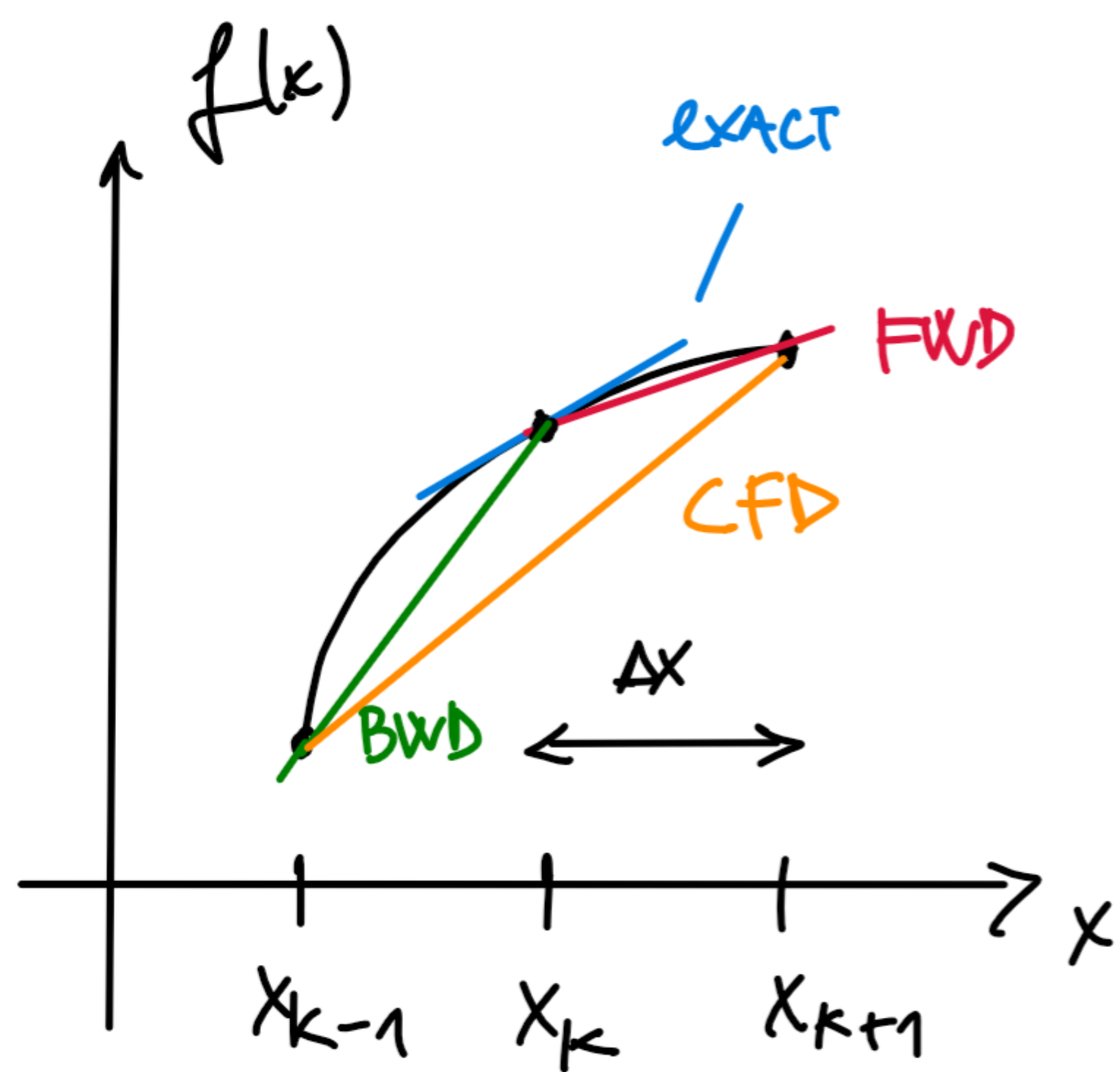
if subtract (1) - (2), using THIRD-ORDER accurate expansions

$$f_{k+1} - f_{k-1} = \cancel{f_k} + \frac{df(x)}{dx} \bigg|_k \Delta x + \frac{1}{2} \frac{d^2 f(x)}{dx^2} \Delta x^2 - \left[ \cancel{f_k} - \frac{df(x)}{dx} \bigg|_k \Delta x + \frac{1}{2} \frac{d^2 f(x)}{dx^2} \Delta x^2 \right] + O(\Delta x^3)$$

$$\frac{df(x)}{dx} \bigg|_k = \frac{f_{k+1} - f_{k-1}}{2 \Delta x} + O(\Delta x^2)$$

CENTERED finite  
difference formula [ACCURACY:  
2<sup>nd</sup>-order]





FORWARD F.D.

$$\left. \frac{df}{dx} \right|_k = \frac{f_{k+1} - f_k}{\Delta x}$$

BACKWARD F.D.

$$\left. \frac{df}{dx} \right|_k = \frac{f_k - f_{k-1}}{\Delta x}$$

CENTERED F.D.

$$\left. \frac{df}{dx} \right|_k = \frac{f_{k+1} - f_{k-1}}{2\Delta x}$$