## Overview

- Error / cond. nr.
- $\oplus P$

Example 2
- Interp
- MC
- Errors / cond.
- F.P.

HW3 due

HW4 timeline

$$\text{Abs} = |x_0 - \tilde{x}| = |\Delta x|$$

with $\Delta x$ over $|x_0 - \tilde{x}|$

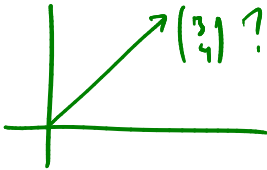$$\text{Rel.} = \frac{\text{Abs.}}{|x_0|} = \frac{|\Delta x|}{|x_0|}$$

# Measuring Error

Why is $|\tilde{x}| - |x_0|$ a **bad** measure of the error?

$$|\tilde{x} - x_0| \longleftarrow \text{Yes.}$$

If $\tilde{x}$ and $x_0$ are vectors, how do we measure the error?

$\begin{pmatrix} 3 \\ 4 \end{pmatrix}$ ?

Need equivalent to the abs. value:

"norm" $\left\| \begin{pmatrix} 3 \\ 4 \end{pmatrix} \right\|_2 = \sqrt[2]{3^2 + 4^2}$

$$\left\| \vec{x_0} - \tilde{\vec{x}} \right\|_2$$

# Sources of Error

What are the main sources of error in numerical computation?

- Truncation error

- Rounding error

$$\text{abs. error} = |x_0 - \tilde{x}| = \text{"trunc. + rounding"}$$

$$\text{abs. error} = 10^{-5} \Big/ \text{both}$$

Scen. 1 $\qquad x_0 = 10^{-6}$

Scen. 2 $\qquad x_0 = 10^0$

$\begin{array}{c} 1.0000\,|\,0 \\ 1.0000\,|\,1 \end{array}$ $\Big\} \; 5 \text{ accu. digit}$

$$\text{rel. error} = \frac{\text{abs. error}}{1} = 10^{-5}$$

# Digits and Rounding

> Establish a relationship between '*accurate digits*' and rounding error.

$$3.14159\ldots$$

$$\rightarrow 3.14$$

rounded to 3 digits

finite precision

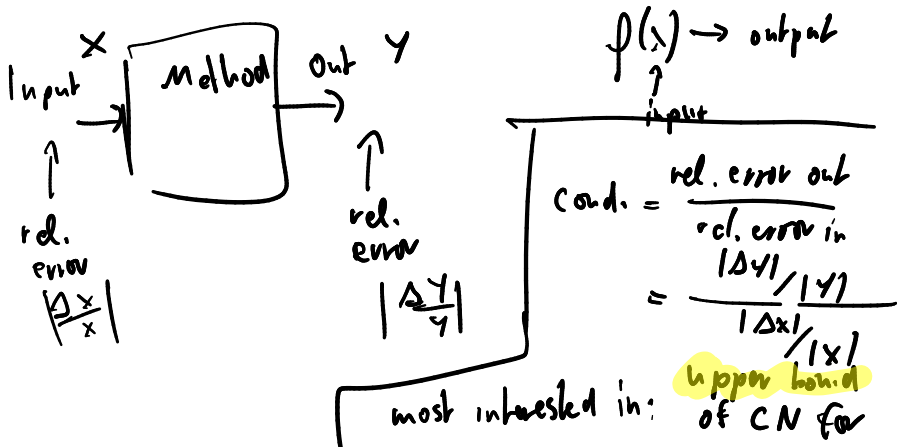rel. error $: \dfrac{|3.14159 - 3.14|}{|3.141591|} \approx 10^{-3}$

# Condition Numbers

> Methods $f$ take input $x$ and produce output $y = f(x)$.
>
> Input has (relative) error $|\Delta x| / |x|$.
>
> Output has (relative) error $|\Delta y| / |y|$.
>
> **Q:** Did the method make the relative error bigger? If so, by how much?



Input $\xrightarrow{x}$ | Method | Out $\xrightarrow{y}$

rel. error $|\frac{\Delta x}{x}|$

rel. error $|\frac{\Delta y}{y}|$

$f(x) \rightarrow$ output

$\uparrow$ input

$$\text{Cond.} = \frac{\text{rel. error out}}{\text{rel. error in}} = \frac{|\Delta y|/|y|}{|\Delta x|/|x|}$$

most interested in: upper bound of CN for

$10^{-3}$ $\xrightarrow{\hspace{2cm}}$ $10^{-1}$,

for this
example:
$CN = 10^2$

all inputs/
outputs

"Good CN": small

"Bad CN": big    (error gets bigger)

$$Abs. \; CN = \left| \frac{Abs. \; error \; in \; out}{Abs. \; error \; in \; in} \right| = \frac{|\Delta Y|}{|\Delta x|}$$

# $n$th-Order Accuracy

Often, *truncation error* is controlled by a parameter $h$.

Examples:

- distance from expansion center in Taylor expansions
- length of the interval in interpolation

A numerical method is called '$n$th-order accurate' if its truncation error $E(h)$ obeys

$$E(h) = O(h^n).$$

# Outline

# Wanted: Real Numbers... in a computer

Computers can represent *integers*, using bits:

$$23 = 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = (10111)_2$$

16    8    4    2    1

How would we represent fractions, e.g. 23.625?

$$23 = 1 \cdot 2^4 + - + 1 \cdot 2^0$$

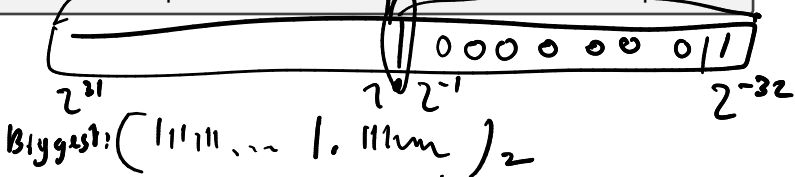$$23.625 = 1 \cdot 2^4 + - + 1 \cdot 2^0 + 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3}$$

$$2^{-3} = .125$$
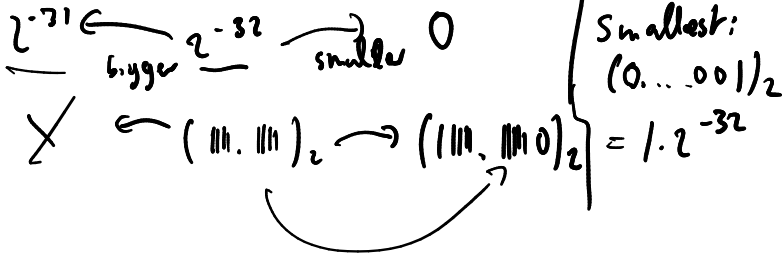
$$23.5$$

$$23.625 = (10111.101)_2$$

# Fixed-Point Numbers

Suppose we use units of 64 bits, with 32 bits for exponents $\geqslant 0$ and 32 bits for exponents $< 0$. What numbers can we represent?

$32$ $32$

$2^{31}$ $2^0$ $2^{-1}$ $2^{-32}$

$0\ 0\ 0\ 0\ \ 0\ 0\ \ 0\ 1$

Biggest: $(111\,111\ldots 1.\,111\text{um})_2$

How many 'digits' of relative accuracy (think relative rounding error) are available for the smallest vs. the largest number?

$2^{-31}$ bigger $2^{-32}$ smaller $0$

$(111.\,111)_2 \longrightarrow (111.\,110)_2$

Smallest:
$(0.\ldots 001)_2 = 1 \cdot 2^{-32}$

exact result: $2^{-32}$

computed result: $2^{-31}$

$\Rightarrow$ rel. error $= \dfrac{|2^{-32} - 2^{-31}|}{2^{-32}} = \dfrac{2^{-32}}{2^{-32}}$

In fixed point:

uneven relative error.

big results: small rel. error

small results: big rel. error

# Floating Point numbers

Convert $13 = (1101)_2$ into floating point representation.

$$13 = \underline{(1.101)_2 \cdot 2^{\underline{3}}}$$
$$= (1101)_2 = (110.1)_2 \cdot 2 = (11.01)_2 \cdot 2^2$$

What pieces do you need to store an FP number?

$1 . 110 \; 11 \quad - \quad 2^{32}$

$(1 . 1 \; 100 \; 11)_2 - 2^{-32}$

**In-class activity:** Floating Point

## Unrepresentable numbers?

Can you think of a somewhat central number that we cannot represent as
$$x = (1.\text{_____})_2 \cdot 2^{-p}?$$