

Lecture 1

Material covered: Up to Chapter 2 in the slides

Symbolic vs Numeric Computing

Symbolic: algebraic manipulation

exact

slow

(takes functions ;
returns functions)

Numerical: Approximate (if we think of real numbers)

fast

(takes numbers ;
returns functions)

Even with symbolic we usually need a number at the end

Packages

Symbolic

SymPy

Mathematica

Sage

Magma

Numeric

Julia

NumPy

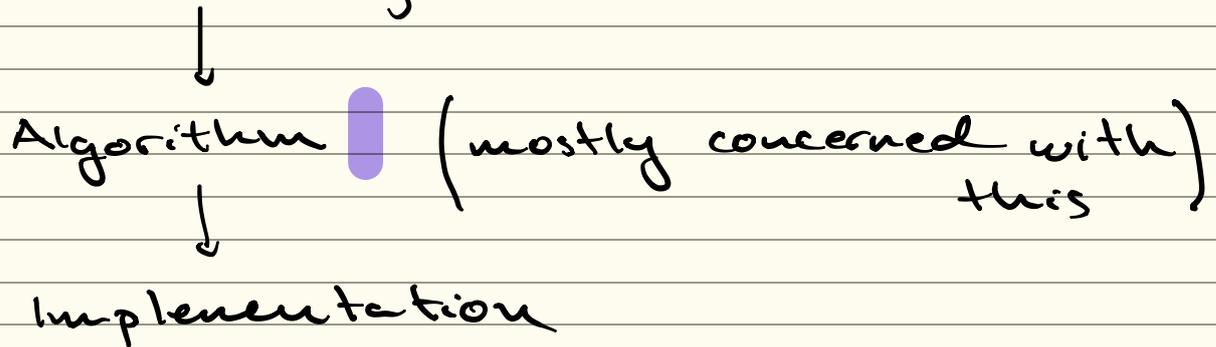
Eigen

BLAS

LAPACK

Overview

* Mathematical formulation



* What is a "good" numerical algorithm?

- stable (small change to input leads to small changes in output)
- robust (not many assumptions)
- fast to compute
- easy to implement

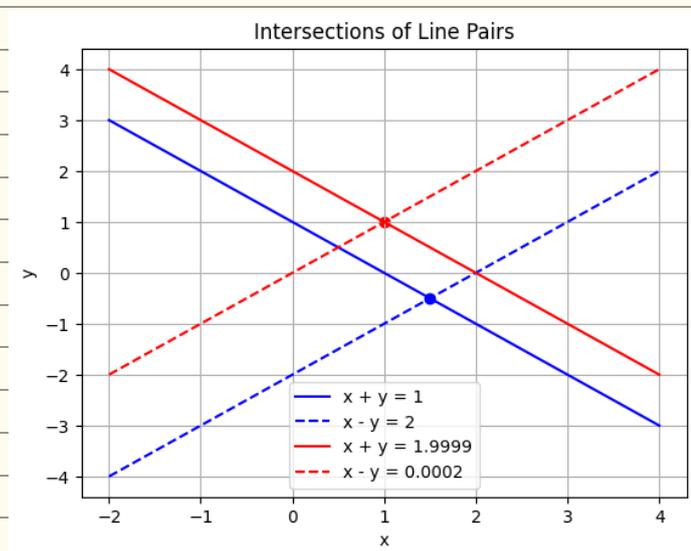
* Sensitive / ill-conditioned problem

- no method can work (so we usually change the problem)

Ex 1 $x + y = 2$
 $x - y = 2$

$x + y = 1.9999$
 $x - y = 0.0002$

Stable



$x = y = 1$

$x = 1.0005$
 $y = 0.99985$

Ex 2

$$x + 0.99y = 1.99$$

$$0.99x + 0.98y = 1.98$$

$$x + 0.99y = 1.9899$$

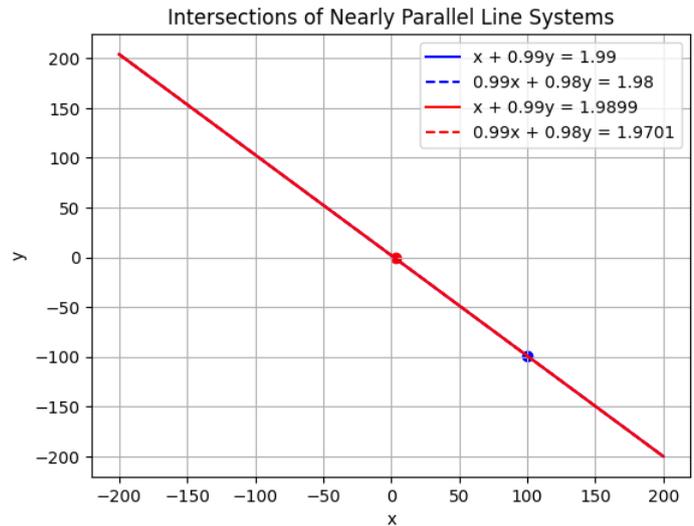
$$0.99x + 0.98y = 1.9701$$

$$x = y = 1$$

$$x = 2.97$$

$$y = -0.99$$

Unstable



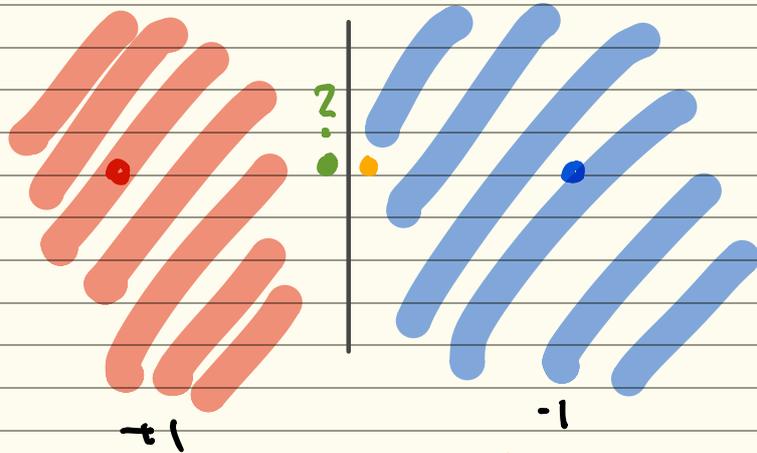
(There are 2 lines)

Ex 3

To which point are we closer to

red $\rightarrow +1$

blue $\rightarrow -1$



very close together $\left\{ \begin{array}{l} \bullet \rightarrow -1 \\ \bullet \rightarrow +1 \end{array} \right.$

Representation of numbers

$$x = \pm 0.d_1d_2\dots d_n \times \beta^e$$

↑
sign bit

β = base

e = exponent

$d_1 \dots d_n$ - mantissa

If $\beta=2$ (binary) $d_1=1$ (almost always)
when is d_1 not 1?

- Overflow vs underflow (see slides)

Alternative representation

$\frac{P}{Q}$ where P & Q are integers

Precision $\pm \frac{1}{2^n}$ (with n bits)
half for P
half for Q

* floating point is more flexible

Rounding Error

Closest number to $x \in \mathbb{R}$

$$f(x) = \begin{cases} x_- & \text{largest representable number} \\ & x \geq x_- \\ x_+ & \text{smallest representable number} \\ & x \leq x_+ \end{cases}$$

$$\rightarrow x_- \leq x \leq x_+$$

Error

Assume $f(x) = x_-$

$$x = (0.1b_2 \dots b_{m+1})_2 \times 2^e$$

$$x_- = (0.1b_2 \dots b_m)_2 \times 2^e$$

$$x_+ = x_- + 2^{-m} = ((0.1b_2 \dots b_m)_2 + 2^{-m}) \times 2^e$$

Absolute error:

$$\max(x - x_-) = \frac{x_+ - x_-}{2} = 2^e \left(\frac{2^{-m}}{2} \right) = 2^{e-m-1}$$

Relative error

$$\frac{x - x_-}{x} \leq \frac{2^{e-m-1}}{\underbrace{2^{e-1}}_x} \leq 2^{-m}$$

x is certainly
bigger than
this, if $b_{m+1} = 0$
($2^e / 2$)

Assume

$$u = 2^{-m}$$

$$x_- = x_- + x - x \geq -ux + x = x(1-u)$$

$$\Rightarrow x_+ \leq x(1+u)$$

$$f(x) \leq x(1+\delta) \quad |\delta| \leq u$$

Computing with representations

Let \odot be an operation

$$\odot \in \{+, -, \cdot, : \}$$

Operations in representations are not associative!

Example: addition

$$(a+b)+c \neq a+(b+c)$$

Cause: Adding two represented numbers does not result in a represented number

Claim: Add small numbers first
(then large ones)

$$\begin{aligned}(a+b)+c &= fl(fl(a+b)+c) \\ &= fl((a+b)(1+\delta_1)+c) \\ &= ((a+b)(1+\delta_1)+c)(1+\delta_2) \\ &= ((a+b+c) + (a+b)\delta_1)(1+\delta_2) \\ &= (a+b+c) \left(1 + \frac{a+b}{a+b+c} \delta_1\right) (1+\delta_2) \\ &= (a+b+c) \left(1 + \frac{a+b}{a+b+c} \delta_1 + \delta_2 + \frac{a+b}{a+b+c} \delta_1 \delta_2\right) \\ &\approx (a+b+c) (1 + \frac{a+b}{a+b+c} \delta_1 + \delta_2) \\ &\quad \delta_1 \delta_2 \ll \delta_1, \delta_2\end{aligned}$$

Likewise

$$a+(b+c) = (a+b+c) \left(1 + \frac{b+c}{a+b+c} \delta_3 + \delta_4\right)$$

if $a < b < c \Rightarrow \boxed{\frac{a+b}{a+b+c} < \frac{b+c}{a+b+c}}$

Addition & Subtraction Error

Subtracting similar numbers

$$a = 0.\overbrace{x \dots x}^m \times 1 \overset{\text{loss to precision}}{\text{ssss}} \cdot 10^e$$

$$b = 0.y \dots y 0 \text{ tttt} \cdot 10^e$$

→ difference may be lost

$$a-b = 0.0 \dots 1 \quad \text{error: } 2^{-m+e}$$

Addition is ok but multiple addition errors accumulate (so many steps → problematic)

Example

Solving $x^2 + 2ax + b = 0$ $a > 0$ $b < a^2$

$$x_2 = \frac{-2a + \sqrt{4a^2 - 4b}}{2} = -a + \sqrt{a^2 - b}$$

Steps $s_1 := a^2$ if $a^2 \gg b$

$$s_2 := s_1 - b$$

$$s_3 := \sqrt{s_2}$$

$$s_4 := s_3 - a$$

s_4 can have a large error

(Possible)

Solution:

$$x_2 = -a + \sqrt{a^2 - b} \cdot \underbrace{\frac{a + \sqrt{a^2 - b}}{a + \sqrt{a^2 - b}}}_1$$
$$= \frac{-b}{a + \sqrt{a^2 - b}}$$

More steps \Rightarrow avoid subtraction of similar terms

A common theme

$$f(x) = x (\sqrt{x+1} - \sqrt{x})$$

Back: if $\sqrt{x+1} \approx \sqrt{x}$

$$f(x) = f(x) \frac{\sqrt{x+1} + \sqrt{x}}{\sqrt{x+1} + \sqrt{x}} = \frac{x}{\sqrt{x+1} + \sqrt{x}}$$

Summation

$$\frac{1}{1 \cdot 2} + \frac{1}{2 \cdot 3} + \dots + \frac{1}{n(n+1)}$$

has a smaller error

than

$$\frac{1}{n(n+1)} + \dots + \frac{1}{2 \cdot 3} + \frac{1}{1 \cdot 2}$$

Integration

$$I_n = \int_0^1 x^n e^{-x} dx$$

$$\textcircled{1} I_n = -\frac{1}{e} + n I_{n-1} \quad I_0 = 1 - \frac{1}{e}$$

$$\textcircled{2} \text{ Backward: } I_{n-1} = \frac{1}{n} I_n + \frac{1}{ne}$$

$\textcircled{2}$ is better than $\textcircled{1}$ \Rightarrow chose any value for

I_n (n large) \Rightarrow error is divided each step by

n (becomes small very quickly)

General Analysis - addition/subtraction

$$x, y \in \mathbb{R} \quad \bar{p} \approx p = x + y$$

$$\begin{aligned}\bar{p} &= f(|f(x) + f(y)|) = f(|x(1+\delta_1) + y(1+\delta_2)|) \\ &= (x(1+\delta_1) + y(1+\delta_2))(1+\delta_3) \\ &= x(1+\delta_1)(1+\delta_3) + y(1+\delta_2)(1+\delta_3) \\ &= x + y + x(\delta_1 + \delta_3 + \delta_1\delta_3) + y(\delta_2 + \delta_3 + \delta_2\delta_3)\end{aligned}$$

Relative error

$$\frac{|\bar{p} - p|}{|p|} = \frac{|x(\delta_1 + \delta_3 + \delta_1\delta_3) + y(\delta_2 + \delta_3 + \delta_2\delta_3)|}{|x + y|}$$

$$\text{if } |x + y| \approx 0 \quad \frac{|\bar{p} - p|}{|p|} \gg 0$$

Depends on size of $|x + y|$

Multiplication & Division

$$p = x \cdot y$$

$$\bar{p} = f(|f(x) \cdot f(y)|) = f(|x(1+\delta_1) \cdot y(1+\delta_2)|)$$

$$= x \cdot y (1+\delta_1)(1+\delta_2)(1+\delta_3)$$

$$\approx x \cdot y (1 + \delta_1 + \delta_2 + \delta_3 + \dots)$$

produkti vec δ_i

$$|\delta_i| \leq \epsilon$$

Relative Error

$$\frac{\bar{p} - p}{p} = \frac{|x \cdot y| (\delta_1 + \delta_2 + \delta_3 + O(\epsilon^2))}{|x \cdot y|}$$

$$= (\delta_1 + \delta_2 + \delta_3 + O(\epsilon^2))$$

Product independent of size of $x \cdot y$