

# Numerical Analysis ?

← Subject of writing numerical approx to problem  
 Idea: Numerical is the heart of studying math which is close to reality but not the same as the desired quantity  
 (approximation)

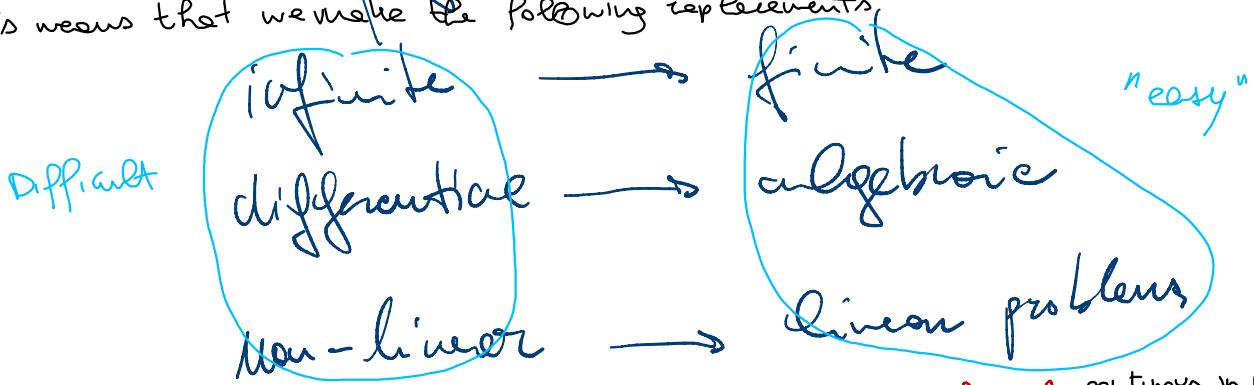
- Algorithms
- Analysis (errors)
- Implementation (python)

We can only approximate **Well posed problems**. i.e.

- ||| 1) there exist a unique solution  
 2) it depends continuously on the data |||

idea behind well posed is that we simplify in a controlled manner

This means that we make the following replacements:



— Example :

Given a function "f", a point  $x_0$ ,  
 can get "f'(x\_0)" (derivative of f in  $x_0$ )

Is it well posed? Yes, because exist a UNIQUE JT ← derivative of f at  $x_0$   
 If f is not derivable  $\Rightarrow$  problem is not well posed  $\Rightarrow$  WE MUST CHECK DOMAIN!

→ use finite difference —

! I'm not claiming to know anything about f except the evaluation of its body

$$f'(x_0)$$

$$\approx$$

$$\frac{f(x_0 + h) - f(x_0)}{h}$$

anything you choose here will be affected by what you can do and what you can't

with small h  
 $\downarrow$

(s.t.  $x_0 + h \in (a, b)$ )

Error : Taylor expansion around  $x_0$ :

Typical ways to estimate error is

$$f(x_0 + h) = f(x_0) + h \cdot f'(x_0) + \frac{h^2}{2} f''(x_0) + \dots$$

When you compute powers of  $2$  ( $2^{**i}, i \dots$ ),  
The accuracy of computer is perfect upto  
 $\sim 10^{-4}$ , and in reality is not perfect upto  
 $10^{14} - 10^{-15}$

In theory, as  $h \rightarrow 0$ , the finite approx should  
get closer to the exact result (exact derivative)

In practice much else happens

$$f(x_0 + h) \approx \sum_{i=0}^n \frac{f^{(i)}(x_0) h^i}{i!}$$

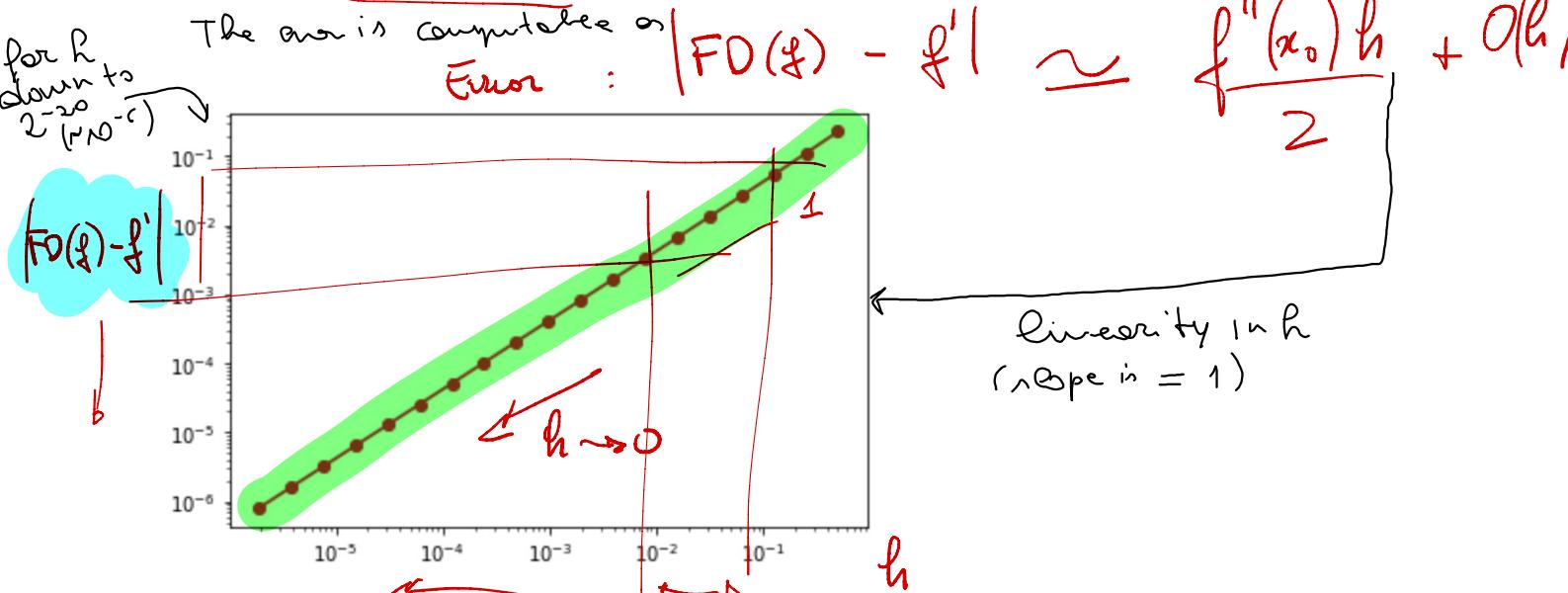
$$FD := \frac{f(x_0 + h) - f(x_0)}{h}$$

for  $n \rightarrow \infty$ ,  $\approx \equiv$   
 provided that function  
 is nicely behaved, so that  
 I can do it for  $n$  functions  
 at each step

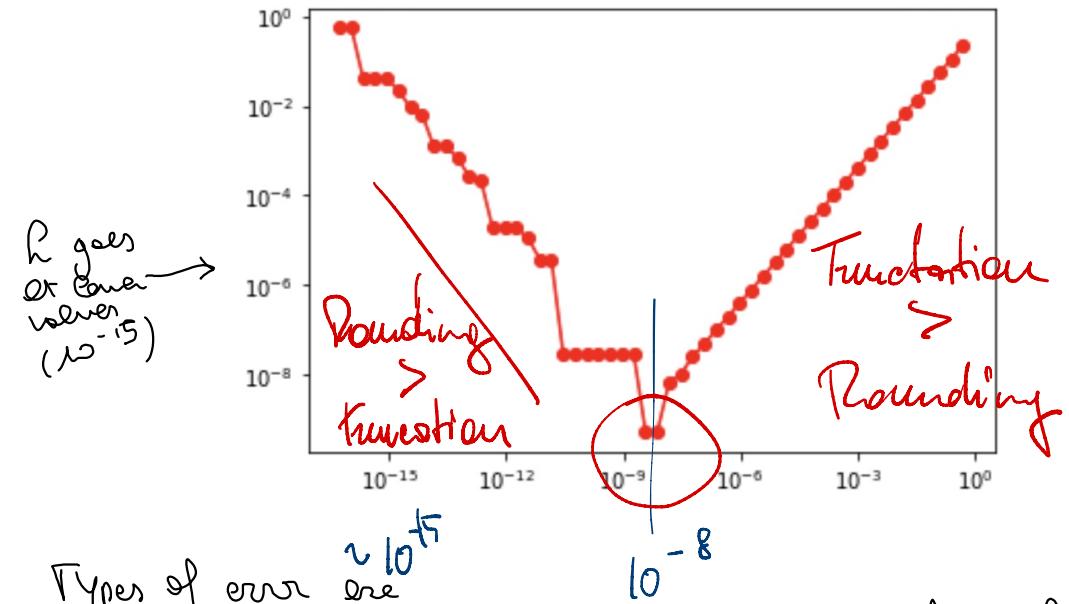
$$f(x_0) + \frac{1}{2} f''(x_0) h + O(h^2)$$

exact derivative  $=$

+ much smaller than  $h$  means that  
 $\lim_{h \rightarrow 0}$  is much smaller compared  
 to  $h$



interesting things happens when  $h$  keeps getting smaller, but the algorithm above doesn't tell us that: it tells us that the error, being  $\sim h$ , will go to zero in a linear fashion as  $h$  goes down



$\sqrt{\epsilon} \sim \text{optimal } h$

$\epsilon$ : precision of machine

(Rounding precision)

If you have error in measurements of data, i.e. the  $x_0$  above, and depends on used algorithms

- Data error (Data) → always present
- Computational Error (Algorithm)
- Truncation error (Algorithm error in Exact Arithmetic) → not the one of PC
- Rounding Error (Finite Arithmetic) → error you make at the algorithm level by rounding off by multiplying exact of the ans

If doesn't take account of the fact that machines operate in a different arithmetic

arithmetic in pc is of power of 2

diff between what you expect from algorithm by putting exact st, and what you get when you run your algorithm (the approximation)

### FINE DIFF:

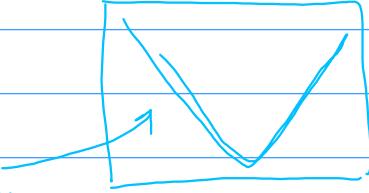
Infinite series - Truncated series

linear  
the graph above is  
**TRUNCATION**  
**Error ONLY**

because trunc. error dominates over rounding error

but, as h drops, rounding error starts becoming more relevant

here it dominates over the trunc error



Pb: Recompute value of a function "F"

$$F: \mathbb{R} \rightarrow \mathbb{R}$$

at a argument  $x$

So we have

•  $x$  : true value of input

what to compute

•  $F(x)$  : desired value of output

affect by rounding,  
data error, ...

Error on  
input data

•  $\hat{x}$  : approximate input (may be "off" for rounding)

Error on  
outcome  
due to  
approximation

•  $\hat{F}(\hat{x})$  : computed output (may be "off" for both  
rounding and truncation)

$$\underbrace{F(\hat{x}) - F(x)}_{\text{effect of rounding on input data}} =$$

ROUNDING OR COMPUTATIONAL  
ERROR DUE TO NUMBER OF DIGITS  
IN INPUT (NUMBER OF DIGITS  
IN APPROXIMATED PROBLEM)  
Truncation error is  
exist on floating point

$$= \underbrace{\hat{F}(\hat{x}) - \hat{F}(x)}_{\text{effect of rounding on algorithm}} + \underbrace{\hat{F}(x) - F(x)}_{\text{effect of truncation}}$$

①

②

or  
Sensitivity of algorithm with  
respect to the original problem

$$= \underbrace{\hat{F}(\hat{x}) - \hat{F}(\hat{x})}_{\text{You are changing the  
algorithm, and  
only more other  
thing is true, but  
you can't remove  
it doesn't slope}} + \underbrace{\hat{F}(\hat{x}) - F(x)}_{\text{Computational error  
(diff between exact and  
expected output)}}$$

③

Sensitivity of original problem  
with respect to perturbations in  
the data

ROUNDING ERROR  
IN EXACT MODELS

Rounding precision :  $\epsilon$   
(machine precision)

the largest number s.t.

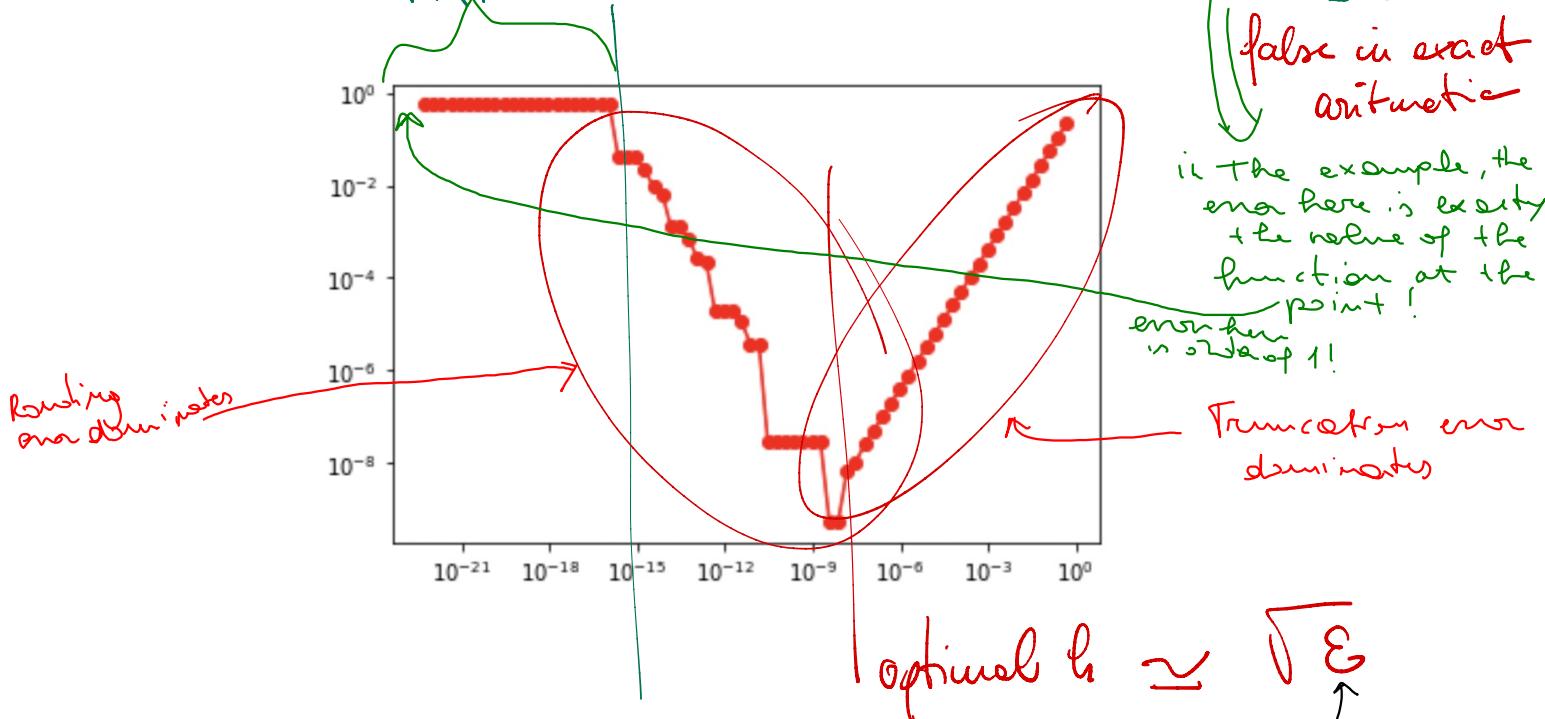
$$\text{fl}(1+\epsilon) = \text{fl}(1)$$

gives optimal  
precision

$\mu < \epsilon$  including  
error dominates

floating point representation of a number

$$f(f(x+h)) = f(f(x)) \xrightarrow{\text{if } f'(x) = 0} f(x+h) - f(x) = 0$$



$$\sin(1) = 0.8414709848078965$$

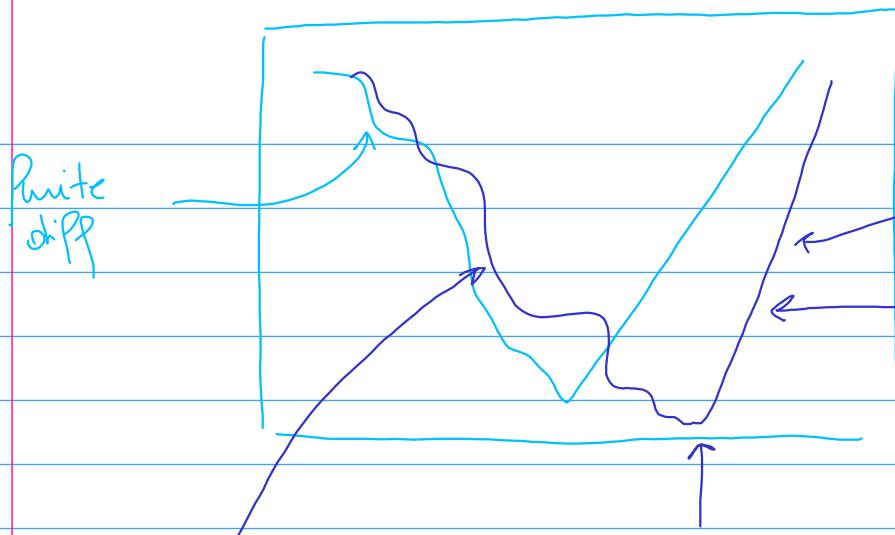
for G4 S.2  
doubles  
(python)  
 $\epsilon \sim 10^{-15}$

$$\frac{f(x+h) - f(x)}{h}$$

This  $h$  must be chosen correctly to prevent overflow or underflow, otherwise you'll get good results for some  $f$  and bad for others.

All of this also depends on where we evaluate  $f$  ( $x_n$ ) and function eval at the same time

! Don't use abs error as a measure of error, but use the RELATIVE Error



Not  
avoidable,  
this gets  
better &  
more linearly

optimal  
 $h$  in row  
 $\sim \varepsilon^{1/3}$

central finite  
difference.  
BETTER, because  
here we have  
another power  
of  $n$ , which  
is power of 2  
(slope  $\ell n n$  is 2<sup>-1</sup>)

# Numerical Analysis - sec. 2

problem can be written  
in a functional approach

8.10.2020

$$y = F(x)$$

where

$$F: \mathcal{X} \longrightarrow \mathcal{Y}$$

$$\text{data} \rightarrow x \longrightarrow y = F(x) \quad \begin{matrix} \text{F is a given function} \\ \mathcal{X} \text{ output/result} \end{matrix}$$

## Functional Evaluation

EXAMPLES

$$1) \mathcal{X} : \mathbb{R}, \mathcal{Y} : \mathbb{R}$$

F: continuous func.

$$y = F(x)$$

### 1) Sum of two numbers:

$$\mathcal{X} : \mathbb{R}^2, \mathcal{Y} : \mathbb{R} \quad : F: \text{sum of numbers.}$$

$\downarrow$  pair of numbers

$$x = (a, b)$$

$$F(x) = F((a, b)) = a+b$$

### 3) Computation of derivative

$$\mathcal{X} : C^1([a, b]) \quad \mathcal{Y} : \mathbb{R}$$

F: evaluation of first derivative in  $x_0 \in [a, b]$

$$y = F(x) := \underbrace{x'}_{\text{function}}(x_0)$$

Only well posed problems are characterized by 2 properties

$$1) \forall x \in \mathcal{X}, \exists! y \in \mathcal{Y} \text{ s.t. } F(x) = y$$

$$2) \exists k \text{ (condition number)} \text{ s.t. } \forall x, \hat{x} \in \mathcal{X}$$

$$\|F(x) - F(\hat{x})\| \leq k \|x - \hat{x}\| \quad \begin{matrix} \text{basically} \\ F \text{ is} \\ \text{continuous} \\ (\text{bounded} \\ \text{for every}) \end{matrix}$$

$\overbrace{\text{each of these}}^{\text{are unique}}$

$\mathcal{Y}$

$\mathcal{X}$

We assume that

## Real Vector Spaces (Normed)

Def (real vector spaces)

RVS: collection of objects for which it makes sense  
"V" to "sum" and "scale"

$$+ : V \times V \longrightarrow V$$

$$(u, v) \longrightarrow w = u + v$$

$$\cdot : V \times \underline{\mathbb{R}} \longrightarrow V$$

such that

$$\forall u, v \in V, \forall \alpha, \beta \in \mathbb{R}$$

$$\alpha u + \beta v = w \in V$$

$\mathbb{R}^2$

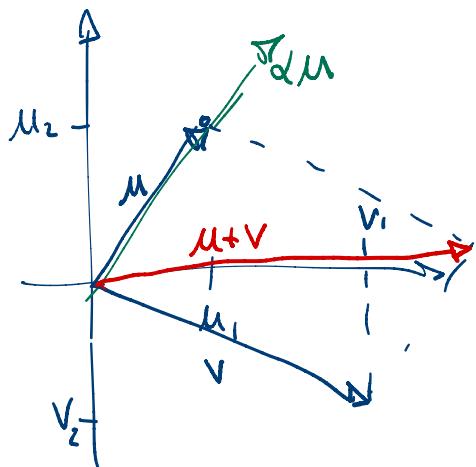
$C^0([a, b])$

$$u = (u_1, u_2)$$

$$v = (v_1, v_2)$$

$$u+v = (u_1+v_1, u_2+v_2)$$

$$\alpha u = (\alpha u_1, \alpha u_2)$$



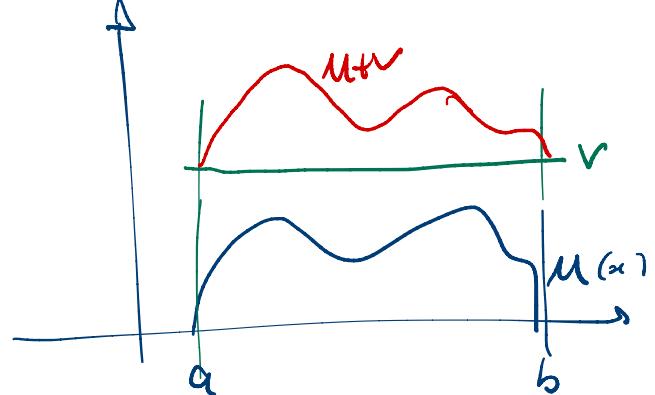
$$u : [a, b] \longrightarrow \mathbb{R}$$

$$v : [a, b] \longrightarrow \mathbb{R}$$

$$w = u+v : [a, b] \longrightarrow \mathbb{R}$$

because  $u, v$   
are defined  
in  $[a, b]$

$$\alpha \longrightarrow u(\alpha) + v(\alpha)$$



Way to measure things

Introduce  $\downarrow$  NORMS

# Norms on Real Vector Spaces

a function  $\|\cdot\| : V \rightarrow \mathbb{R}_0^+$  positive or zero

$$1) \|u\| > 0 \quad \forall u \in V$$

$$2) \|u+v\| \leq \|u\| + \|v\| \quad \text{Triangle Inequality}$$

$$3) \|\alpha u\| = |\alpha| \|u\|$$

$$\text{optional } 4) \|u\| = 0 \iff u = 0$$

without ④ it is a semi-norm

We will use

$L_p$  norms on  $\mathbb{R}^n$

$L_p$  norms on  $S \subset \mathbb{R}^d$

$\|\cdot\|_*$  operational norms induced by  $\|\cdot\|_V$

Defn:  $\mathbb{R}^n$ :  $L_p$  norm  $\|u\|_p := \left( \sum_{i=1}^n |u_i|^p \right)^{\frac{1}{p}}$   $\downarrow p \rightarrow \infty$

$\|u\|_\infty := \max_i |u_i|$   $\ell_\infty$  norm

$S \subset \mathbb{R}^d$

$L_p$  norms  $\|u\|_p := \left( \int_S |u|^p \right)^{\frac{1}{p}}$   $\downarrow p \rightarrow \infty$  L<sub>p</sub>-norm

$\|u\|_\infty := (\text{ess}) \sup_{x \in S} |u(x)|$

$\|\cdot\|_*$

induced by vector space norm  $\|\cdot\|_V$  on  $V$

$\|A\|_* := \sup_{0 \neq x \in V} \frac{\|A(x)\|_W}{\|x\|_V}$

You def the norms of domain and codomain of  $A$ , and then you can def the operat norm

↑ used to measure how functional, functions work, or what is their measure in a space

$\ell_p$  norm of Matrices in  $\mathbb{R}^{n \times m}$

which is an obj of this kind

$$A : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$A^{n \times m} \xrightarrow{x \in \mathbb{R}^m} x \in \mathbb{R}^n$$

$$A \cdot x = y \in \mathbb{R}^n$$

$\ell_p$  norm of the matrix  
functional  
norm because it measures  
the norm of a function

$$A \in \mathbb{R}^{h \times m}$$

$$\|A\|_p := \sup_{0 \neq x \in \mathbb{R}^n} \frac{\|Ax\|_p}{\|x\|_p} = \|A\| \text{ induced by } \|\cdot\|_p \text{ on } \mathbb{R}^n$$

$$A(x) = y$$

$$\sum_{j=1}^m A_{ij} x_j = y_i$$

$$A \in L(\mathbb{R}^m, \mathbb{R}^n)$$

$$\|\cdot\|_* \text{ norm on } L(\mathbb{R}^m, \mathbb{R}^n) \equiv \mathbb{R}^{n \times m}$$

Why we need them? Every quantity related to a well posed problem can be put in a way between two spaces and their norms.

### Well posed problems

$$1) \quad \forall x \in \mathcal{X}, \exists! y \in \mathcal{Y} \text{ s.t. } F(x) = y$$

Normed Vector Space

Normed Vector Space

$$2) \quad \exists k \text{ (condition number)} \text{ s.t. } \forall x, \hat{x} \in \mathcal{X}$$

$$\|F(x) - F(\hat{x})\| \leq (k_{\text{Abs}}) \|x - \hat{x}\|$$

Here we have norms

$$3) \quad \text{Relative cond. number (only if } x \neq 0, F(x) \neq 0)$$

$$\text{Krel st. } \forall x, \hat{x}$$

$\uparrow$   
INDEPENDENT ON X

$$\frac{\|F(x) - F(\hat{x})\|_Y}{\|F(x)\|_Y} \leq \kappa_{rel} \frac{\|x - \hat{x}\|_X}{\|x\|_X}$$

Prop

A problem is well posed if  $\kappa_{rel/abs}$  is

"small"  $\leftrightarrow$  ~~Krel/abs must be controllable~~

with respect to what  $\rightarrow$  relative depends on problem handled

Example : sum of two numbers : ( $L_1$  norm)

$$X: \mathbb{R}^2 \quad Y: \mathbb{R}$$

$$x \in \mathbb{R}^2 \rightarrow \|x\|_1 := |x_1| + |x_2| \quad \|y\|_1 = |y|$$

We have that

$$\frac{\|F(x) - F(\hat{x})\|_1}{\|F(x)\|_1} = \frac{|x_1 - \hat{x}_1 + x_2 - \hat{x}_2|}{|x_1 + x_2|} \leftarrow \| \cdot \|_{Y=1}$$

$$\frac{\|x - \hat{x}\|_1}{\|x\|_1} = \frac{|x_1 - \hat{x}_1| + |x_2 - \hat{x}_2|}{|x_1| + |x_2|}$$

$$\Delta y, \quad \Delta x \quad \Delta x = x - \hat{x} \quad \Delta y = F(x) - F(\hat{x})$$

We now find  $k_{abs}$ :  $\exists k_{abs} \mid \forall x_1, x_2, |\Delta y| \leq k \|\Delta x\| ?$

$$|\Delta y| \leq k \|\Delta x\| \leq k(|\Delta x_1| + |\Delta x_2|)$$

Assume that  $\Delta x \neq 0$

$$\frac{|\Delta g|}{|\Delta x|} \leq \kappa$$

$$\frac{|\Delta x_1 + \Delta x_2|}{|\Delta x_1| + |\Delta x_2|} \leq \frac{|\Delta x_1| + |\Delta x_2|}{|\Delta x_1| + |\Delta x_2|} \leq 1$$

From pt 5 of absolute  $\kappa$ , this is perfect: you perturb the dots, and the perturbation of the result follows perfectly the pt b of the inputs ( $\kappa_{abs} = 1$ )

$$\kappa_{abs} = 1$$

but ~~not~~ local st. Indeed:

$$\frac{|\Delta y|}{|\Delta x|}, \frac{|x|}{|y|} \leq \frac{|x|}{|y|} \leq \kappa_{rel}$$

$\kappa_{abs} \leq 1$

we cannot control this number!

because

$\Rightarrow$  if you try to sum 2 numbers with closing values and opposite signs, you're dealing with an ill-posed problem, but the sum of 2 numbers is not guaranteed to be well-conditioned.

Example:

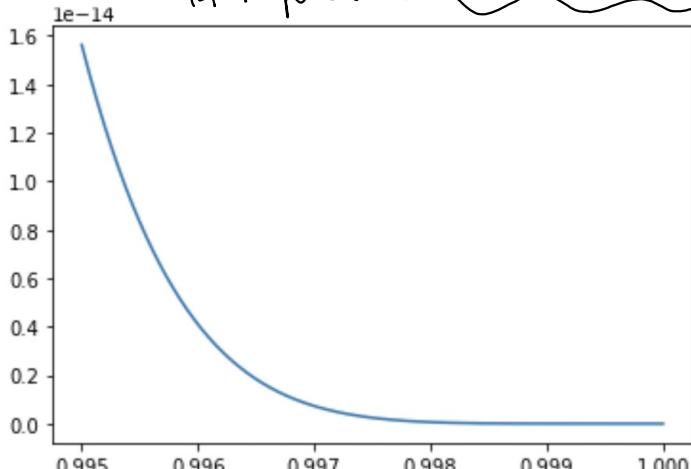
$$\frac{|x|}{|y|} = \frac{|x_1| + |x_2|}{|x_1 + x_2|}$$

can be as large as we want

$$(1-x) \text{ in } [1-\epsilon, 1]$$

$\epsilon$  is small.

if I plot

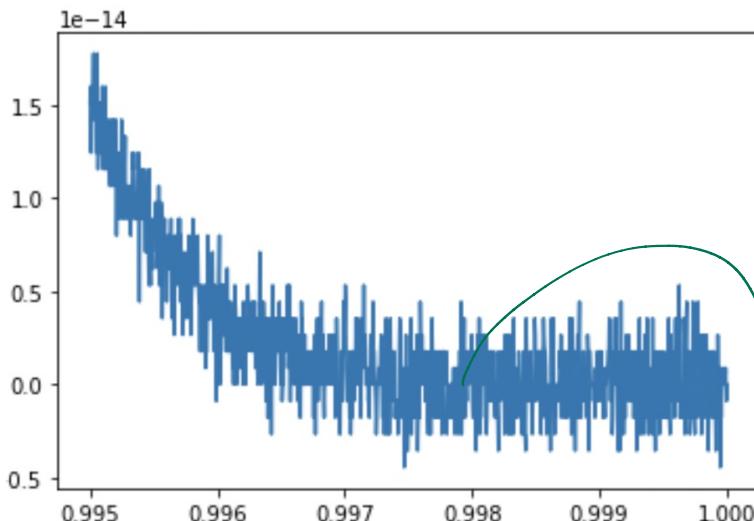


```
x = linspace(.995, 1, 1029)
```

```
y = (x-1)**6
```

but if I round the power I get a bad result

$$y = x^{**6} - 6*x^{**5} + 15*x^{**4} - 20*x^{**3} + 15*x^{**2} - 6*x + 1$$



due to addition of very close numbers (in value) but with  $\neq$  signs  $\rightarrow$  we say that this is an ill posed problem  
 $x - x = \text{floating point errors}$   
(rounding errors)

Consider:

If  $x_1 = -(x_2 + \varepsilon)$

so what we obtain from machine epsilon

$$\frac{|x_1|}{|y|} = \frac{|x_2 + \varepsilon| + |x_2|}{|-x_2 - \varepsilon| + |x_2|} \leq \frac{2|x_2| + |\varepsilon|}{|\varepsilon|}$$

if  $x_2$  not close to zero then

$$\frac{|x_1|}{|y|} \leq \frac{2|x_2|}{|\varepsilon|} + 1$$

so, relating this to the Krel quantity  $\frac{|\Delta y|}{|y|}$  if  $|\varepsilon|$  is too small we get a problem --

$$\frac{|\Delta y|}{|y|} \cdot \frac{|x|}{|\Delta x|} \leq K_{\text{rel}} \Rightarrow$$

which means  $K_{\text{rel}} = \infty$  to agree with its def

$$\frac{|\Delta y|}{|\Delta x|} \leq K_{\text{rel}}$$

→ It has to be greater for every chosen  $\varepsilon$ )

In this one, it may also depend on  $x$

Sum of 2 numbers is not a "RELATIVELY"  
well posed problem ( $\leftrightarrow$  in relative terms  
 $\leftrightarrow$  k rel) relative change in result is caused by  
relative change in input, what you get is  
not controllable

But it is an "ABSOLUTELY" well posed problem

THIS DEPENDS ON  
THE EVALUATION POINT

$\uparrow$  difference!  
 $\downarrow$   
if my pt will be 0  
for non zero  $x \leftarrow$  difference!  
 $\Rightarrow$  k rel  $\rightarrow \infty$   
 $\Rightarrow$  ill posed problem

! Rounding errors are  
MAGNIFIED by  
big old train  
numbers

1) 2! solution

2) pb. is stable  $\Leftrightarrow$  well posed  $\Leftrightarrow$  finite cond. Num.

## $X, Y$ Normed Vector Spaces

$$F: X \rightarrow Y$$

]! solution  $\Rightarrow F$  is an INVERTIBLE FUNCTION  
positive and possibly infinite  
used whenever is required to study stability of a problem and implies that there is no more redundancy in data

We have

1)  $k_{abs} \in \mathbb{R}^+ \cup \{+\infty\}$

s.t.  $x, \hat{x} \in X$  LIPSCHITZ PROPERTY,  $\|F(x) - F(\hat{x})\|_Y \leq k_{abs} \|x - \hat{x}\|_X$

which is  $\Leftrightarrow$  to  $k_{abs} \geq \frac{\|F(x) - F(\hat{x})\|_Y}{\|x - \hat{x}\|_X}$  choose the smallest number greater than or equal to  $k_{abs}$

If  $k_{abs} \neq \infty \Rightarrow$  problem is stable (well conditioned)

1)  $k_{rel} \in \mathbb{R}^+ \cup \{+\infty\}$

s.t.  $x, \hat{x} \in X$

Similarly as before, the smallest number greater than or equal to

$k_{rel} : \sup_{x, \hat{x} \in X}$

$$\frac{\|F(x) - F(\hat{x})\|_Y}{\|F(x)\|_Y} \leq k_{rel} \frac{\|x - \hat{x}\|_X}{\|x\|_X}$$

denominates as the relativity

I split the norm because I have inequalities but then I can obviously the definition and put  $k_{rel}$  even to this

$$\frac{\|F(x) - F(\hat{x})\|_Y}{\|x - \hat{x}\|_X} \frac{\|x\|_X}{\|F(x)\|_Y}$$

If  $F$  is not linear  $\Rightarrow$  has in necessarily the Lipschitz constant of  $F$

$$k_{abs} \sup_{x, \hat{x} \in X} \|F^{-1}(F(x))\|_X$$

use def. of operator norm applied to  $F^{-1}$ : can do this because  $F$  is invertible

$$\frac{\|x\|_X}{\|F(x)\|_Y}$$

If operator is linear

Example:

$$Ax = y \Rightarrow y = F(x) \quad y = F(x)$$

$$\delta x = x - \hat{x} \quad \delta y = y - \hat{y}$$

Absolute stability:

$$A\hat{x} = \hat{y} \Rightarrow \hat{y} = F(\hat{x})$$

$$A\delta x = \delta y$$

$$\|Ax - A\hat{x}\| = \|A(x - \hat{x})\| = \|\delta y\|$$

I can collect the  $A$  thanks to the linearity of the problem

Exhibit property of operator norm of a matrix (follows from its definition)

$$\|Sg\| \leq \|A\| \|Sx\|$$

end we get  $K_{abs}$

Finding  $K_{abs}$  gives (by applying the def of  $K_{abs}$  above)

$$\frac{\|Sg\|}{\|y\|} \leq \|A\| \|A^{-1}\| \frac{\|Sx\|}{\|x\|}$$

Condition number of A

(this in def: condition number of a matrix is  $\|matrix\| \|matrix^{-1}\|$ )

it will depend on the norm you use

Prop

A problem is (absolutely) relatively well posed if  
( $K_{abs}$ )  $K_{abs} < +\infty$ .

What is an approximation?

Sequence of problems  $F_n$  ( $n \in \mathbb{N}$ )

applied to a sequence of "approximated data"  $x_n$   
(implied approximated sequence of domains  $X_n$ )

s.t.

$$\lim_{n \rightarrow \infty} |F_n(x_n) - F(x)| = 0 \quad \leftarrow$$

CONVERGENCE OF THE  
PROBLEM: this is really  
what we want to check

Note 1:

$$X_n \rightarrow X$$

Simple case:  $X_n \subset X$

2) Note 2:  $F_n \rightarrow F$

$\equiv$  thin in the algorithm

consistency  $\leftrightarrow$  Approximated  
method converges to  
the exact  
problem  
that we have

Assume that  $X_n = X$

$F_n \rightsquigarrow$  I can apply the  
algorithm to all pt  
of the "continuous  
space"

2)  $\xrightarrow{\text{can be rewritten}} \lim_{n \rightarrow \infty} |F_n(x) - F(x)| \rightarrow 0$

I'm asking that the algo  
applied to exact input  
minus the exact problem  
(in modulus) goes to zero  
for increasing  $n$

LAX-RICHTMEYER THEOREM:

For consistent pb  $(\lim_{n \rightarrow \infty} |F_n(x) - F(x)| = 0)$

# CONVERGENCE

$\Leftrightarrow$   
if  
 $F$

# STABILITY

$y \leftarrow$  output

Proof  $x^*$   $\xrightarrow{\text{input}}$

If  $F$  is well posed  $\Rightarrow F'(y_0) = \hat{x}$

$$F'(\hat{y}_0) = \hat{x} \xrightarrow{F} F_n \xrightarrow{\text{evaluate it with algorithm}} \cdot y_n$$

can be expressed of  
 $x_{n+1}$  or  $x_n + \Delta x$

if you want to prove convergence you must show that  
Convergence  $\lim_{n \rightarrow \infty} |F_n(x_n) - F(x)| \rightarrow 0$

We have:

$$\left| F_n(x_n) - F(x_n) + F(x_n) - F(\hat{x}) + F(\hat{x}) - F(x) \right| = \dots \rightarrow 0$$

adding and subtracting

consistency  $\rightarrow 0$       stability  $\leq K|x_n - \hat{x}_n|$       stability  $\leq K|\hat{x} - x|$

All the errors falls into 3 categories

If we control this, convergence is implied by stability of the problem.

Stability errors (of continuous pb)  
 " " " (of discrete pb)

consistency error

If continuous problem is stable, and the discrete problem is stable and consistent  $\Rightarrow$  we have convergence

Example:

Interpolation

(Polynomial interpolation)

We consider the use of

One dimension.

and we study the following problem

Pb: Find a Polynomial of order  $n$  that coincides with the input function at (given)  $n+1$  points

CONSIDERATIONS

1. in  $\mathcal{X}$  we use the  $L^\infty$  norm:

We want to understand what  $F$  is and what's the condition number of the problem

$$\|u\|_\infty := \max_{x \in [a,b]} |u(x)|$$

They are not inputs!  
They are part of the problem itself

$\mathbb{P}^n([a,b])$  is finite dimensional, while  $\mathcal{Y}$  is not because the outcome is a polynomial space of dimension  $n+1$ .

2. The space  $\mathcal{Y}$  is  $\mathbb{P}^n([a,b])$ : we use again  $\|\cdot\|_\infty$  as a norm.

Note 1:  $\mathcal{Y}$  is finite dimensional, and we can write linear combination

$$\text{it as } \mathcal{Y} = \mathbb{P}^n([a,b]) = \text{Span}\{v_i\}_{i=0}^n \text{ s.t. } p \in \mathcal{Y} \Leftrightarrow \exists! \{p^i\}_{i=0}^n \in \mathbb{R}^{n+1} \text{ such that } p(x) = \sum_{i=0}^n p^i v_i(x)$$

linear independence of  $v_0, v_1, \dots, v_n \Leftrightarrow p^i = 0 \quad i=0, \dots, n \Leftrightarrow p(x) = 0 \quad \forall x \in [a,b]$

Polyomial interpolation: (Given  $n+1$  points  $\{x_i\}_{i=0}^n$ )  
The algorithm takes a function (continuous) and returns a polynomial  $F: \mathcal{X} \equiv C^0([a,b]) \rightarrow \mathcal{Y} \equiv \mathbb{P}^n([a,b])$

Where  $p(x_i) = u(x_i) \quad i=0, \dots, n$  the given points

$$\sum_{j=0}^n V_j p^j = u(x_i) \quad \text{conditions}$$

contravariant coeff.  
and invert order reading to Einstein notation

$$\sum_{j=0}^n V_{ij} p^j = u_i \quad \Leftrightarrow \quad V_p = u$$

contravariant coeffs

$$V_{ij} = V_j(x_i)$$

Vandermonde Matrix

what we doing: going from a vector of values of a pt in  $C^0([a,b])$  to a vector of coefficients of a basis in  $\mathbb{P}^n([a,b])$  (where  $V_i(x) = \text{pow}(x, i)$ )

Einstein notation:

$$V_{ij} p^j \equiv \sum_{j=0}^n V_{ij} p^j$$

we set  $(V_{ij})^{-1} = V^{ij}$  so that

$$V^{ij} V_{jk} = \delta_k^i = \begin{cases} 1 & \text{if } i=k \\ 0 & \text{if } i \neq k \end{cases}$$

so if needed

$$V_{ij} p^j = u_i$$

$$p^j = V^{ji} u_i \quad \Leftrightarrow \quad p = V^{-1} u$$

We define the "Reduced problem":  $X = \mathbb{R}^n$

where

$X$ : set of values of  $u$  in  $x_i$

with norm  $\|u\|_2 := \left( \sum_{i=0}^n u_i^2 \right)^{\frac{1}{2}}$  and same for  $y$  polynomial

The original problem can be cast in this problem with 2 additional functions

We want to find Condition number of  $F$ :  $u \rightarrow V u = p$

that is input  $\{u_i\}_{i=0}^n \rightarrow$  output  $\{p^j\}_{j=0}^n$

If I was suppose to perturb the input, starting from

$$\|p - \hat{p}\|_2 \leq K_{abs} \|u - \hat{u}\|_2$$

We know ALREADY that the error is given by the norm of the matrix (find out starting from  $\|p - \hat{p}\|_2$  and using)

$$\|p - \hat{p}\|_2 \leq \|V\| \|u - \hat{u}\|_2$$

L' As long as  $V$  is invertible, problem\* is absolutely

The op. norm is defined here as:

$$\|V^{-1}\|_* := |\lambda(V^{-1})| = \max_i |\lambda_i(V^{-1})| = \|V^{-1}\|_2$$

operator induced by the 2-norm

To make the problem better conditioned, we need to minimize  $K_{abs} = \|V^{-1}\|_2$  that is  $V \equiv I$

If you want to solve the prob of interpolation  
 $\Rightarrow$  solve a linear system  
 $\Rightarrow$  condition number of problem unrelated to condition number of  $V$

LINEAR PROBLEM now!

then we have

$$\nabla = I \rightarrow v_j(x_i) = s_{ji}$$

which means we have a

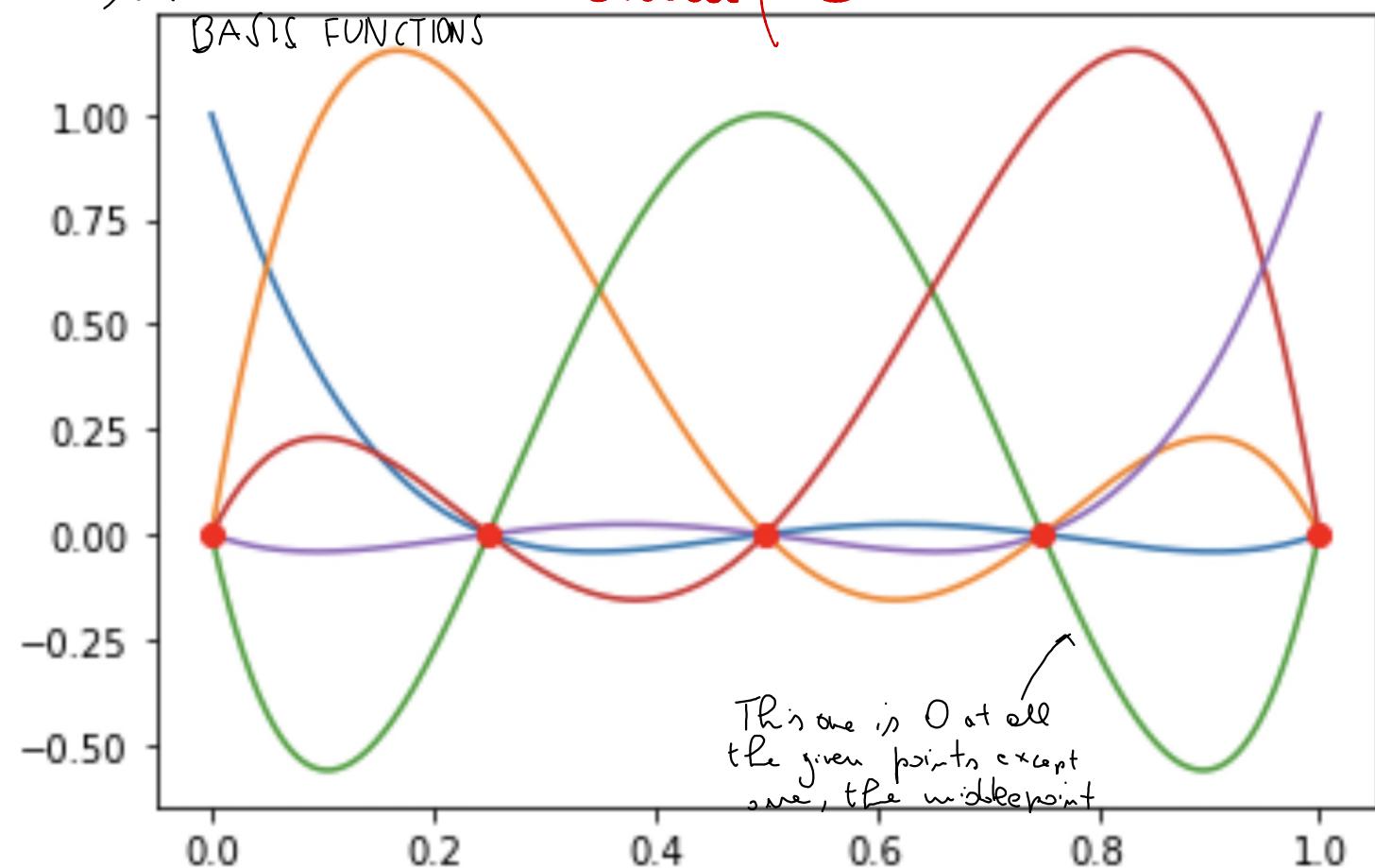
Polynomial that is zero at every point  $x_i$  except on  $x_5$  where it is 1. We call this polynomial

$$e_i := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x - x_j)}{(x_i - x_j)}$$

poly of order n  
=>  $e_i(x_j \neq x_i) = 0, e_i(x_i) = 1$

Lagrange Basis

example with  $n = 5$



This example shows you how to choose the basis function such that

$$\nabla = I$$

And you can use directly these basis functions to construct the polynomial interpolation

$$p(x) = \sum_{j=0}^n u(x_j) v_j(x)$$

where the coefficient are exactly

$$P_j = u(x_j)$$

because  $\nabla^{-1} = I$

If you have columns of  $\nabla$  that are very similar to one another (f.e. because they get closer and closer to one another)  $\nabla$  will not be invertible (the condition number will be very bad!)

The condition number of Vandermonde matrix will deteriorate as we increase the number of points (simply because, the more points we have the similar all the basis functions will get!)

# Numerical Analysis — lec 4 L+H —

General interpolation problem from  $\mathcal{X} := C^0([a, b])$  to the space  $\mathcal{Y} := V := \text{span}\{v_i\}_{i=0}^n$  with fixed support points  $\{x_i\}_{i=0}^n$

$$\text{II: } C^0([a, b]) \longrightarrow V$$

$$u(x) \longrightarrow p(x) := \sum_{i=0}^n p^i v_i(x)$$

Such that  $p$  and  $u$  coincides in  $\{x_i\}_{i=0}^n$

indeed  
this is a linear system

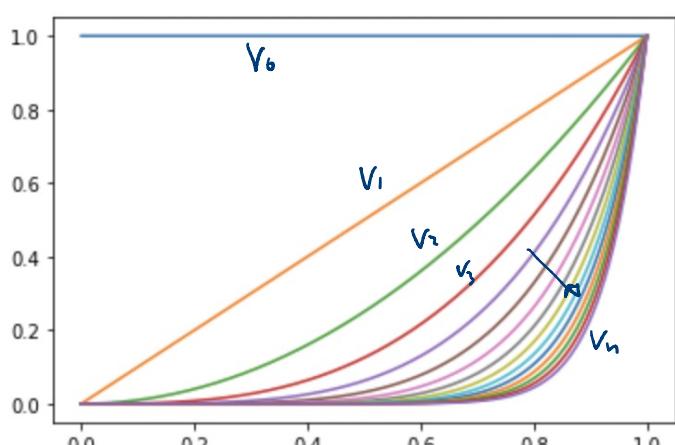
$$\sum_{j=0}^n p^j v_j(x_i) = u(x_i) \quad \equiv$$

$$\underline{V}_{ij} := v_j(x_i) \quad \underline{p}_i = \underline{u}$$

$$\{\underline{p}_i^j = p^j\}$$

$$\{\underline{u}_i^j = u(x_i)\}$$

POWER BASIS IN  $[0, 1]$   
Example  $V : [a, b] = [0, 1]$



$$\text{Linearity implies } \text{II}(u+v) = \text{II}(u) + \text{II}(v)$$

$$\text{because } \text{II}(u) = (\mathbb{V}^{-1} u)^i v_i(x) \quad \text{II}(v) = (\mathbb{V}^{-1} v)^i v_i(x)$$

$$\text{II}(u+v) = (\mathbb{V}^{-1} (u+v))^i v_i(x)$$

the choice of this number of points is important (in principle you could have even higher  $n$ )  
then it would be possible to satisfy exactly the relationship which the points coincide at some higher  $n$  of points

The basis functions  $v_i \in C^0([a, b])$   
(basis functions can be chosen among polynomials, cos/sin, sigmoid)

$$v_i = \text{pow}(x, i) = x^i$$

if we fix  $n+1=5$   $\underline{p} := [1, 0, 0, 2, 3]$

$$p(x) = 1 + 2x^3 + 3x^4$$

$$= p^i v_i(x)$$

$$\underline{u} + \underline{v} = (\underline{u} + \underline{v})$$

$$(\underline{u} + \underline{v})_i = u(x_i) + v(x_i)$$

then here we use

Condition number for interpolation problem

We use as norm in  $C^0([a, b])$   $\|\cdot\|_\infty$ , the same in  $V$ :  
 explicit to be fact that the interval problem is the definition of the Chebyshev problem

$$\|\mathbb{I}(u) - \mathbb{I}(\hat{u})\|_\infty = \|\mathbb{I}(u - \hat{u})\|_\infty \leq \|\mathbb{I}\|_* \|u - \hat{u}\|_\infty$$

What is  $\|\mathbb{I}\|_*$ ?

Let's apply the definition

$$\|\mathbb{I}\|_* := \sup$$

domain of  
the function  
"solved" by  
the problem

$$u \in C^0([a, b])$$

$$\|\mathbb{I}(u)\|_\infty = \sup$$

$$u \in \mathcal{X}$$

norm of  
 $V$

norm of  
matrix  $V$

$P$   
is

$$\|\mathbb{I}(x_j) - v_i(x_j)\|_\infty$$

$$P = \mathbb{I}M$$

$$\leq \|\mathbb{V}^{-1}\|_\infty \max_i \|v_i\|_\infty$$

$$\leq \sup_{u \in \mathcal{X}} \|\mathbb{V}^{-1}\|_\infty \|u\|_\infty \max_i \|v_i\|_\infty$$

$$\text{because for some } n$$

$$\text{the value of } u$$

$$\text{of the elements}$$

$$\text{of } u \text{ at } x_j \text{ is a}$$

$$\text{linear combination}$$

$$\text{of the basis functions } v_i$$

$$\text{so } \|u\|_\infty = \max_i \|v_i\|_\infty$$

$$\text{and } \|\mathbb{V}^{-1}\|_\infty = \text{cond}(\mathbb{V})$$

$$\text{all norms are equivalent}$$

$$\Rightarrow \|\mathbb{V}^{-1}\|_\infty = \text{cond}(\mathbb{V})$$

$$\text{and number of the problem}$$

$$\text{determines if matrix determinant}$$

$$\text{when basis function becomes too large}$$

$$\text{Two things: } \{x_i\}_{i=1}^n, \{v_i\}_{i=0}^n$$

$$\Rightarrow \min \text{cond}(\mathbb{V}^{-1})$$

$$\text{are}$$

$$\text{the Lagrange basis functions:}$$

$$\ell_i := \prod_{j=0}^{n-1} \frac{(x - x_j)}{(x_i - x_j)}$$

$$\text{L} \rightarrow \mathbb{V} = \mathbb{I} \ell_i$$

$$\text{by construction}$$

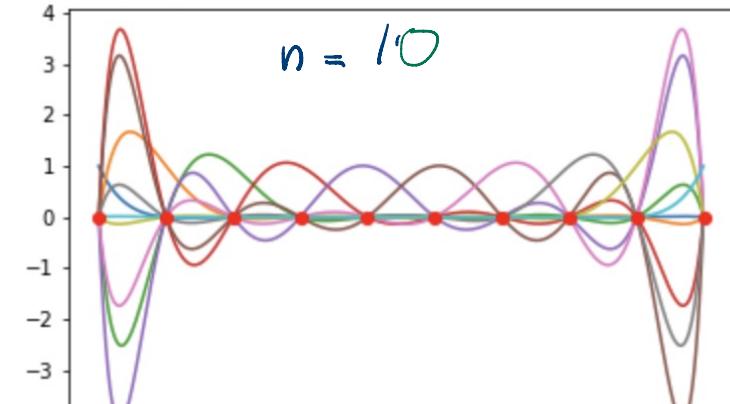
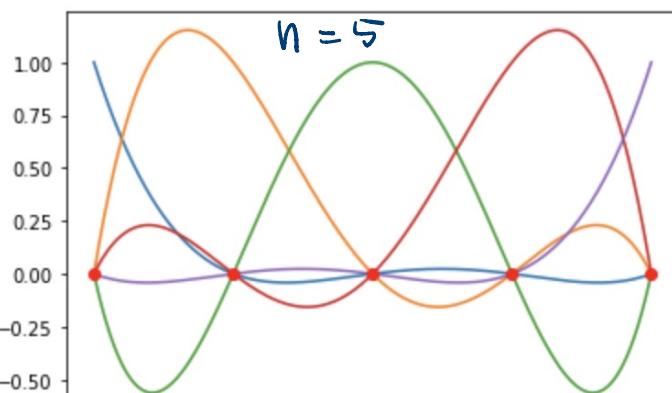
$$\text{cond}(\mathbb{V})$$

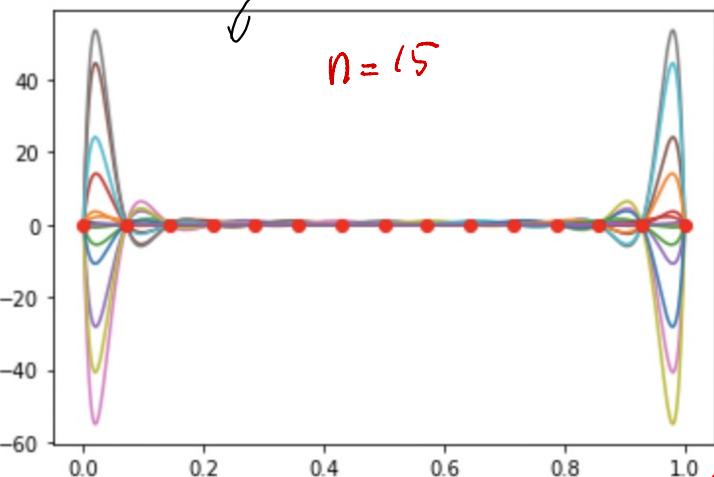
$$\text{for power basis}$$

$$\text{is}$$

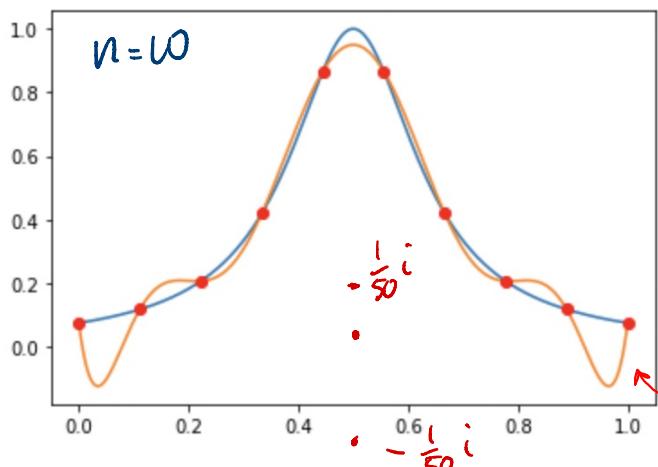
$$\mathbb{V}_{ij} := \ell_j(x_i) = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

$$\text{We should choose } x_i \text{ s.t. } \max_i \|\ell_i\|_\infty \text{ is small}$$





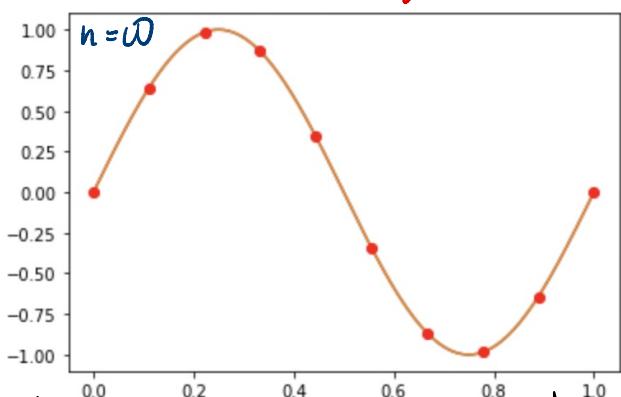
then you start getting bad oscillations  
 ↳ Is interpolation  
 a good strategy  
 for approximating  
 a function?



The orange line represents using  $n$  LAGRANGE interpolation  
 $\sum_{i=0}^n u(x_i) l_i(x) = p(x)$  with Lagrange basis

$$\text{blue line is the interpolated function } f(x) = \frac{1}{1 + 50(x - \frac{1}{2})^2}$$

BAD BEHAVIOR NEAR (away from the center while at the center things are good but this is FUNCTION DEPENDENT  
 $f(x) = \sin(2\pi x)$ )



for sine results are perfect

Can we do better?

What can we say for smooth functions about interpolation properties of polynomials

Thus  $f \in C^{n+1}([a, b])$  if  $x_i \in (a, b)$ ,  $\forall x \in (a, b) \exists \xi \in (a, b)$   
 where  $n+1$  points  $\Rightarrow$  we have  $n+1$  point derivatives

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

where  $\omega(x)$  is the characteristic polynomial of  $\{x_i\}_{i=1}^n$

$$\omega(x) = \prod_{i=0}^n (x - x_i)$$

can we make these two numbers small?

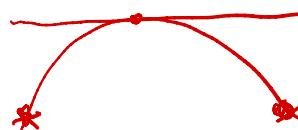
First consequence:  $\|I^n f - f\|_\infty \leq \|w\|_\infty \|f^{(n+1)}\|_\infty$   
 If you control the  $n+1$  derivative of the function ANYWHERE  $\rightarrow$  you have a bound on the norm of the  $n$  function

Proof: Let  $x$ , define  $G(t)$  s.t. it is equal to 0 in all points, in which the polynomials and the function coincide plus the point  $x$  where the polynomials and the function coincide. Then  $G(t) = (f(t) - p(t)) w(x) - (f(x) - p(x)) w(t)$

$$\text{where } p(t) = \sum_{i=0}^n f(x_i) e_i(t) = (\text{If})(t) \quad \text{~\textcircled{1}~ interpolation function}$$

$G(t)$  has  $n+2$  zeros:  $\{x_i\}_{i=0}^n \cup \{x\}$

According to Rolle's theorem:  $\exists \xi \in (a, b)$  s.t.  $\frac{d^{n+1} G(\xi)}{dt^{n+1}} = 0$



With two zeros  $\Rightarrow \exists \xi \in (a, b)$  s.t.  $f'(\xi) = 0$

writing explicitly the derivative

$$\frac{d^{n+1}}{dt^{n+1}} G(\xi) = f^{(n+1)}(\xi) \cdot w(x) - (f(x) - p(x))(n+1)! = 0$$

$\frac{d^{n+1}}{dt^{n+1}} (w(t))$

$\Rightarrow f(x) - (\text{If})(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} w(x) \quad \square$  QED

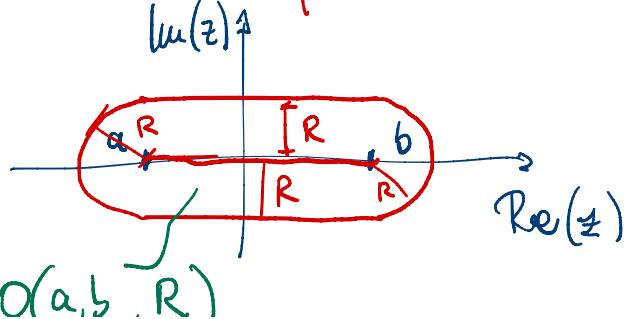
How to apply this?

Theorem 2: If  $f$  is analytically extendible in an oval  $O(a, b, R)$  with  $R > 0$

Then  $\|f^{(n+1)}\|_\infty \leq \frac{(n+1)!}{R^{n+1}} \|\tilde{f}\|_{L^\infty(O(a,b,R))}$

Analytically extendible means:  $\tilde{f}$

Take  $z$ , and replace  $t$  with  $z \in \mathbb{C}$



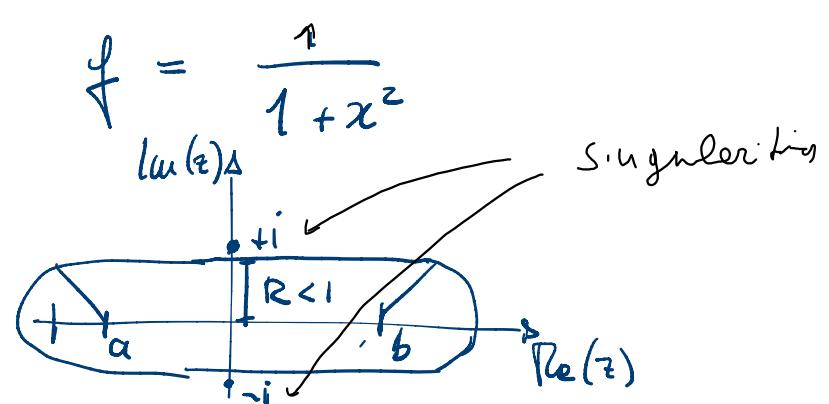
we are asking that  $f \in C^\infty$  in  $O(a, b, R)$

EXAMPLE

Runge function:

extension to  $\mathbb{C}$

$$f(z) = \frac{1}{1+z^2}$$



$f$  is analytical in  $O(-\infty, +\infty, 1-\varepsilon)$   $\forall \varepsilon > 0$

What are the consequences?

$$\|f^{(n)}\| \leq \frac{(n+1)!}{R^{n+1}} \|f\|_\infty$$

$R < 1$  stabilizing  
at finite singularities

So, by the previous theorem:

$$\|\rho - f\|_\infty \leq \|f^{(n+1)}\|_\infty \frac{\|\omega\|_\infty}{(n+1)!} \leq \frac{(n+1)!}{(n+1)!} \frac{\|\omega\|_\infty}{R^{n+1}} \|f\|_\infty$$

where

$$\|\omega\|_\infty = \left\| \prod_{i=0}^n (x-x_i) \right\| \leq \underbrace{|b-a|^{n+1}}_{\text{since } x \in (a,b)} \quad x \in (a,b)$$

$$\|\rho - f\|_\infty \leq \left( \frac{|b-a|}{R} \right)^{n+1} \|f\|_\infty$$

$\hookrightarrow$  good only when  $|b-a| < R$ , otherwise  $\|\rho - f\|_\infty$  is NOT BOUNDED

$\Rightarrow$  Runge:

$$\frac{1}{1+x^2} \quad \text{in } (-1, 1)$$

everything ok

Can we do better?

Best Approximation. Definition:

Restrict ourselves

$\mathcal{K}$  Banach, reflexive, strictly convex  
✓ subspace of  $\mathcal{K}$

We call  $p$  the best approximation in  $V$  of  $f$  in  $\mathcal{K}$   
when :

$$\|f-p\| \leq \|f-q\| \quad \forall q \in V$$

in other words

$$\|f-p\| = E(f) = \inf_{q \in V} \|f-q\|$$

Theo:  $\exists! p$  satisfying these  
certain (mandatory conditions)

How does polynomial interpolation compare to best approximation?

We can say that

$$\|f - I^n f\|_\infty = \|f - p + p - I^n f\|_\infty =$$

add on subtract best approx

$$= \|f - p + I^n(p - f)\|_\infty$$

$p$  is a poly of od n  
 $\Rightarrow$  it's poly interpolation  
is itself

$$\text{triangle inequality} \leq \|f - p\|_\infty + \|I^n\|_* \|p - f\|_\infty$$

$$\text{property of norms} \leq (1 + \|I^n\|_*) \|p - f\|_\infty$$

For Lagrange interpolation:

$$\|I^n_*\| := \sup_{\substack{u \in C^0([a, b]) \\ =}} \frac{\left\| \sum_{i=0}^n u(x_i) L_i(x) \right\|_\infty}{\|u\|_\infty} \leq$$

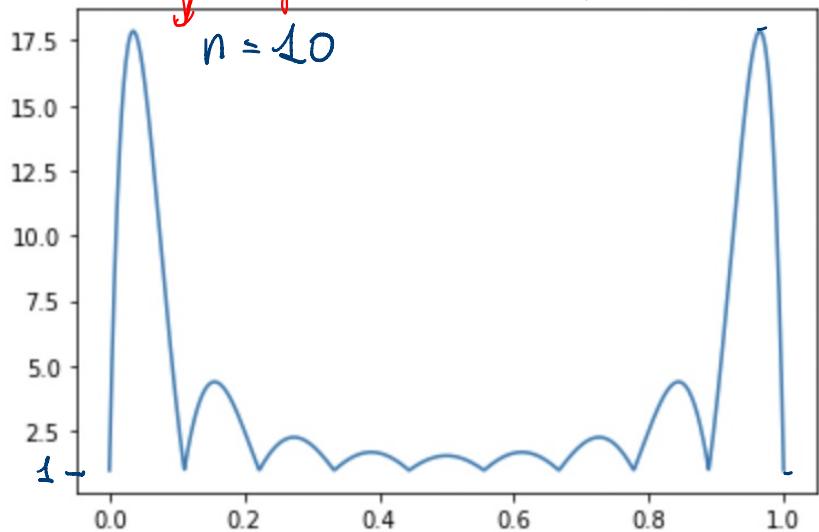
how away we  
are from best  
interpolation

$$= \left\| \sum_{i=0}^n |L_i(x)| \right\|_\infty \quad \text{def } L(x) := \sum_{i=1}^n |L_i(x)|$$

LEBESGUE Funktion

$$\Rightarrow \|\mathcal{I}_*^n\| = \|\mathcal{L}\|_\infty$$

Lebesgue functions are such that  $\Rightarrow$  bad results. Indeed  $\|\mathcal{L}\|_\infty$  for equispaced points



$\|\mathcal{L}\|_\infty$  for equispaced points

$$\|\mathcal{L}\|_\infty \leq \frac{2^{n+1}}{\pi} \text{ grows quickly}$$

A collection of points  $\{x_i\}_{i=0}^n$  is a list of  $n+1$  points with increasing size -

Endos

$\nexists$  collection of points  $\{x_i\}_{i=0}^n$

$\exists c > 0$  s.t.

Faber

$$\|\mathcal{L}\|_\infty \geq \frac{2}{\pi} \log(n+1) - c$$

of independently of how you choose the points, at least you get logarithmic with the number of the points

$\nexists$  collection of points  $\{x_i\}_{i=0}^n$

$\exists \delta$  s.t.

$$\lim_{n \rightarrow \infty} \|\mathcal{I}^n f - g\|_\infty \rightarrow \infty$$

( $\hookrightarrow$  you can always find a pt that destroys all)

The best is  $\alpha := \{x_i\}_{i=0}^n$  s.t.

$$\|\mathcal{L}^\alpha\|_\infty \leq \|\mathcal{L}^\beta\|_\infty$$

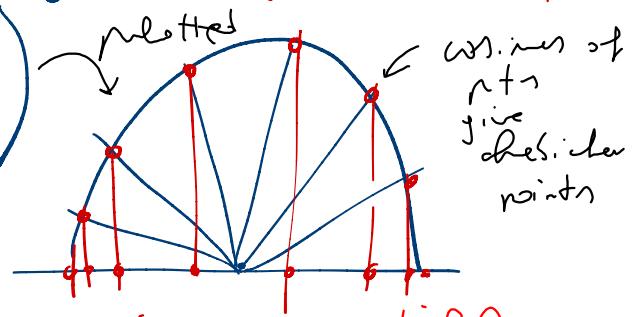
$\nexists \beta = \{x_i\}_{i=0}^n$   
such  $\alpha := \{x_i\}_{i=0}^n$

are called Chebyshev points:

between  $[-1, 1]$

defined as  $x_i = \cos\left(\frac{(2i+1)\pi}{2n+2}\right)$

$$\|\mathcal{L}^\alpha\|_\infty \leq \frac{2}{\pi} \log(n+1) + 1$$



Best case scenario

still grows exponentially

$$\frac{2}{\pi} \log(n+1) - c \leq \|\mathcal{L}^\alpha\|_\infty \leq \frac{2}{\pi} \log(n+1) + 1$$

$\Rightarrow$  INTERPOLATION IS GOOD BUT NOT  $\pi$  SO WEBS

# Nonlinear equations



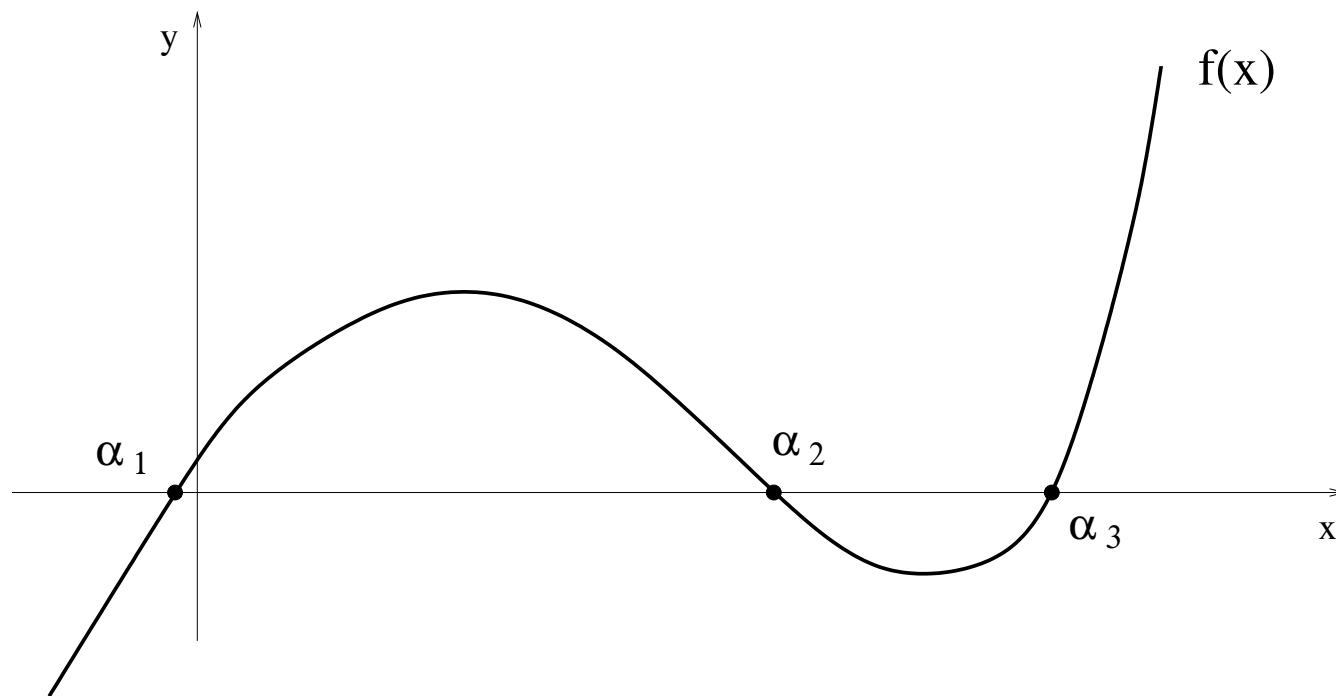
## Numerical Analysis

Profs. Gianluigi Rozza - Luca Heltai

2019-SISSA mathLab Trieste

# Nonlinear equations

**Objective:** Find the root of scalar (or vector) non-linear functions, i.e., find  $\alpha \in \mathbb{R}$  such that  $f(\alpha) = 0$ .



# Examples and motivation

**Example 1** (Interest rates). We want to compute the mean interest rate  $I$  of a portfolio over several years. We invest  $v = 1000$  Euro every year. After 5 years we end up with  $M = 6000$  Euro.

The relation between  $M$ ,  $v$ ,  $I$  and the number of years  $n$  is

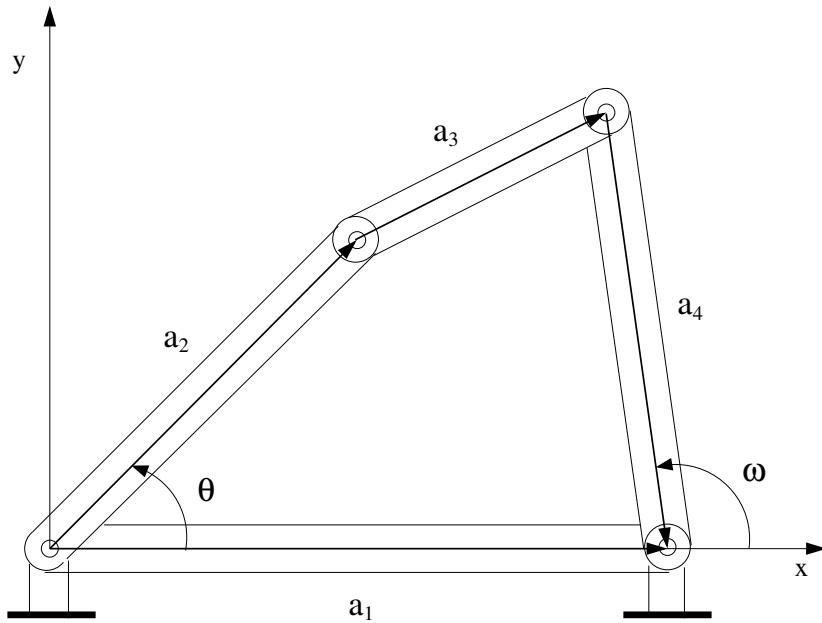
$$M = v \sum_{k=1}^n (1 + I)^k = v \frac{1 + I}{I} [(1 + I)^n - 1]$$

This can be rewritten as: *find  $I$  such that*

$$f(I) = M - v \frac{1 + I}{I} [(1 + I)^n - 1] = 0 \quad (1)$$

Therefore we have to solve a nonlinear equation in  $I$ , for which we can't find an analytical solution.

**Example 2** (Rods system). Let us consider the mechanical system represented by four rigid rods



For any admissible angle  $\omega$ , we want to compute the angle  $\theta$  between  $\mathbf{a}_1$  and  $\mathbf{a}_2$ .

Thanks to the vector identity

$$\mathbf{a}_1 - \mathbf{a}_2 - \mathbf{a}_3 - \mathbf{a}_4 = 0$$

and keeping  $\mathbf{a}_1$  on the  $x$ -axis, we can derive the following equation: involving  $\omega$  and  $\theta$ :

$$\frac{a_1}{a_2} \cos(\omega) - \frac{a_1}{a_4} \cos(\theta) - \cos(\omega - \theta) = -\frac{a_1^2 + a_2^2 - a_3^2 + a_4^2}{2a_2 a_4} \quad (2)$$

where  $a_i$  is the length of the  $i$ th rod. Equation (2) is nonlinear and can be solved only for particular values of  $\omega$ . For a general  $\omega$  it is not possible to find an analytic solution.

**Example 3** (State equation of a gas). We want to determine the volume  $V$  occupied by a gas at temperature  $T$  and pressure  $p$ . The state equation (i.e. the equation that relates  $p$ ,  $V$  et  $T$ ) is

$$\left[ p + a \left( \frac{N}{V} \right)^2 \right] (V - Nb) = kNT ,$$

where  $a$  and  $b$  are two coefficients that depend on the specific gas,  $N$  is the number of molecules which are contained in the volume  $V$  and  $k$  is the Boltzmann constant. We need therefore to solve a non-linear equation whose root is  $V$ .

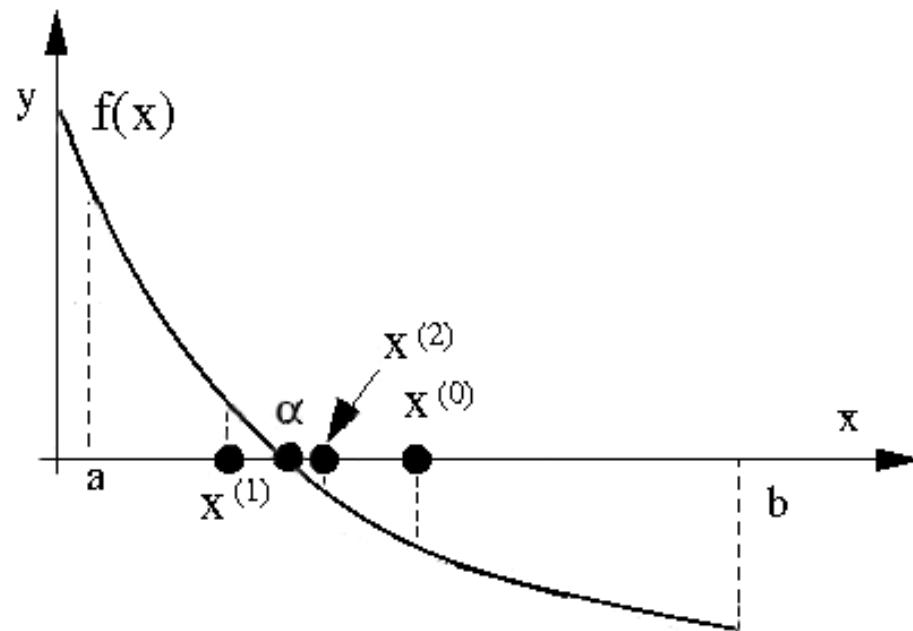


# Bisection method

(Book Sec. 2.2)

This method is used to compute the root of a **continuous** function  $f$ , i.e., the point  $\alpha$  such that  $f(\alpha) = 0$ . We can build a sequence  $x^{(0)}, x^{(1)}, \dots, x^{(k)}, (x^{(0)}$  such that  $\lim_{k \rightarrow \infty} x^{(k)} = \alpha$ .

We assume that  $f : (a, b) \rightarrow \mathbb{R}$  and  $a < b$ . If  $f(a)f(b) < 0$ , since  $f$  is continuous, we know that there exists (at least) one root  $\alpha$  of  $f$  in the interval  $(a, b)$ .



Then

1. We set  $a^{(0)} = a$ ,  $b^{(0)} = b$  and  $x^{(0)} = \frac{a^{(0)} + b^{(0)}}{2}$ ,
2. if  $f(x^{(0)}) = 0$ , then  $x^{(0)}$  is the zero.
3. if  $f(x^{(0)}) \neq 0$ , then:
  - (a) if  $f(x^{(0)})f(a^{(0)}) > 0 \Rightarrow$  the zero  $\alpha \in (x^{(0)}, b^0)$  and we define  $a^{(1)} = x^{(0)}$ ,  $b^{(1)} = b^{(0)}$  and  $x^{(1)} = (a^{(1)} + b^{(1)})/2$
  - (b) if  $f(x^{(0)})f(a^{(0)}) < 0 \Rightarrow$  the zero  $\alpha \in (a^{(0)}, x^{(0)})$  and we define  $b^{(1)} = x^{(0)}$ ,  $a^{(1)} = a^{(0)}$  et  $x^{(1)} = (a^{(1)} + b^{(1)})/2$

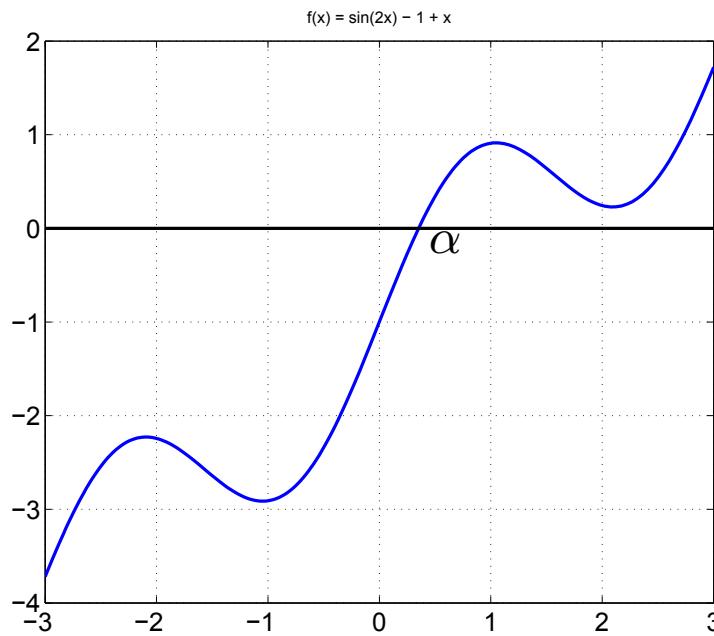
By the divisions of this type, we construct the sequence  $x^{(0)}, x^{(1)}, \dots, x^{(k)}$  that satisfies for all  $k$ ,

$$|e^{(k)}| = |x^{(k)} - \alpha| \leq \frac{b - a}{2^{k+1}},$$

because we divide the interval by 2 at every step, starting from step 0

**Example 4.** We want to find the zero of the function  $f(x) = \sin(2x) - 1 + x$ . We draw the graph of the function  $f$  using the following commands in Matlab/Octave:

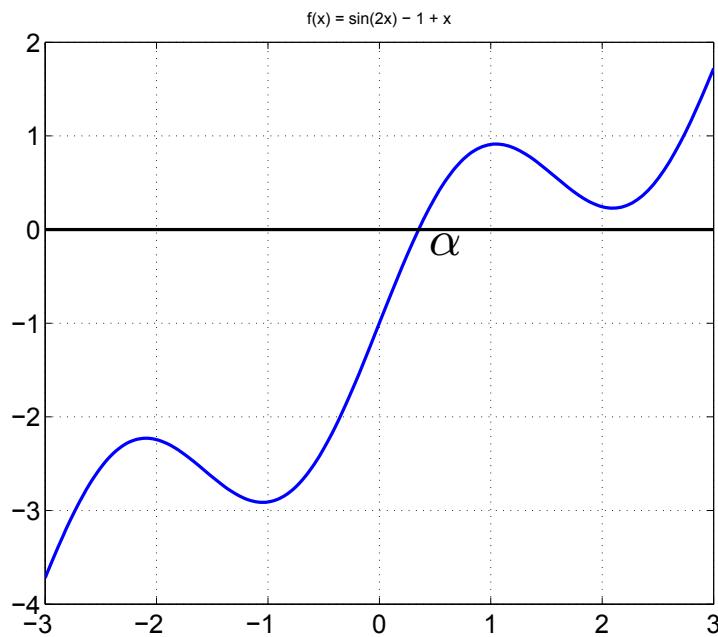
```
>> f = @(x) sin(2*x) - 1 + x;
>> x=-3:0.1:3;
>> plot(x,f(x)); grid on;
```



If we apply the bisection method in the interval  $[-1, 1]$  with a tolerance  $10^{-8}$  and maximum number of iterations  $k_{max} = 1000$

```
>> [zero,res,niter]=bisection(f,-1,1,1e-8,1000);
```

We find the value  $\alpha = 0.352288462$  after 27 iterations.



## APRIORI KNOWLEDGE OF # OF MAX ITERATION

Assume stopping criterion is based on size of interval.

At  $k$ :

$$\frac{b-a}{2^k} \leq \varepsilon \quad \rightarrow \text{solve for } k: \quad k \geq \log_2 \frac{b-a}{\varepsilon} = \frac{\ln(b-a) - \ln(\varepsilon)}{\ln 2}$$

$\uparrow$   
set tolerance

# Newton's method

(Chapt. 2.3 of the book)

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a differentiable function.

Let  $x^{(0)}$  be an initial guess. Let us consider the equation  $y(x)$  which passes through the point  $(x^{(k)}, f(x^{(k)}))$  and which has the slope  $f'(x^{(k)})$ :

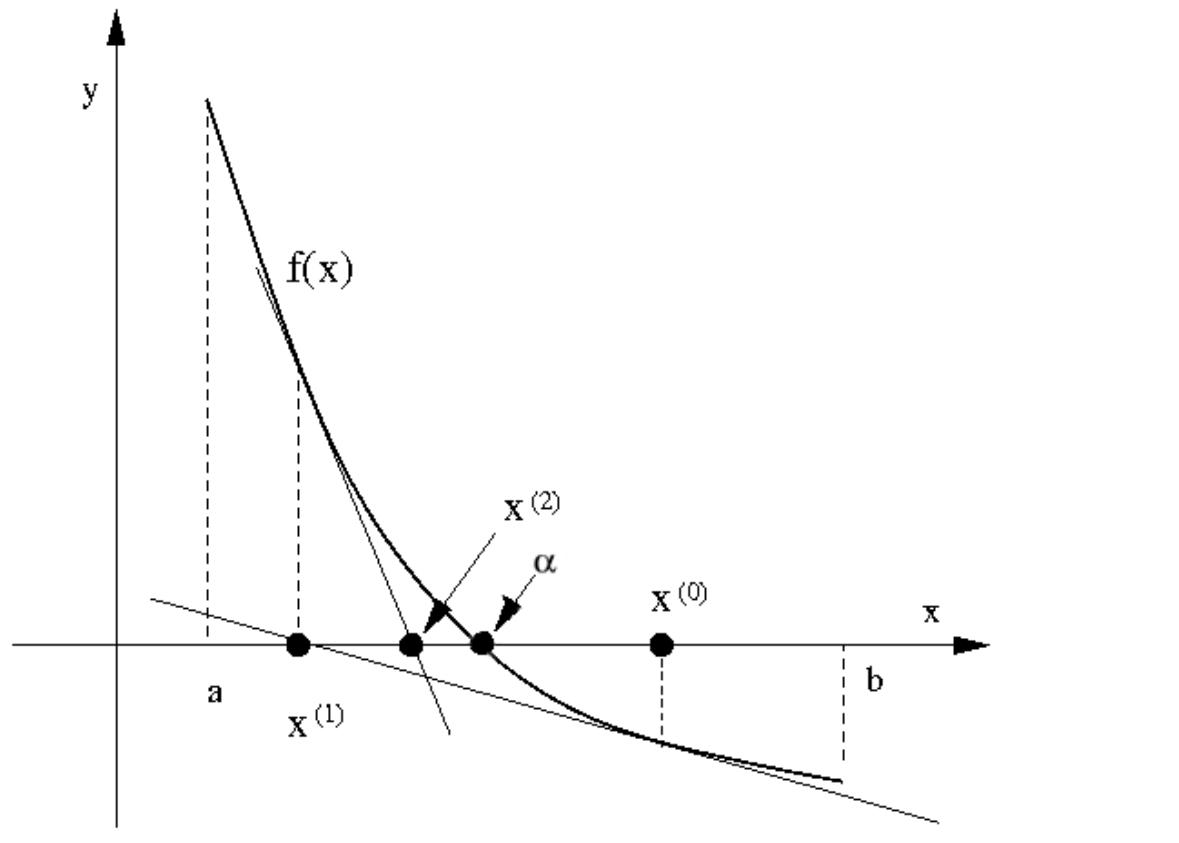
$$y(x) = f'(x^{(k)})(x - x^{(k)}) + f(x^{(k)}).$$

We define  $x^{(k+1)}$  by the point where this line intersects the axis  $x$ , i.e.  $y(x^{(k+1)}) = 0$ . We deduce that:

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k = 0, 1, 2, \dots \quad (3)$$

# Newton's method

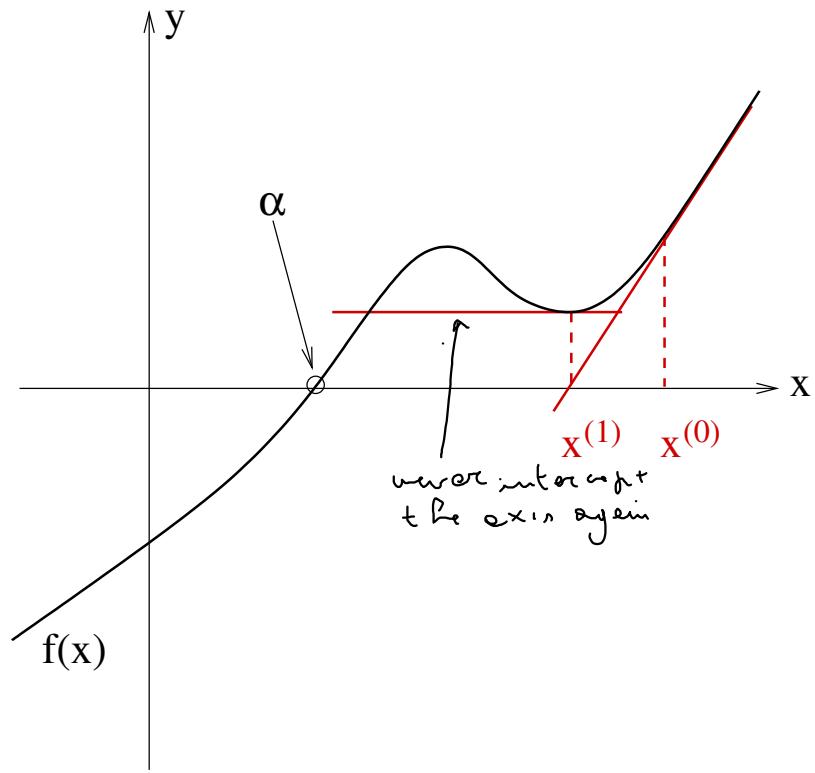
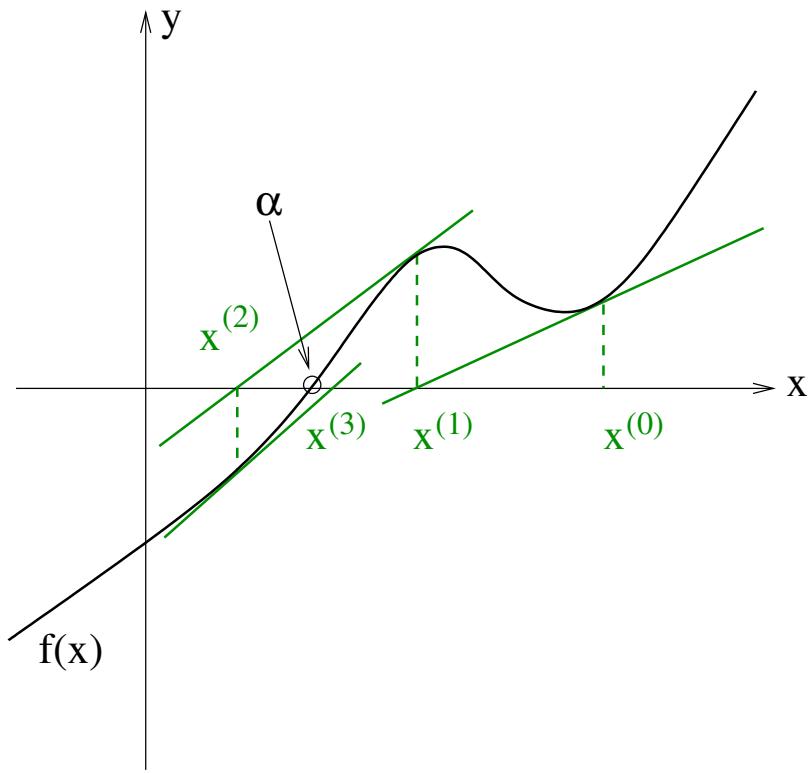
Starting from the point  $x^{(0)}$ , the sequence  $\{x^{(k)}\}$  converges to the root of  $f$



# Convergence?

Does this method always converge?

- it depends on the **property of the function**;
- and on the **initial guess**.



# Fixed point iterations.

(Chapt. 2.4 in the book)

A general method for finding the roots of a nonlinear equation  $f(x) = 0$  is the transformation<sup>of the problem in</sup> an equivalent problem  $x - \phi(x) = 0$ , where the auxiliary function  $\phi : [a, b] \rightarrow \mathbb{R}$  must have the following property :

$$\phi(\alpha) = \alpha \quad \text{if and only if} \quad f(\alpha) = 0.$$

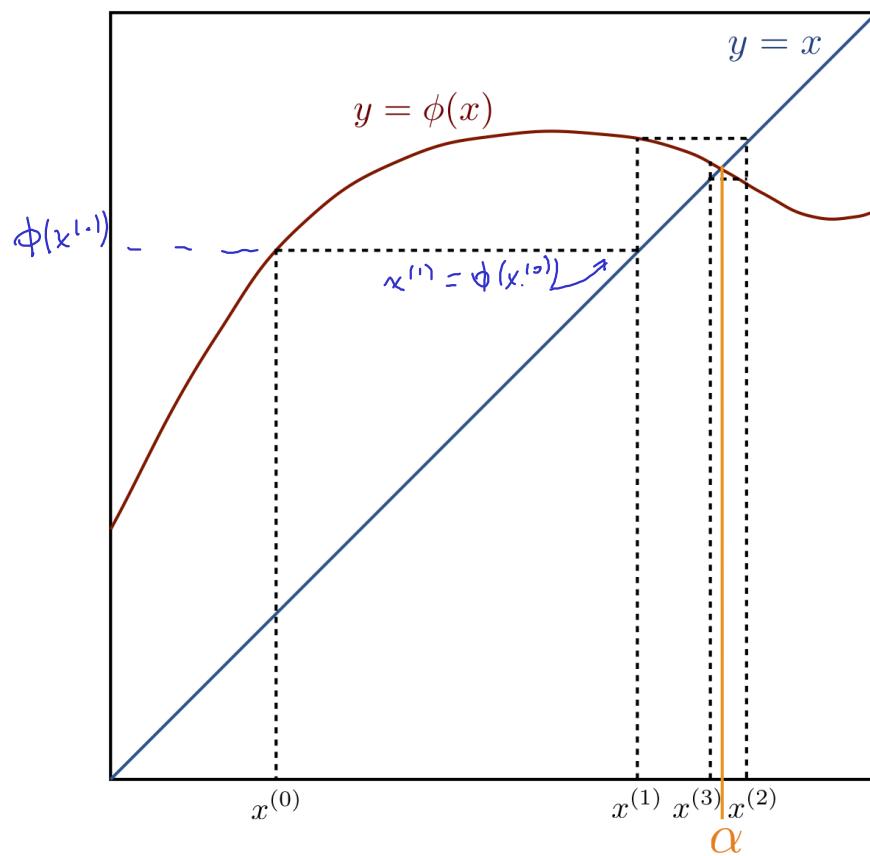
*↙ and  $\phi$  is called CONTRACTION*

The point  $\alpha$  is called *a fixed point* of  $\phi$ . Searching the zeros of  $f$  is reduced to the problem of determining the fixed points of  $\phi$ .

**Idea :** It could be computed by the following algorithm:  $x^{(k+1)} = \phi(x^{(k)})$ ,  $k \geq 0$ . Indeed, if  $x^{(k)} \rightarrow \alpha$  and if  $\phi$  is continuous on  $[a, b]$ , then the limit  $\alpha$  satisfies  $\phi(\alpha) = \alpha$ .

$$x^{(k+1)} = \phi(x^{(k)}) \xrightarrow{x^{(k)} \rightarrow \alpha} \phi(\alpha) \xrightarrow{\text{follows from the definition of ft } \phi} \alpha$$

Starting from the point  $x^{(0)}$ , the sequence  $\{x^{(k)}\}$  converges to the fixed point  $\alpha$

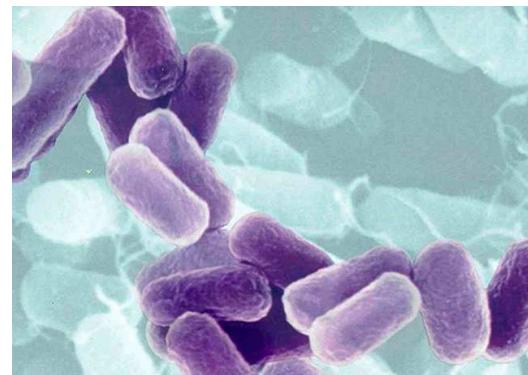


## Example 5 (Population dynamics).

In the survey population (e.g. bacteria), we want to establish a link between the number of individuals in the generation  $x$  and the number of individuals in the next generation  $x^+$ :

$$x^+ = \phi(x) = xR(x), \quad (4)$$

where  $R(x)$  represents the rate of growth (or decay) of the considered population.



Several models are available for  $R(x)$ :

- The Malthusian growth model (Thomas Malthus 1766-1834),

$$x^+ = \phi_1(x) = xR_1(x) \text{ with } R_1(x) = r, \quad r \text{ is a positive constant}$$

- the model of growth with limited resources (Pierre François Verhulst, 1804-1849),

$$x^+ = \phi_2(x) = xR_2(x) \text{ with } R_2(x) = r/(1 + x/K), \quad r > 0, K > 0$$

$K$  may be some kind of measurement for  
 available resources  $\rightarrow x/K$  is the number  
 of members of population per resource

that improve the Malthausian growth model by taking into account the growth of a population is limited by the resources.

- the predator/prey model

$$x^+ = \phi_3(x) = xR_3(x) \text{ with } R_3(x) = rx/(1 + (x/K)^2)$$

that represents the change of the Verhulst model by presence of an antagonist population.

The dynamics of a population is defined by a iterative process, starting from a given initial guess ( $x^{(0)}$ ),

$$x^{(k+1)} = \phi(x^{(k)}), \quad k \geq 0,$$

where  $x^{(k)}$  represents the number of individuals in  $k$ -th generation.  
 In addition, the steady states (equilibriums)  $x^*$  of a considered population are identified by the following problem,

$$x^* = \phi(x^*), \tag{5}$$

or equivalently,

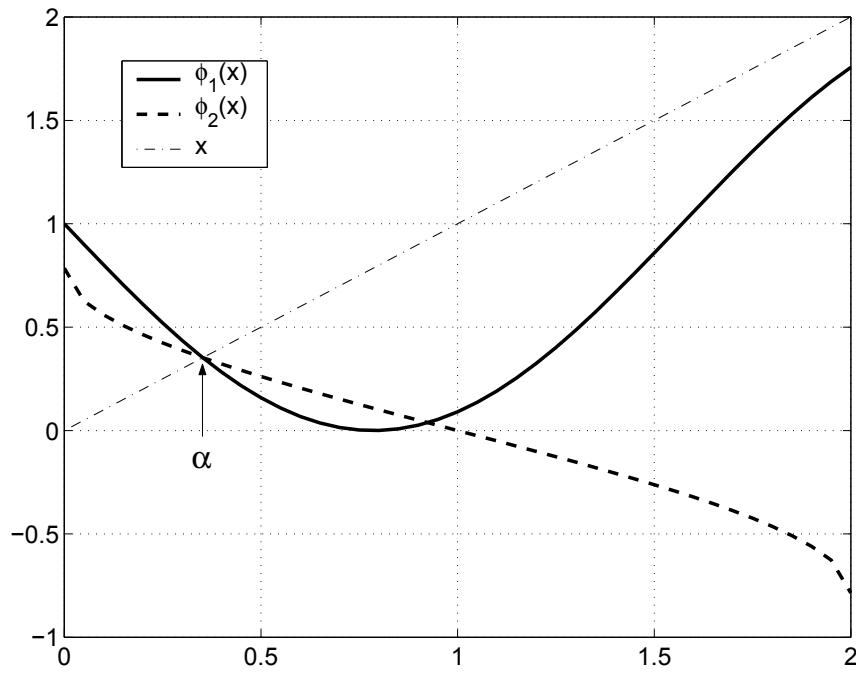
$$x^* = x^* R(x^*), \quad \text{i.e.} \quad R(x^*) = 1. \tag{6}$$

In both cases, we ask to solve a non-linear problem. In particular the problem is called the fixed point problem.

**Example 6.** We consider the equation  $f(x) = \sin(2x) - 1 + x = 0$ . We can rewrite it in two different fashions:

$$x = \phi_1(x) = 1 - \sin(2x)$$

$$x = \phi_2(x) = \frac{1}{2} \arcsin(1 - x), \quad 0 \leq x \leq 1$$



**Proposition 1.** (*Global convergence*)

that is: the pair  $(x, \phi(x))$   
 belongs to a square

1. Assume that  $\phi(x)$  is continuous on  $[a, b]$  and such that  $\phi(x) \in [a, b]$  for all  $x \in [a, b]$ ; then there exists at least one fixed point  $\alpha \in [a, b]$  of  $\phi$ .

2. If  $\exists L < 1$  such that  $|\phi(x_1) - \phi(x_2)| \leq L|x_1 - x_2| \quad \forall x_1, x_2 \in [a, b]$ ,

then there exists a unique fixed point  $\alpha \in [a, b]$  and the sequence

$x^{(k+1)} = \phi(x^{(k)})$ ,  $k \geq 0$  converges to  $\alpha$ , for any initial guess  $x^{(0)} \in [a, b]$ .

*Proof.*

$$\phi(x) \in [a, b] \quad \forall x \in [a, b]$$

because  $\phi(x)$  and  $x$  are continuous

1. The function  $g(x) = \phi(x) - x$  is continuous in  $[a, b]$  and, thanks to assumption made on the range of  $\phi$ , it holds  $g(a) = \phi(a) - a \geq 0$  and  $g(b) = \phi(b) - b \leq 0$ . By applying the theorem of zeros of continuous functions, we can conclude that  $g$  has at least one zero in  $[a, b]$ , i.e.  $\phi$  has at least one fixed point in  $[a, b]$ .
2. Indeed, should two different fixed points  $\alpha_1$  and  $\alpha_2$  exist, then

Applying I property of  $\phi$

$$|\alpha_1 - \alpha_2| = |\phi(\alpha_1) - \phi(\alpha_2)| \leq L|\alpha_1 - \alpha_2| < |\alpha_1 - \alpha_2|,$$

Roll?

which cannot be. There exists a unique fixed point  $\alpha \in [a, b]$  of  $\phi$ . □

Let  $x^{(0)} \in [a, b]$  and  $x^{(k+1)} = \phi(x^{(k)})$ . We have

$$0 \leq |x^{(k+1)} - \alpha| = \underbrace{|\phi(x^{(k)}) - \phi(\alpha)|}_{\text{1}} \leq L|x^{(k)} - \alpha| \leq \dots \leq L^{k+1}|x^{(0)} - \alpha|,$$

↗ by applying 1 and 2  
 iteratively

i.e.

$$\frac{|x^{(k)} - \alpha|}{|x^{(0)} - \alpha|} \leq L^k.$$

Because  $L < 1$ , for  $k \rightarrow \infty$ , we obtain

$$\lim_{k \rightarrow \infty} |x^{(k)} - \alpha| \stackrel{\text{dk, because } |x^{(0)} - \alpha| \text{ is constant}}{\leq} \lim_{k \rightarrow \infty} L^k = 0.$$

So,  $\forall x^{(0)} \in [a, b]$ , the sequence  $\{x^{(k)}\}$  defined by  $x^{(k+1)} = \phi(x^{(k)})$ ,  $k \geq 0$  converges to  $\alpha$  when  $k \rightarrow \infty$ .

### Remark 1.

If  $\phi(x)$  is differentiable in  $[a, b]$  and

$\exists K < 1$  such that  $|\phi'(x)| \leq K \ \forall x \in [a, b]$ ,

*just a matter  
of scale*

then the condition 2 of the proposal (1) is satisfied. This assumption is stronger, but is more often used in practice because it is easier to check.

**Definition 1.** For a sequence of real numbers  $\{x^{(k)}\}$  that converges,  $x^{(k)} \rightarrow \alpha$ , we say that the convergence to  $\alpha$  is **linear** if exists a constant  $C < 1$  such that, for  $k$  that is large enough

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|.$$

If exists a constant  $C > 0$  such that the inequality

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|^2$$

is satisfied, we say that convergence is **quadratic**.

In general, the convergence is **with order  $p$** ,  $p \geq 1$ , if exists a constant  $C > 0$  (with  $C < 1$  when  $p = 1$ ) such that the following inequality is satisfied

  
 =>  
 $\overset{C>1}{\text{dove}} \text{mt}$   
 converg

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|^p.$$

**Proposition 2.** (*Local convergence - Theorem 2.1 in the book*)

Let  $\phi$  be a continuous and **differentiable** function on  $[a, b]$  and  $\alpha$  be a fixed point of  $\phi$ . If  $|\phi'(\alpha)| < 1$ , then there exists  $\delta > 0$  such that, for all  $x^{(0)}$ ,  $|x^{(0)} - \alpha| \leq \delta$ , the sequence  $\{x^{(k)}\}$  defined by  $x^{(k+1)} = \phi(x^{(k)})$  converges to  $\alpha$  when  $k \rightarrow \infty$ .

Moreover, it holds

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{x^{(k)} - \alpha} = \phi'(\alpha).$$

Note that, if  $0 < |\phi'(\alpha)| < 1$ , then for any constant  $C$  such that  $|\phi'(\alpha)| < C < 1$ , if  $k$  is large enough, we have:

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|.$$

Proof: expand  $\phi(x^{(k)})$  around  $\alpha$ , proceed as in proof of next part

$\frac{x^{(k+1)} - \alpha}{x^{(k)} - \alpha}$  will reach  $\phi'(\alpha)$   
 from above or below  
 $\Rightarrow$  a number, it will be  $\leq C$

**Proposition 3.** (Proposition 2.2 in the book) *use the result of the previous preparation*

Let  $\phi$  be a twice differentiable function on  $[a, b]$  and  $\alpha$  be a fixed point of  $\phi$ . Let us consider that  $x^{(0)}$  converges locally to  $\alpha$ . If  $\phi'(\alpha) = 0$  and  $\phi''(\alpha) \neq 0$ , then the fixed point iterations converges with order 2 and

QUADRATIC CONVERGENCE

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{(x^{(k)} - \alpha)^2} = \frac{\phi''(\alpha)}{2}.$$

*Proof.* Using the Taylor series for  $\phi$  with  $x = \alpha$ , we have

$$x^{(k+1)} - \alpha = \phi(x^{(k)}) - \phi(\alpha) = \underbrace{\phi'(\alpha)}_0 (x^{(k)} - \alpha) + \frac{\phi''(\eta)}{2} (x^{(k)} - \alpha)^2$$

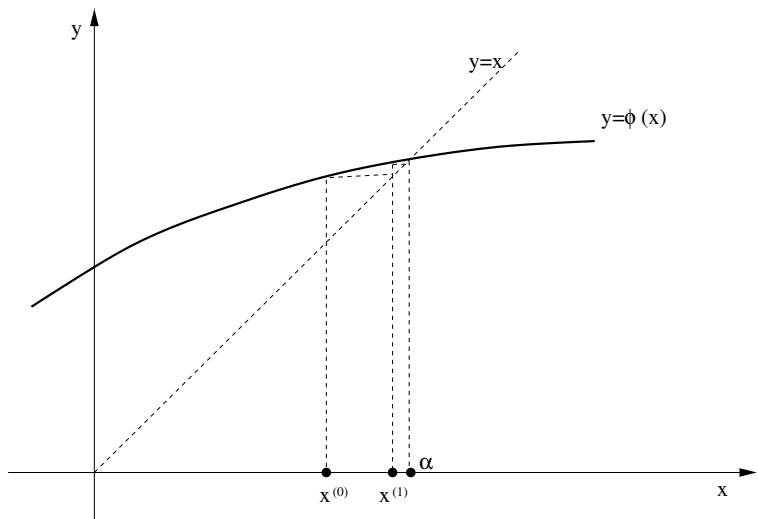
where  $\eta$  is between  $x^{(k)}$  and  $\alpha$ . So, we have

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{(x^{(k)} - \alpha)^2} = \lim_{k \rightarrow \infty} \frac{\phi''(\eta)}{2} = \frac{\phi''(\alpha)}{2}.$$

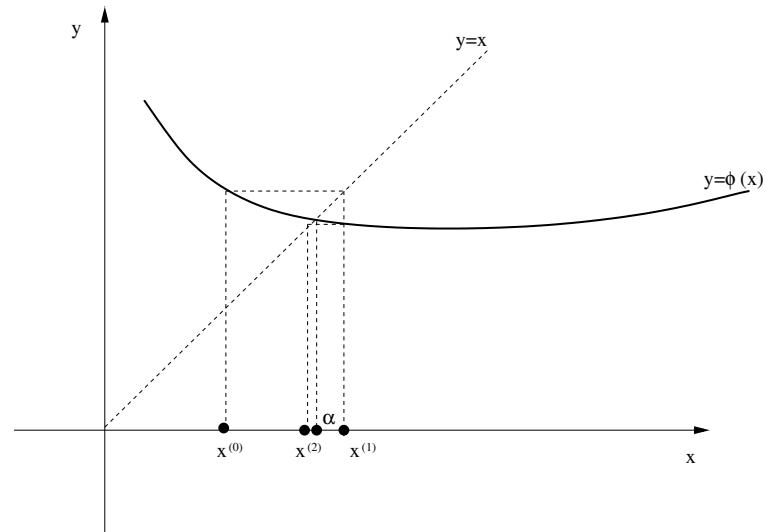
If we didn't have this, Taylor would be meaningless since  $x^{(k+1)}$  and  $x^{(k)}$  could be anywhere but near  $\alpha$ .

Some examples on how the value  $| \phi'(\alpha) |$  influences the convergence  
**Convergent cases:**

$$0 < \phi'(\alpha) < 1,$$



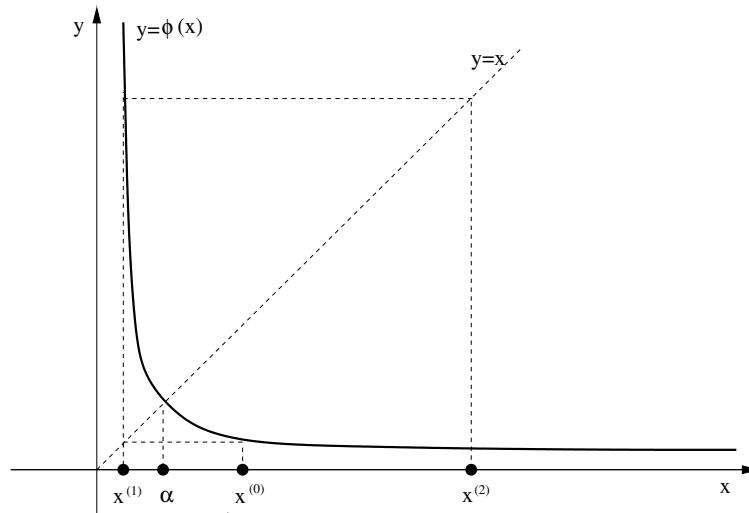
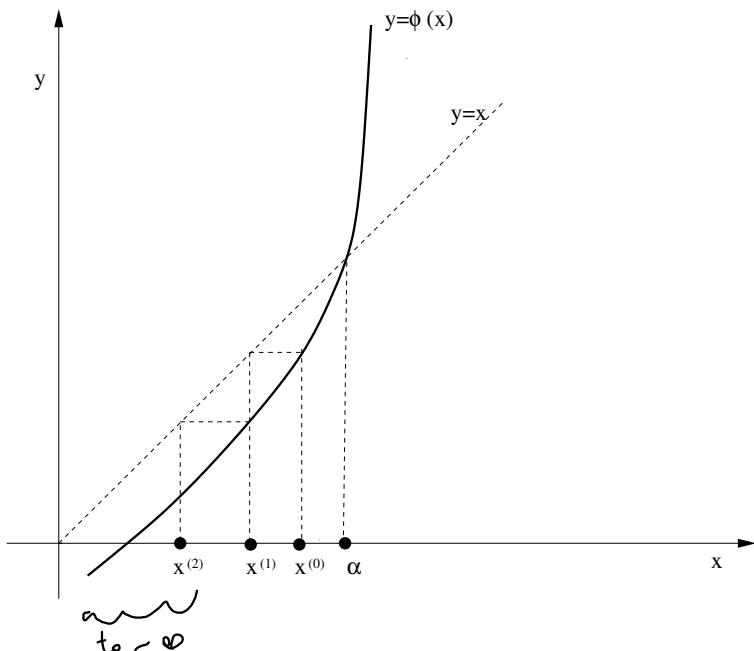
$$-1 < \phi'(\alpha) < 0.$$



## Divergent cases:

$$\phi'(\alpha) > 1,$$

$$\phi'(\alpha) < -1.$$



**Example. 5 (suite)** We apply the fixed point iterations on functions  $\phi_2(x) = rx/(1 + x/K)$  and  $\phi_3(x) = rx^2/(1 + (x/K)^2)$  that represent the Verhulst model and predator/prey model respectively, with  $K = 1.5$  and  $r = 2$ . We consider the starting point  $x^{(0)} = 1.0$ .

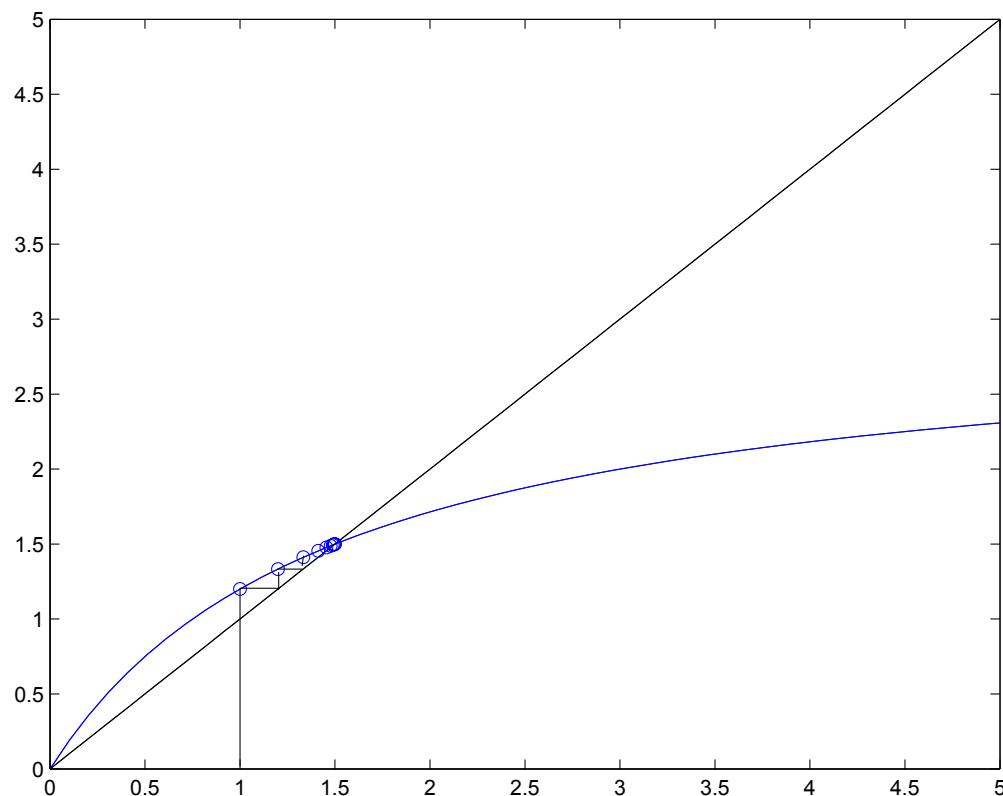
```

>> phi2=@(x) x.*(2./(1+(x./1.5)));
>> phi3=@(x) x.*(2*x./(1+(x./1.5).^2));
>> x=linspace(0,5,50);
>> figure I ght be plotted I pt t-be plotted
>> plot(x,phi2(x), 'b', x, x, 'k');
>> [p2,res2,niter2]=fixpoint(phi2,1,1e-6,1000);
>> figure
>> plot(x,phi3(x), 'b', x, x, 'k'); A for which to tollerance find fixed pt max # steps
>> [p3,res3,niter3]=fixpoint(phi3,1,1e-6,1000);

```

We find the stationary points  $\alpha_2 = 1.5$  and  $\alpha_3 = 3.9271$ .

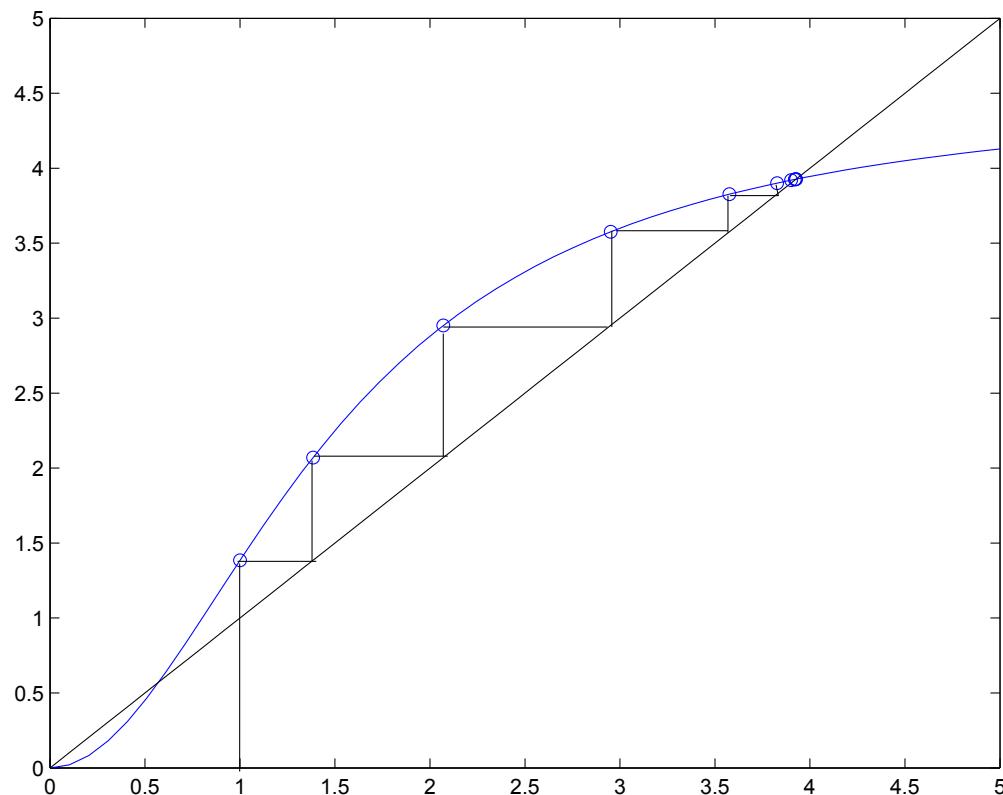
Function  $\phi_2(x)$ :



Careful to  
what first index  
is! (0 or 1!)

$x^{(0)} = 1.0000,$ $x^{(1)} = 1.2000,$ $x^{(2)} = 1.3333,$ $x^{(3)} = 1.4118,$	$ x^{(0)} - \alpha_2  = 0.5000$ $ x^{(1)} - \alpha_2  = 0.3000$ $ x^{(2)} - \alpha_2  = 0.1667$ $ x^{(3)} - \alpha_2  = 0.0882$
--	--

Function  $\phi_3(x)$ :



$$\begin{aligned} x^{(0)} &= 1.0000, \\ x^{(1)} &= 1.3846, \\ x^{(2)} &= 2.0703, \\ x^{(3)} &= 2.9509, \end{aligned}$$

$$\begin{aligned} |x^{(0)} - \alpha_3| &= 2.9271 \\ |x^{(1)} - \alpha_3| &= 2.5424 \\ |x^{(2)} - \alpha_3| &= 1.8568 \\ |x^{(3)} - \alpha_3| &= 0.9761 \end{aligned}$$

**Example. 6 (cont)** We have used the fixed point algorithms using the two functions  $\phi_1$  and  $\phi_2$  with initial value  $x^{(0)} = 0.7$ . Remember that both have the same fixed point  $\alpha$ .

easier, but  $\rightarrow x = \phi_1(x) = 1 - \sin(2x)$  because  $|\phi'_1(x)| > 1$   
 not working!

more difficult  $\rightarrow x = \phi_2(x) = \frac{1}{2} \arcsin(1 - x), \quad 0 \leq x \leq 1$   
 but works

! FIRST THING  
 TO CHECK  
 is I deriv!

```

>> [p1,res1,niter1]=fixpoint(phi1,0.7,1e-8,1000);
>> [p2,res2,niter2]=fixpoint(phi2,0.7,1e-8,1000);
  
```

The fixed point algorithm with the first function does not converge, while with the second one it converges to  $\alpha = 0.352288459558650$  in 44 iterations. Indeed,  $\phi'_1(\alpha) = -1.5237713$  and  $\phi'_2(\alpha) = -0.65626645$ .

# More about the Newton method.

↑ Take the best of Newton method and extend it in the whole fixed point framework

Very good convergence property (QUADRATIC!)

The Newton method is a fixed point method:  $x^{(k+1)} = \phi(x^{(k)})$  for the function

$$\boxed{\phi(x) = x - \frac{f(x)}{f'(x)}}.$$

Let  $\alpha$  be a zero of  $f$ , i.e. such that  $f(\alpha) = 0$ . Note that  $\underbrace{\phi'(\alpha) = 0}$ , when  $f'(\alpha) \neq 0$ . Indeed,

$\phi'(x) = 1 - \frac{[f'(x)]^2 - f(x)f''(x)}{[f'(x)]^2}.$

↓  
 if  $f'(\alpha) = 0$   
 UNLUCKY CASE

↓  
 $\phi'(\alpha) = 1 - \frac{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2 - \overbrace{f(\alpha)}^{\text{if } f(\alpha)=0} \overbrace{f''(\alpha)}^{\text{if } f''(\alpha)\neq 0}}{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2}$   
 $= 1 - \frac{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2}{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2} = 0$

so good because we obtain fixed point  
 with II procedure,  
 (it's not iteration),  
 and it is consistent with  
 Newton method (II order)

use it one algorithm  
 ↓  
 merge fixed point  
 taking the property of  
 Newton method

↓  
 QUADRATIC CONVERGENCE  
 to the fixed point  
 which is "justified" by  
 obtaining  $\phi'(\alpha) = 0$

It is important that  $\phi''(\alpha) \neq 0$

If this fails to happen  
failure is of Newton  
method (it is going to  
be linear)



Solve by change a little  
bit the  $\phi$  ft  
(DEFLATION METHOD)

{  
modify Newton scheme  
by making sure that  
(the II derivative  $\neq 0$ )

**Theorem 1.** If  $f$  is twice differentiable,  $f(\alpha) = 0$  and  $f'(\alpha) \neq 0$ , then there exists  $\delta > 0$  such that, if  $|x^{(0)} - \alpha| \leq \delta$ , the sequence defined by the Newton method converges to  $\alpha$ .

Moreover, the convergence is quadratic; more precisely

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{(x^{(k)} - \alpha)^2} = \frac{\overbrace{f''(\alpha)}^{\text{this represents } \phi''(\alpha)}}{2\overbrace{f'(\alpha)}^{\text{the second derivative}}}$$

*Proof.* The property of convergence comes from the Proposition 2, while the quadratic convergence is a consequence of the Proposition 3, because

$$f'(\alpha) = 0 \text{ and } \frac{\phi''(\alpha)}{2} = \frac{f''(\alpha)}{2f'(\alpha)}.$$

$\Rightarrow$  We prove quadratic convergence thanks to inferring from the functions, without the second derivative of  $\phi(x) = x - \frac{f(x)}{f'(x)}$

**Definition 2.** Let  $\alpha$  be a zero of  $f$ .  $\alpha$  is said to have multiplicity  $m$ ,  $m \in \mathbb{N}$ , if  $f(\alpha) = \dots = f^{(m-1)}(\alpha) = 0$  and  $f^{(m)}(\alpha) \neq 0$ .

A zero that has multiplicity  $m = 1$  is called simple zero.

✓ !  $\frac{x^{(k+1)} - \alpha}{x^{(k)} - \alpha}$  IS UNCOMPUTABLE, because

every ft or pty related to  $\alpha$  is  
UNCOMPUTABLE, since  $\alpha$  is NOT known

**Remark 2.** If  $f'(\alpha) = 0$ , the convergence of the Newton method is linear, not quadratic. We can use the modified Newton method:

$$x^{(k+1)} = x^{(k)} - m \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k = 0, 1, 2, \dots \quad (7)$$

where  $m$  is the multiplicity of  $\alpha$ .

at least = 2 since  $f'(\alpha) = 0$

If the multiplicity  $m$  of  $\alpha$  is unknown, there are other methods, the *adaptive methods*, which can recover the quadratic order of convergence.

Amplifying  $f$  by  $m$  or reduce  $f'$  by  $m$   
 $\Rightarrow$  IMP: you can try to avoid updating  $f'$   
 at each iteration and delay the update  
 of the most expensive computational part  
 not free!  
 you may  
 ↳ ADAPTIVITIES ↳ ACCURACY

# A stopping criterion for Newton

When to stop the Newton method? A good stopping criterion is the **control of the increment** : the iterations is completed when

$$|x^{(k+1)} - x^{(k)}| < \epsilon \quad (8)$$

where  $\epsilon$  is a fixed tolerance.

Indeed, if we denote  $e^{(k)} = \alpha - x^{(k)}$  is the error of the iteration  $k$ , we have

$$e^{(k+1)} = \alpha - x^{(k+1)} = \phi(\alpha) - \phi(x^{(k)}) \stackrel{\text{1 order Taylor expansion}}{=} \phi'(\xi^{(k)}) e^{(k)},$$

we bind error  
 at following  
 iteration

C between  $x^{(k)}$  and  $\alpha$

where  $\xi^{(k)}$  is between  $x^{(k)}$  and  $\alpha$ , and

$$x^{(k+1)} - x^{(k)} \stackrel{?}{=} \alpha - x^{(k)} - \alpha + x^{(k+1)} = e^{(k)} - e^{(k+1)} = \left(1 - \phi'(\xi^{(k)})\right) e^{(k)}. \quad (9)$$

Assuming that if  $k$  is large enough, we have  $\phi'(\xi^{(k)}) \approx \phi'(\alpha)$  and knowing that the Newton method for  $\phi'(\alpha) = 0$ , if  $\alpha$  is a simple zero, we find the estimation

$$|e^{(k)}| \approx |x^{(k+1)} - x^{(k)}|.$$

In this case, one can be goodly approx or bad approx by difference of subsequent iteration

The error that we commit when we adopt the criterion (8) is smaller than the fixed tolerance.

↓  
 NEWTON:  
 nice jfns  
 in perf. mean  
 ce (probabilic  
 convergence)  
 and nice re-  
 met for our  
 estimation

# Stopping criteria: the general case

In general, for all discussed methods, we can use two different stopping criteria: the iterations is completed when

$$|x^{(k+1)} - x^{(k)}| < \epsilon \quad (\text{control of the increment}),$$

or

$$|f(x^{(k)})| < \epsilon \quad (\text{control of the residual}),$$

where  $\epsilon$  is a fixed tolerance.

$f(x^{(n)})$  is the RESIDUAL, because  
 $f(x) = 0$ , so there we should  
 guarantee the consistency of the  
 methodology

Using fixed point iterations we obtain the following estimation:

$$e^{(k)} \approx \frac{1}{(1 - \phi'(\alpha))} (x^{(k+1)} - x^{(k)}).$$

*as a fit of  $\phi'$*

*for  $k$  large enough*

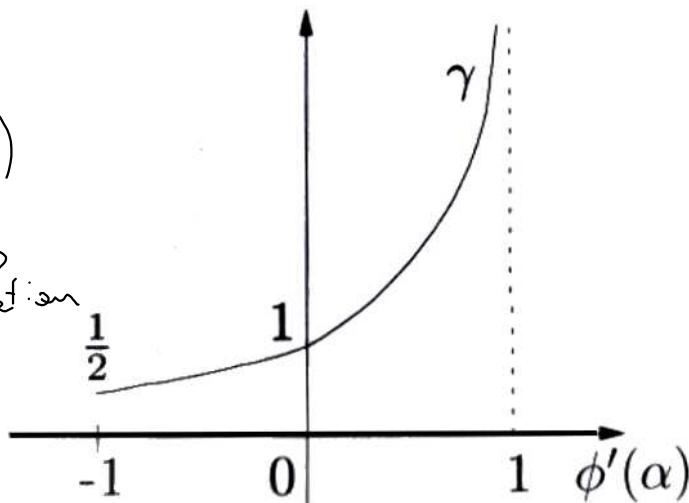
We can plot the graph of  $\frac{1}{(1 - \phi'(\alpha))}$  and comment on the relevance of the stopping criterium based on the increment:

- if  $\phi'(\alpha)$  is near to 1 the test is not satisfactory (esymptotic behavior)
- for methods of order 2 ( $\phi'(\alpha) = 0$ ), the criterium is optimal,  $\leftarrow$  Newton  $\Rightarrow$  best approximation
- if  $-1 < \phi'(\alpha) < 0$  the criterium is still all right.

*error is  $\frac{1}{2}$  of increment*

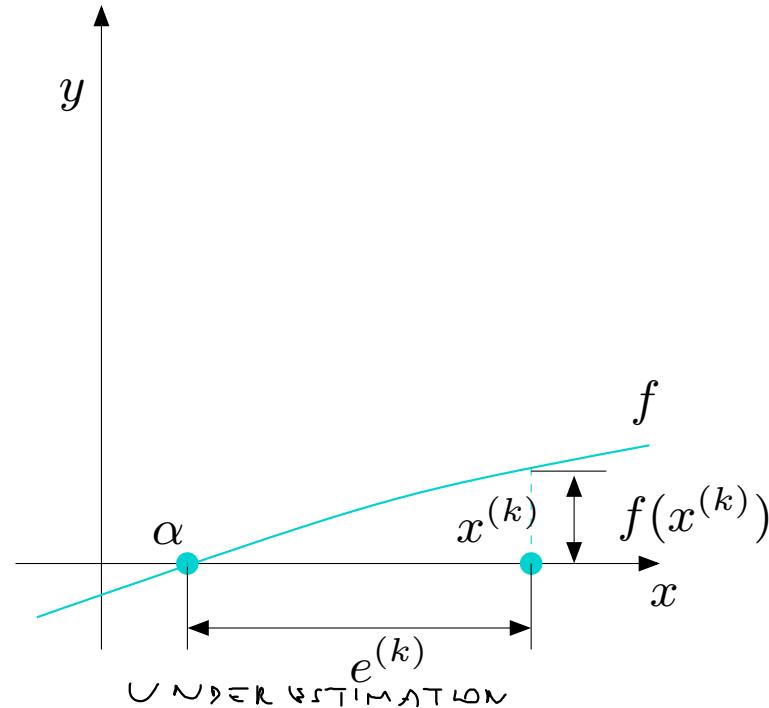
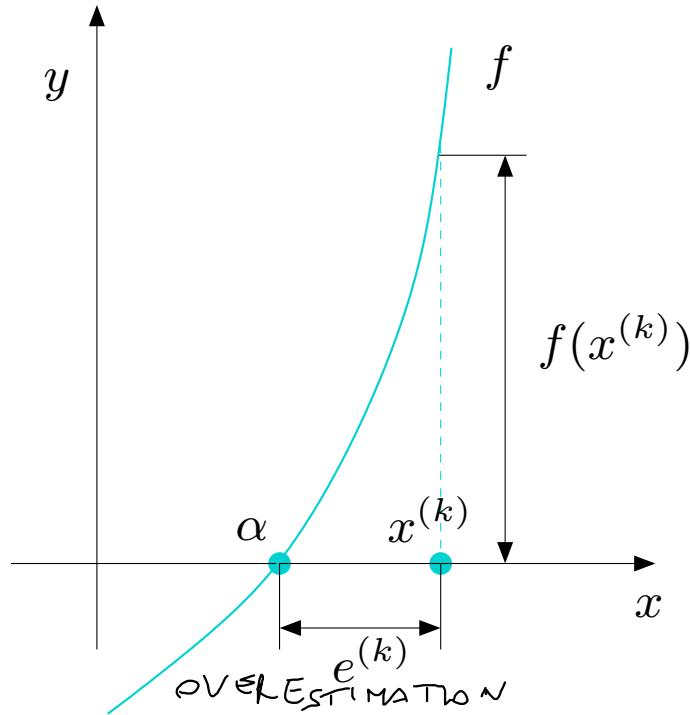
$\Rightarrow$  good, because it means we are in a

*SAFE ZONE* (good when things are worse than reality)



# Stopping Criteria

The stopping criterium based on the control on the residual  $|f(x^{(k)})| < \epsilon$  is satisfactory only if  $|f'| \simeq 1$  near the root  $\alpha$ . Otherwise it is too strong (if  $|f'| \gg 1$ ) or too weak (if  $|f'| \ll 1$ ):



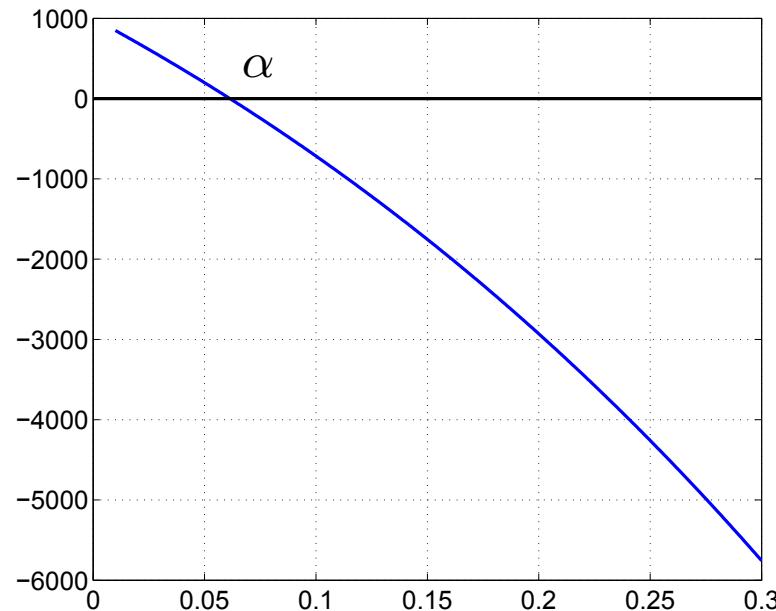
Two cases where the residual is a bad estimator of the error:  $|f'(x)| \gg 1$  (left),  $|f'(x)| \ll 1$  (right) with  $x$  near to  $\alpha$

my DEA... if divide  
 $f(x^{(n)})$  by  $f'(x^{(n)})$   
 non-iteration

# Applications

**Example. I (suite)** We draw the graph of  $f(I) = M - v \frac{1+I}{I} [(1+I)^n - 1]$  on the interval  $[0.01:0.3]$  with  $M = 6000$ ,  $v = 1000$  and  $n = 5$ :

```
>> f=@(x) 6000-1000*(1+x).*((1+x).^5 - 1)./x;
>> I = [0.01:0.001:0.3];
>> grid on;plot(I,feval(f,x));
```



The root of  $f$  is between 0.05 and 0.1.

We can apply the bisection method on the interval  $[0.05, 0.1]$  with a tolerance  $10^{-5}$

```
>> [zero,res,niter]=bisection(f,0.05,0.1,1e-5,1000);
```

The approximate solution after 12 iterations is  $\bar{x} = 0.061407470703125$ .

We can apply the Newton method with initial guess  $x^{(0)} = 0.05$

```
>> df=@(x) 1000*((1+x).^5.* (1-5*x) - 1)./(x.^2);
>> [zero,res,niter]=newton(f,df,.05,1e-5,1000);
```

The result is approximately the same, but we need only 3 iterations

The interest rate is 6.14%.

*! If you use Newton Method : You have to give also the I derivative of the nonlinear equation*

**Example. 2 (cont)** We would like to plot the angle  $\theta$  as function of  $\omega$  for  $0 \leq \omega \leq \pi$  with  $a_1 = 10\text{ cm}$ ,  $a_2 = 13\text{ cm}$ ,  $a_3 = 8\text{ cm}$ ,  $a_4 = 10\text{ cm}$ .  
 For each  $\omega$ , we have to solve the nonlinear problem

$$f(\theta) = \frac{a_1}{a_2} \cos(\omega) - \frac{a_1}{a_4} \cos(\theta) - \cos(\omega - \theta) + \frac{a_1^2 + a_2^2 - a_3^2 + a_4^2}{2a_2a_4} = 0. \quad (10)$$

To start with, we plot the graph of  $f(\theta)$  for  $\omega = \pi/3$ :

```
>> F = @(x, a1, a2, a3, a4, omega) ...  

    (a1/a2)*cos(omega) - (a1/a4)*cos(x) - cos(omega-x) ...  

    + ( (a1.^2 + a2.^2 - a3.^2 + a4.^2) / (2*a2*a4) );  

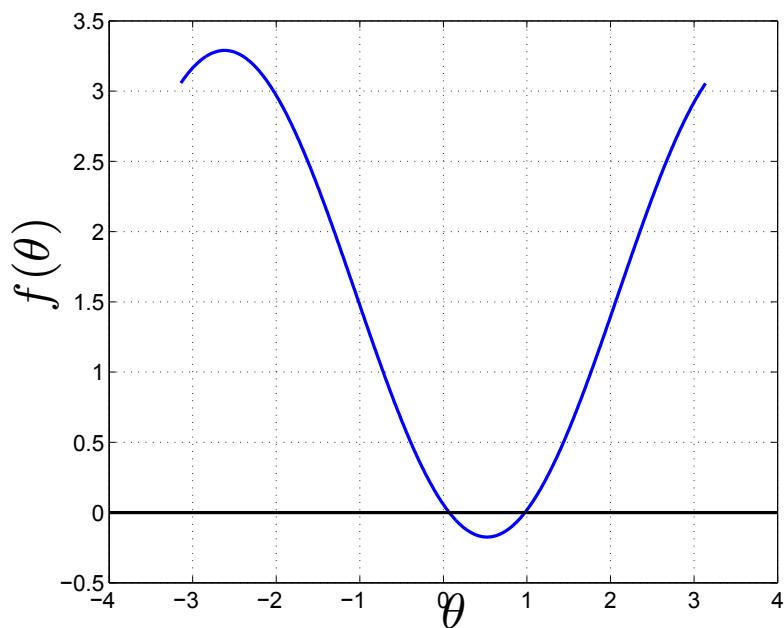
>> a1=10; a2=13; a3=8; a4=10; omega=pi/3;  

>> f = @(x) F(x,a1,a2,a3,a4,omega);  

>> x = [-pi:0.01:pi];  

>> plot( x, f(x) ); grid on;
```

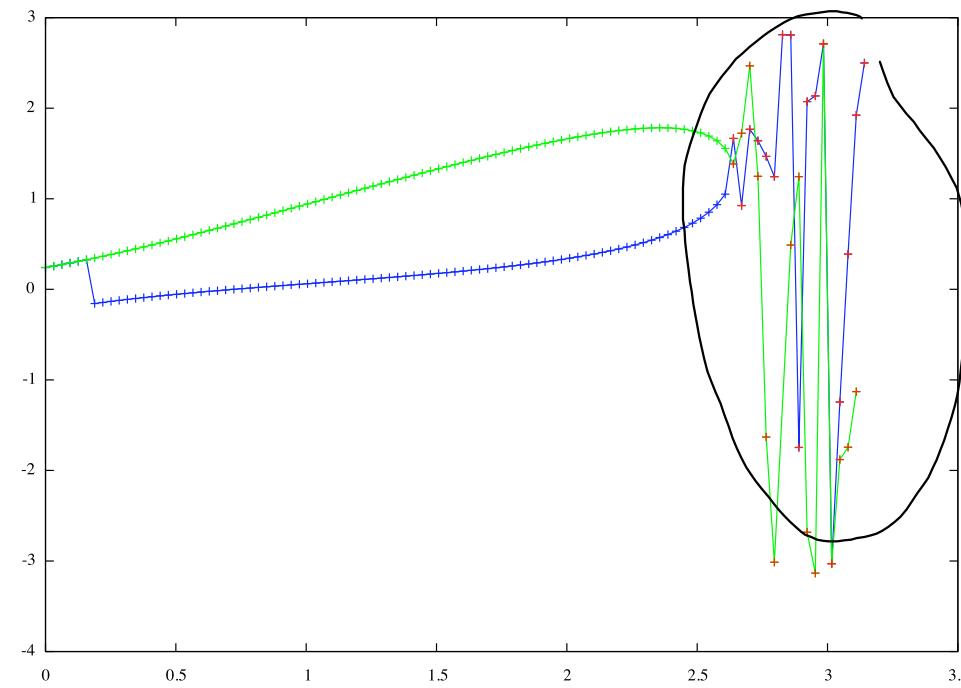
Function  $f$  with  $\omega = \pi/3$ .



We have two roots, which means that we have two possible configurations.

Now we chose 101 different values for  $\omega$ ,  $\omega_k = k \frac{\pi}{100}$ ,  $k = 0, \dots, 100$  and for each of them solve the nonlinear problem (10) with the Newton method. Since we know that we may have two distinct solutions, we give two different initial guesses to our Newton method. For example  $\theta_{01} = -0.1$  and  $\theta_{02} = 2/3\pi$

The following figure shows the solutions as function of  $\omega$ . For  $\omega > 2.6358$ , the Newton algorithm does not converge anymore. In fact, with these values there is no configuration possible



no solution,  
but INCONSISTENT  
BEHAVIOR OF  
MECHANICAL  
SYSTEM  
oscillation on  
the to numerical  
wise : NOT PHYSICAL

Here are the Matlab/Octave commands that we have used:

```
>> n=101; x01=-0.1; x02=2*pi/3; nmax=100;
>> dF = @(x,a1,a2,a3,a4,w) a1/a4*sin(x)-sin(w-x);
>> for k=1:1:n
    omega(k) = (k-1)*pi/100;
    f = @(x) F(x,a1,a2,a3,a4,omega(k));
    df = @(x) dF(x,a1,a2,a3,a4,omega(k));
    [theta1(k),res,niter] = newton(f,df,x01,1e-5,nmax);
    [theta2(k),res,niter] = newton(f,df,x02,1e-5,nmax);
end
>> plot(omega,theta1,'b:',omega,theta2,'g-')
```

**Example. 3(suite)** We consider the carbon dioxide ( $\text{CO}_2$ ), for which  $a = 0.401 \text{ Pa m}^6$  and  $b = 42.7 \cdot 10^{-6} \text{ m}^3$ .

We search the volume occupied by  $N = 1000$  molecules of  $\text{CO}_2$  in temperature  $T = 300 \text{ K}$  and pressure  $p = 3.5 \cdot 10^7 \text{ Pa}$ . We know that the Boltzmann constant is  $k = 1.3806503 \cdot 10^{-23} \text{ Joule K}^{-1}$ .



We draw the graph of the function

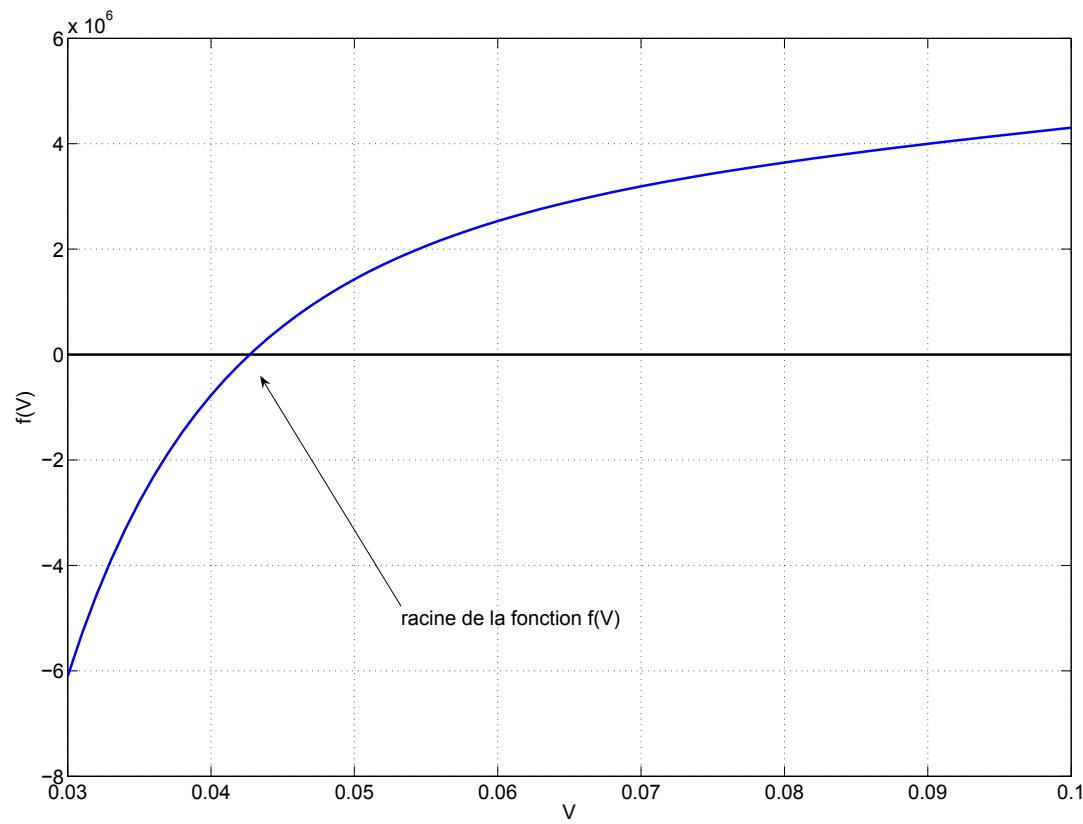
$$f(V) = \left[ p + a \left( \frac{N}{V} \right)^2 \right] (V - Nb) - kNT$$

for  $V > 0$ . We do not consider  $V < 0$  (it does not have physical meaning), because  $V$  is the volume of gas.

We use the commands in Matlab/Octave:

```
>> a=0.401; b=42.7e-6; p=3.5e7; T=300; N=1000; k=1.3806503e-23;
>> f = @(x,p,T,a,b,N,k)(p+a*((N./x).^2)).*(x-N*b)-k*N*T;
>> x=[0.03:0.001:0.1];
>> plot(x,f(x,p,T,a,b,N,k))
>> grid on
```

We obtain the graph of the function  $f(V)$ :



We see that there is a zero for  $0.03 < V < 0.1$ . If we apply the bisection method on the interval  $[0.03, 0.1]$  with a tolerance  $10^{-12}$ :

```
[zero,res,niter]=bisection(f,0.03,0.1,1e-12,1000,p,T,a,b,N,k);
```

then we find, after 36 iterations, the value  $V = 0.0427$ .

If we use the Newton method with the same tolerance, starting from the initial point  $x^{(0)} = 0.03$ ,

```
>> df = @(x,p,T,a,b,N,k) -2*a*N^2/(x^3)*(x-N*b)+(p+a*((N./x).^2));
>> [zero,res,niter]=newton(f,df,0.03,1e-12,1000,p,T,a,b,N,k);
```

then we find the same solution after 6 iterations.

The conclusion is that the volume  $V$  occupied by the gas is  $0.0427 \text{ m}^3$ .

# The rope method

→ way to avoid I derivative  
could be e Jacobian  
 IP in N-dimension  
 $N > 1$

This method is obtained by replacing  $f'(x^{(k)})$  by a fixed  $q$  in the Newton method:

$$x^{(k+1)} = x^{(k)} - \frac{1}{q} f(x^{(k)}), \quad k = 0, 1, 2, \dots \quad (11)$$

We can take, for example,  $q = f'(x^{(0)})$  or  $q = \frac{f(b) - f(a)}{b - a}$ , in the case when we search a zero in the interval  $[a, b]$ .

**Example. 6 (suite)** We apply the rope method and the Newton method to find the zero of  $f$ .

The rope method in the interval  $[-1, 1]$ , with  $x^{(0)} = 0.7$  :

```
>> [zero,res,niter]=chord(f,-1,1,0.7,1e-8,1000)
```

We find the result after 15 iterations.

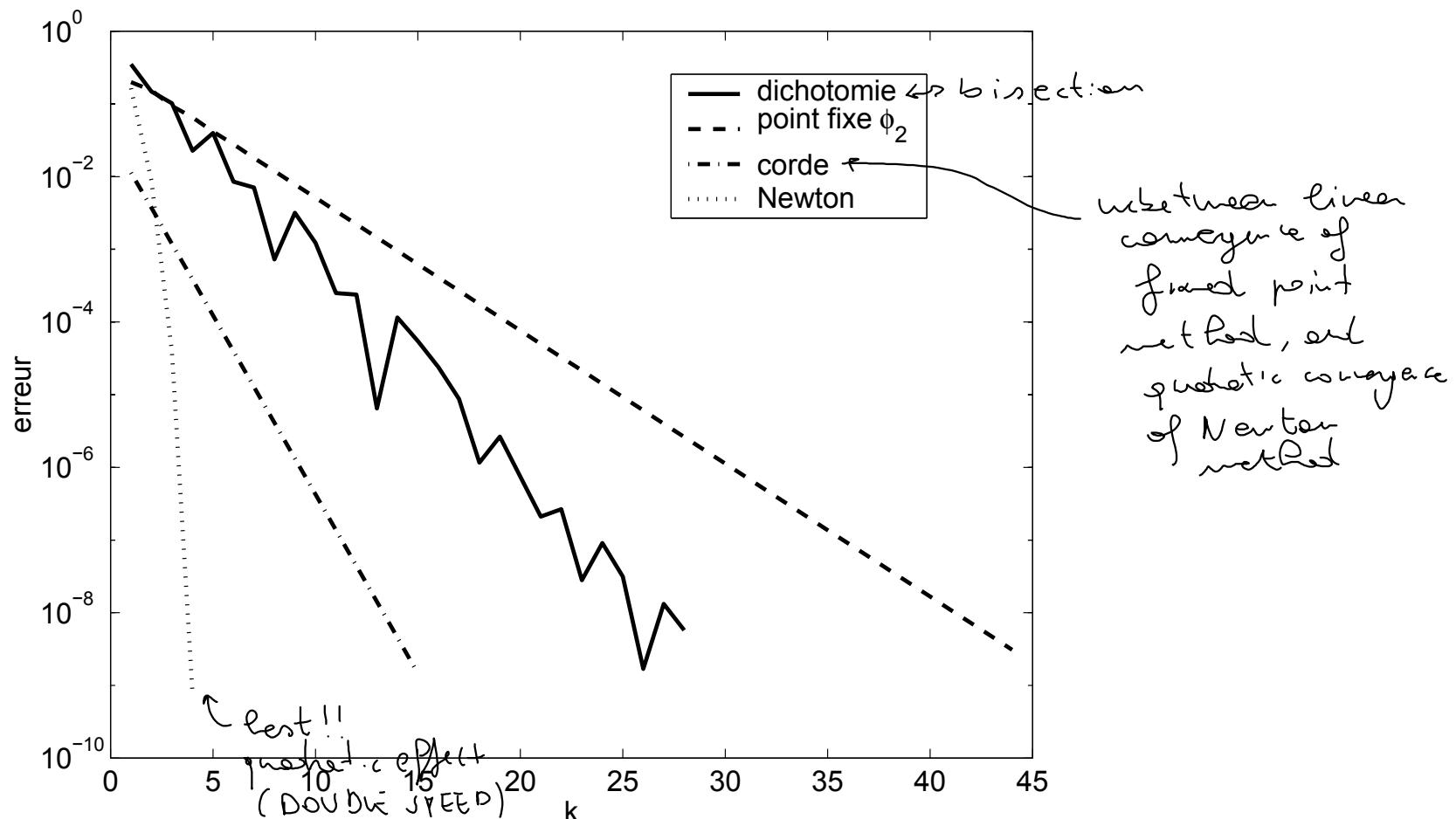
The Newton method with the same  $x^{(0)}$  :

```
>> df = @(x) 2*cos(2*x) + 1;
>> [zero,res,niter]=newton(f,df,0.7,1e-8,1000);
```

We find the result after 5 iterations. Much faster.

Exercice  
question

Values of the errors plotted versus the number of iterations for 4 methods: bisection, fixed point  $\phi_2$ , rope and Newton. There is logarithmic scale on the axis  $y$ .



**Remark 3.** *The rope method is also a fixed point method for*

$$\phi(x) = x - \frac{1}{q}f(x).$$

*So, we have  $\phi'(x) = 1 - \frac{1}{q}f'(x)$  and thanks to the Proposition 2, we obtain that the method converges if the following condition is satisfied:*

$$| 1 - \frac{1}{q}f'(\alpha) | < 1 .$$

we want  $n$  to be ~~at~~  
~~of interpolation~~  $\Rightarrow$  ~~at~~  $n+1$   
~~at less in P^n~~

Recap on interpolation : 1) a set of interpolation points  $\{x_i\}_{i=0}^n$   
 2) a set of basis functions  $\{v_i\}_{i=0}^n$

Interpolation of  $f \in C^0([a, b])$   $x_i \in [a, b] \forall i$

on the space  $V = \text{span}\{v_i\}_{i=0}^n$  is defined as the solution to :

find  $p$  in  $V$  s.t.  $p(x_i) = f(x_i) \forall i$

Given  $\{x_i\}_{i=0}^n$  best choice for  $\{v_i\}_{i=0}^n$  is Lagrange basis

$$\Rightarrow v_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^{n} \frac{(x - x_j)}{(x_i - x_j)}$$

So the interpolation poly becomes

$$p = \sum_{i=0}^n v_i(x) f(x_i) := L^n f$$

Next step was:

Try to find the "best" set of  $\{x_i\}_{i=0}^n$

Remember that  $L^n$  is a smooth poly.  
 Every smooth poly  $\underbrace{\text{last approx}}$

$$\|L^n f - f\|_\infty = \|L^n(f - p) + p - f\|_\infty \leq (\|L^n\|_\infty + 1) \|p - f\|_\infty$$

$$\|L^n\|_\infty \leq \|L^n(\{x_i\}_{i=0}^n)\|_\infty$$

We try to min the ~~min~~ of the absolute values of the basis ft ( $\Leftrightarrow$  Lebesgue

$$\lambda := \sum_{i=0}^n |v_i|$$

$\hookrightarrow$  Lebesgue function

For Chebyshev points:

$$\frac{2}{\pi} \log(n+1) - c \leq \|\lambda\|_\infty \leq \frac{2}{\pi} \log(n+1) + 1$$

! True only if the function is continuous

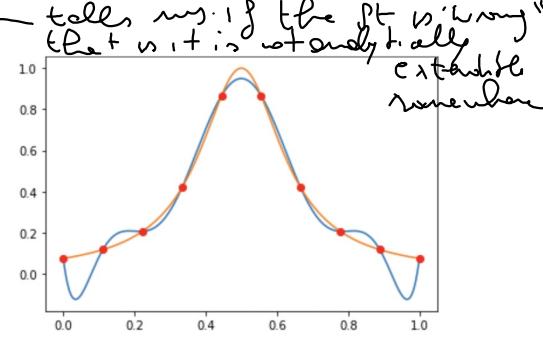
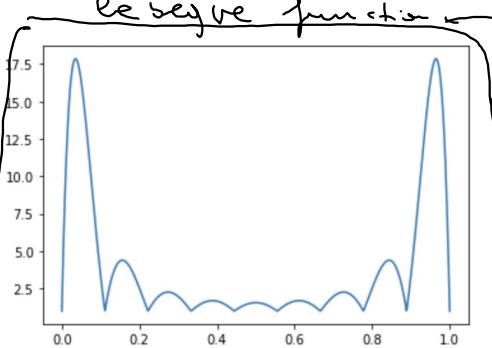
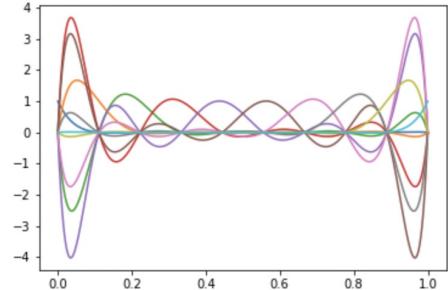
If you add more regularity, you get

Better things (Result for analytically  
extendible f) with continuous derivatives  
up to (the n-th derivative)

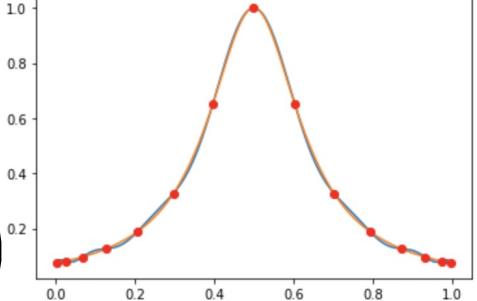
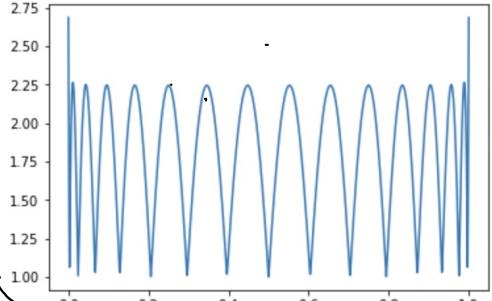
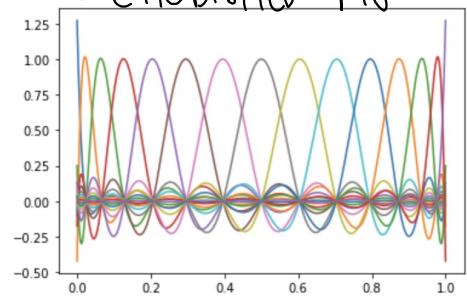
! INTERPOL IS EXACT AT POINT OF  
INTERPOL

=> try to get exact approx every where  
in some way

USING EQUISPACED PTS



USING CHEBYSHEV PTS



If interpolation is not good, what can we do in  $C^0([a,b])$ ?

Ans: Weierstrass Approximation Theorem

$\forall f \in C^0([a,b])$ ,  $\forall \epsilon > 0$ ,  $\exists p \in P^n$ , s.t.  $\|f-p\|_{\infty} < \epsilon$

~~No matter what  $f$ , you can always find a  $p$  such that  $p$  approx  $f$~~

Proof: Let  $B_n$  be a sequence of linear operators s.t.

- 1)  $B_n$  positive linear operators  $C([a,b]) \rightarrow P^n([a,b])$
- 2)  $B_n q \rightarrow q$  pointwise  $\forall q \in P^2([a,b])$  that is only for quadrat. polys

Then:

$B_n f$  converges uniformly to  $f$   $\forall f \in C^0([a,b])$

N.B.:  $B_n$  is "positive" if  $q(x) \geq 0 \Rightarrow (B_n q)(x) \geq 0$  in  $[a,b]$   
(if  $q(x) \leq 0$  it remains negative if  $B$  is positive)

Part 1: .  $\forall f \in C^0([a,b])$ ,  $\forall x_0 \in (a,b)$ :

- construct  $q^\pm \in P^2$  s.t.  $q^-(x_0) < f(x_0) < q^+(x_0)$
- use  $B_n q^\pm \rightarrow q^\pm$
- use  $B_n(f - q^-)(x_0) > 0$        $B_n(q^+ - f)(x_0) > 0$

$$f \in C^0([a, b]) \Rightarrow \forall \varepsilon > 0, \exists \delta > 0$$

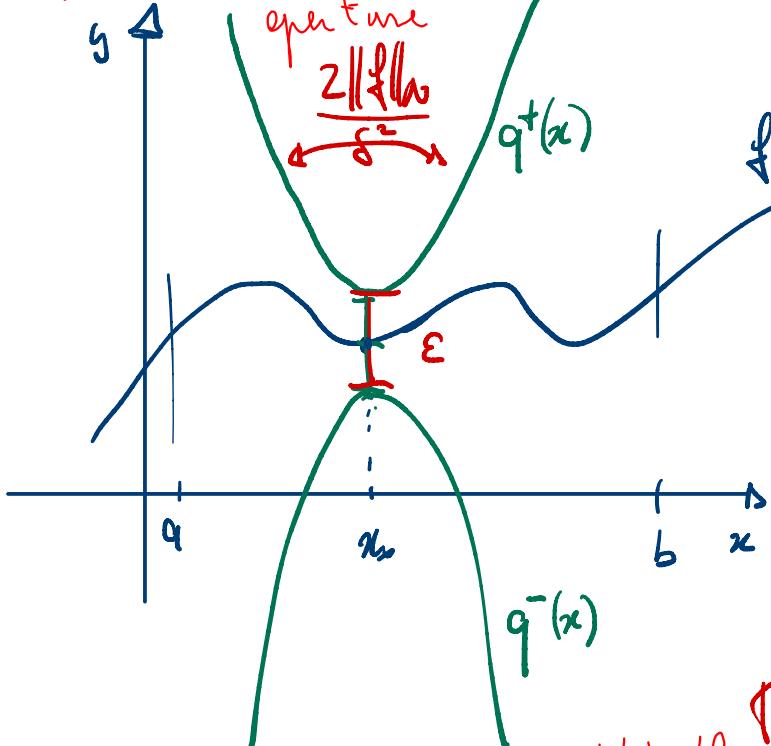
} def of continuity!

$$|x_1 - x_2| \leq \delta \rightarrow |f(x_1) - f(x_2)| \leq \varepsilon$$

Define

$$q^\pm := f(x_0) \pm \left( \frac{\varepsilon}{2} + \frac{2\|f\|_\infty}{\delta^2} (x - x_0)^2 \right)$$

Busting up in this way guarantees that they are always exact below &



for varying  $x_0$

$$q^\pm(x) = a^\pm x + b^\pm x + c^\pm$$

$a^\pm, b^\pm, c^\pm$  depend on

$x_0, \|f\|_\infty, \delta, \varepsilon$

Take  $M = \max_{x_0 \in [a, b]} (|a^\pm|, |b^\pm|, |c^\pm|)$

I take the max over the entire interval

$M$  depends on  $\|f\|_\infty, \delta, \varepsilon$ , but Not on  $x_0$

Choose  $N$  large enough s.t.

( $\exists N$  s.t. this is possible)  
Since we know that  $x_n$  converges uniformly for all of the two (HYPOTHESIS)

$$\|B_n x_n^i - x_i^i\|_\infty \leq \frac{\varepsilon}{6M}$$

$\forall n \geq N$   
 $\forall i = 0, 1, 2$

$$\Rightarrow \forall x_0 \quad * \|B_n q^\pm - q^\pm\| \leq \frac{\varepsilon}{2}$$

$M$  controls every coeff  
 $\Rightarrow$  this allows to say that  $\frac{\varepsilon}{6M} \leq \frac{\varepsilon}{2}$  (?)

$$\text{in } x_0 \quad f(x_0) - \varepsilon \leq q^-(x_0) - \frac{\varepsilon}{2} \leq B_n q^-(x) \leq B_n f(x_0)$$

And repeat for the reversed inequality:

$$\text{if } x_0, \exists N \text{ s.t., } \forall n \geq N \quad B_n f \leq B_n q^+ \leq q^+ + \frac{\varepsilon}{2} \leq f(x_0) + \varepsilon$$

$$\|B_n f - f\| \leq \varepsilon$$

How to define  $B_n$ ? Let  $(a, b) = [0, 1]$

Trick : Consider  $\underbrace{\left( (1-x) + x \right)^n}_{\stackrel{=1}{\text{defines bin. functions}}} = \sum_{i=0}^n \binom{n}{i} x^i (1-x)^{n-i}$  poly's of  $n$

$$= \sum_{i=0}^n v_i(x)$$

$v_i(x) = \binom{n}{i} x^i (1-x)^{n-i}$

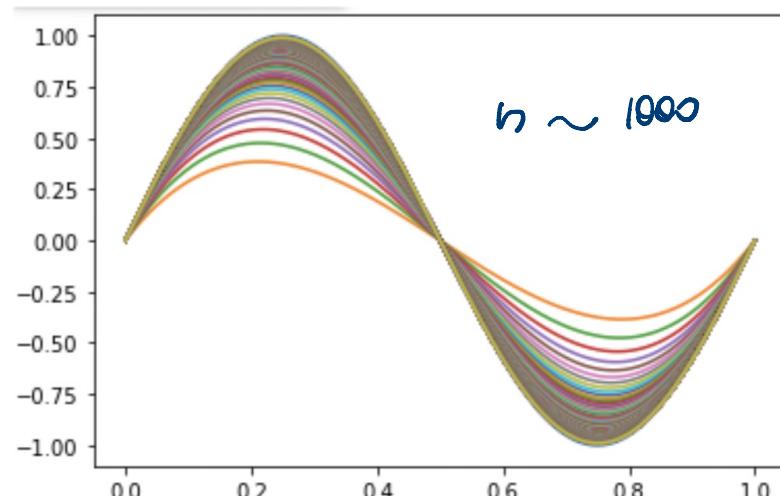
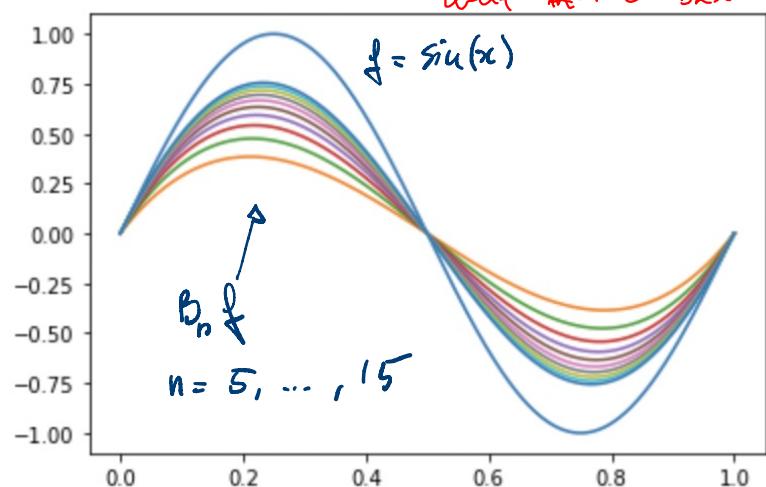
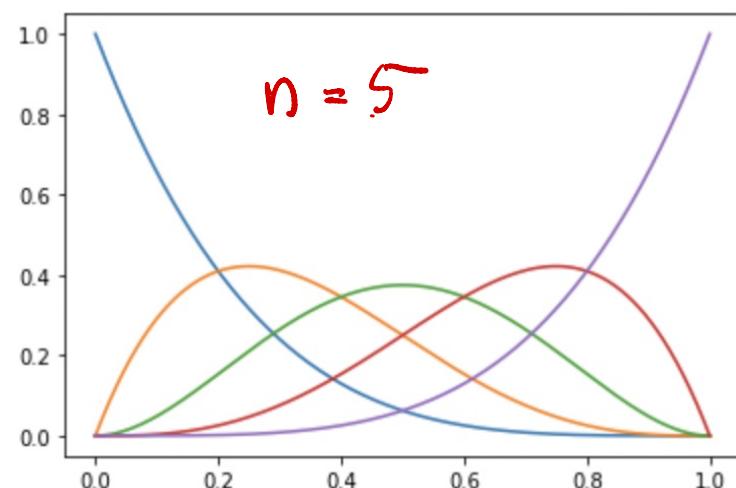
Bernsteine polynomials.

$B_n f := \sum_{i=0}^n v_i(x) f\left(\frac{i}{n}\right)$

$\rightarrow$  but always positive in  $[0, 1]$   
 or also from 0 in middle  
 interval to 0 at ends  
 because  $v_i$  positive  
 Difference between  
 interpolation and this  
 $f \geq 0 \Rightarrow B_n f \geq 0$   
 you cannot  
 say that  
 the  $v_i$  are  
 interpolating  
 at any pt but  
 0 and 1

$B_n x = x \quad B_n 1 = 1$

$B_n x^2 = \left(\frac{n-1}{n}\right)x^2 + \frac{1}{n}x \neq x^2$   
 for large  $n$   
 but  $n \rightarrow \infty \Rightarrow B_n x^2 \rightarrow x^2$



We don't  
 have a quick  
 convergence

## Lecture 2 PROJECTION BASE APPROX METHOD 10. Nov. 2020

1) Approximate and measure using  $L^2$  norm (vs  $L^\infty$ ) derived from a scalar product

→ Scalar product :  $u, v \in L^2([0,1])$ ,

$$(u, v) := \int_0^1 u \cdot v \, dx \quad \|a+b\|^2 = (a+b, a+b) = \|a\|^2 + \|b\|^2 + 2(a, b)$$

$$\leadsto \|u\|_{L^2([0,1])}^2 \equiv \|u\|^2 := (u, u) = \int_0^1 |u|^2 \, dx$$

let's assume we have  $V := \text{span}\{v_i\}_{i=0}^n$ , and  $\|\cdot\|$ .

The best approximation of  $u \in L^2([0,1])$  in  $V$  satisfies:

$$p \in V \text{ s.t. } (u-p, v) = 0 \quad \forall v \in V$$

$\curvearrowleft (u-p, v) = 0 \quad \forall v \in V \iff p \text{ is. B.A. of } u \text{ in } V$

Proof: Assume  $p$  is best approx  $\Rightarrow (u-p, v) = 0 \quad \forall v \in V$   
by def of B.A.

$$\|u-p\|^2 \leq \|u-p + tq\|^2 \quad \forall t > 0, \forall q \in V$$

$$\|u-p + tq\|^2 - \|u-p\|^2 \geq 0 \quad \begin{matrix} \text{if } v \text{ in any direction from} \\ \text{get larger distance from } p \text{ to } u \end{matrix}$$

$$(a+b)^2 - (a-b)^2 = 4(a, b)$$

$$\left\| \underbrace{(u-p + \frac{t}{2}q)}_{a+b} + \left( \frac{t}{2}q \right) \right\|^2 - \|u-p\|^2 \geq 0$$

$$a = u-p + \frac{t}{2}q$$

$$b = \frac{t}{2}q$$

$$2 \not\perp (u-p + \frac{t}{2}q, \frac{t}{2}q) \geq 0$$

Linearity of scalar prod in 1 argument

$$0 \leq 2(u-p, q)t + t^2(q, q)$$

We have a symmetry for:

$$q \rightarrow -q$$

$$-\frac{t}{2}\|q\|^2 \leq (u-p, q) \leq \frac{t}{2}\|q\|^2$$

$$\forall t > 0$$

$$\Rightarrow (u-p, q) = 0$$

$$\forall q \in V$$

Assume  $\star$   $(u - p, v) = 0 \quad \forall v \in V \Rightarrow p$  is B.A.

$$\|u - v\|^2 = \|u - p + p - v\|^2 = \|u - p\|^2 + \|p - v\|^2 + 2(u - p, p - v)$$

$$\Rightarrow \|u - p\|^2 \leq \|u - v\|^2 \quad \forall v \in V$$

that is  $p \in$  B.A. of  $\frac{u-p}{\|u-p\|}$

by  $\star$   
and the fact  
that  $u-p \perp p$

B.A. <sup>(weak)</sup> Find  $p \in V$  s.t.

$$\text{1) } (p, v) = (u, v) \quad \forall v \in V$$

given  
Write  $p(x) = \sum p^j v_j(x)$

$$\Rightarrow 2) \quad (p^j v_j, v_i) = (u, v_i)$$

MASS  
MATRIX

$$\underline{\underline{M}} \underline{p} = \underline{U}$$

$$M_{ij} p^j = U_i$$

$$M_{ij} = (v_i, v_j)$$

$$U_i = (u, v_i)$$

computation  $\Rightarrow$  solve the prob  
 $p^j = \underline{\underline{M}}^{-1} \underline{U}_i$

Example :

$$V = \text{span} \left\{ \text{pow}(x, i) \right\}_{i=0}^n$$

$$p = \sum_j p^j x^{(j)} \text{ power basis}$$

then  $M_{ij} = \int_0^1 x^i \cdot x^j dx = \int_0^1 x^{(i+j)} dx = \frac{1}{i+j+1}$

Here <sup>"called"</sup> Hilbert Matrix : cond. ( $H$ )  $\approx \frac{(1+\sqrt{2})^{4n}}{\sqrt{n}}$

it is Tenu-  
ely condi-  
tioned  
matrix  
If  $i=j$  very  
large  $\Rightarrow$  column  
should be same  
 $(i+j+1) \times i = 1+2+..$

What's the best possible "H" matrix?

Identity :  $M_{ij} = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$

Def B.A. in a space with norm defined  
for scalar product, select basis functions

and try to compute mass matrix



B.A. in a finite dim space (spanned by  
some basis  $\{f_i\}$ ) what you obtain is

$$\underline{M} \underline{p} = \underline{U}$$

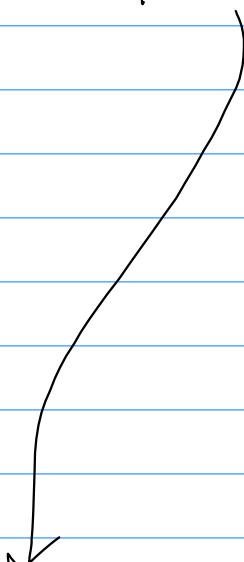
$\Rightarrow$  Best pt is the one in which you  
don't have to write any matrix  $\Rightarrow$  possible if  $M=11$

$$\Rightarrow p_i = (u, v_i)$$



Simplest procedure

GRAM SCHMIDT PROCEDURE



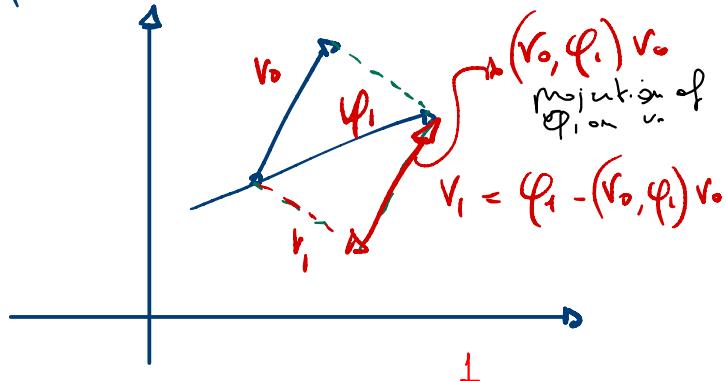
Let's try to find an orthonormal basis

Set  $v_0 = 1 \rightarrow \left( \int_0^1 (1)^2 dx \right)^{\frac{1}{2}} = 1$

$$q_{i+1} = x v_i - \sum_{j=0}^{i-1} (x v_i, v_j) v_j$$

and normalize

$$v_{i+1} = \frac{q_{i+1}}{\|q_{i+1}\|}$$



Possible problems:

$$\begin{aligned} \|q_i\|^2 &= \|x v_0\|^2 + (x v_0, v_0)^2 \|v_0\|^2 - 2(x v_0, v_0)^2 \\ &= \|x v_0\|^2 - (x v_0, v_0)^2 \end{aligned}$$

This can be very small!!

We have another solution which does not ask same norm:  $\vartheta_{ij} \neq \delta_{ij}$   
but it is diagonal:

Bonnet recursion:

$$P_0 = 1 \quad P_1 = x$$

$$(n+1) P_{n+1}(x) = (2n+1) x P_n(x) - n P_{n-1}(x)$$

Are orthogonal in  $[-1, 1]$ , and

$$P_n(1) = 1 \quad \forall n$$

dim (?)  
(start from fact that  $P_0, P_1$  are OK)

"The only" thing which is left to do is to compute

$$\int_0^1 u P_i dx$$

Approximate integrals.

Between functions you def the error as  $\| P(x) - u(x) \|$  where  $\| \cdot \|$  is  
the norm used.

## Lecture 13

12.11.2020

Recap of Lec 12: Given  $u \in L^2([0,1])$ , find  $\rho \in V = \text{span}\{v_i\}_{i=0}^n$   
 s.t.  $(\rho, v) = (u, v) \quad \forall v \in V$

$$\Rightarrow \sum_j (\rho^T v_j, v_i) = (u, v_i) \quad i = 0, \dots, n$$

rewritable in a system

$$\underline{\underline{\rho}} = \underline{\underline{v}}$$

$$M_{ij} := (v_i, v_j)$$

$$v_i = (u, v_i)$$

$$(a, b) := \int_0^1 a \cdot b \, dx$$

Induced or norm

$$\Rightarrow \|a\|^2 = (a, a) \quad \Rightarrow \|a\| = \left( \int_a^b a^2 \, dx \right)^{\frac{1}{2}}$$

Today: Approximate

$$\int_a^b u(x) \, dx$$

For now:  $u \in C^0([a,b])$

Rough approx schemes for  $C^0([a,b])$  functions

For example:

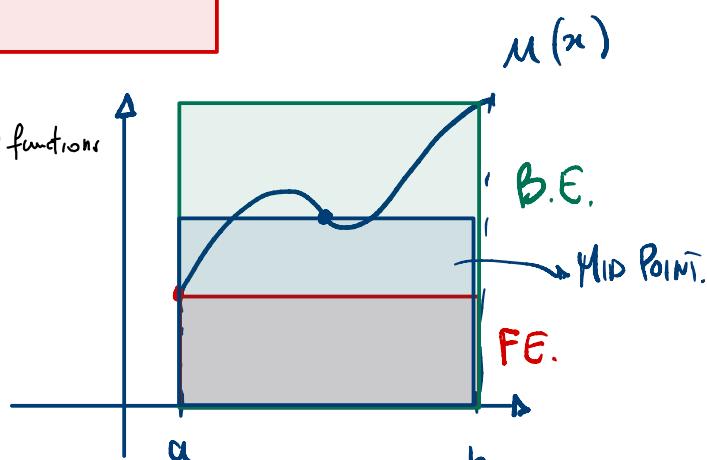
$\begin{cases} \text{FE} & \text{(Forward Euler)} \\ \text{BE} & \text{(Backward Euler)} \\ \text{MIDPOINT} & \end{cases}$
--

We have, for each scheme

$$\text{FE: } \int_a^b u(x) \, dx \approx u(a)(b-a)$$

$$\text{BE: } \int_a^b u(x) \, dx \approx u(b)(b-a)$$

$$\text{MIDPOINT: } \int_a^b u(x) \, dx \approx u\left(\frac{a+b}{2}\right)(b-a)$$



you don't use just one of them, but you apply all three, one for each special case (such as particular sub intervals of  $[a,b]$ )

Interpolatory Quadrature Rules:

$$I(u) := \int_a^b u(x) \, dx$$

$$I_n(u) := \sum_{i=0}^n u(q_i) w_i$$

↑ weights  
quadrature points

$$\approx \int_a^b u(x) \, dx = I(u)$$

if you choose correctly both  $q_i$  and  $w_i$

If we want to integrate we can choose  $\{q_i\}_{i=0}^n$  and  $\{w_i\}_{i=0}^n$  2(n+1) unknowns

Most popular choice: Given  $\{q_i\}_{i=0}^n$ , construct  
interpolatory quadrature rule  $L_u^n$  and integrate (exactly)  $L_u^n$

Lagrange interpolation of  $u$  function

$$I(u) = \int_0^1 (L_u^n)(x) dx$$

where  $L_u^n = \sum_{i=0}^n u(q_i) e_i(x)$  choose the weights  $e_i = \prod_{j=0, j \neq i}^n \frac{(x - q_j)}{(q_i - q_j)}$

$$\Rightarrow I(u) = \int_0^1 \sum_{i=0}^n u(q_i) e_i(x) dx = \sum_{i=0}^n u(q_i) \int_0^1 e_i(x) dx$$

$$= \sum_{i=0}^n u(q_i) w_i$$

number independent in the function  $f$  in trying to integrate

FE, B.E., Mid Point are all interpolatory quadrature rules with  $P^0$  polynomial approximation  
or that is with only 1 quadrature pt

### Trees 1

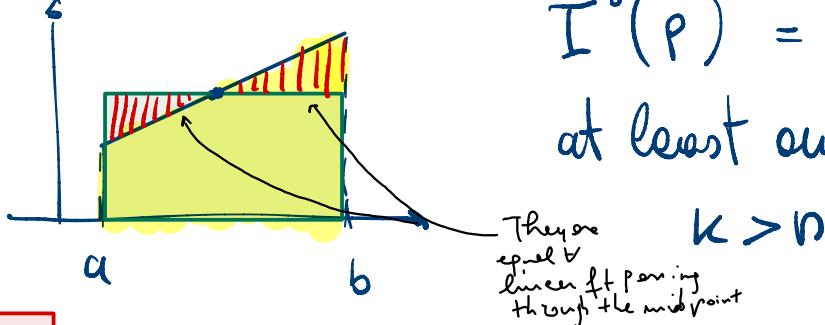
Interpolatory quadrature rules with  $n+1$  points are exact at least for polynomials of order  $\leq n$

Proof:  $\forall p \in P_{odd}^{n+1}$ ,  $L_p^n = p$  by construction

$$\Rightarrow I(p) = I^n(p) \quad \forall p \in P^n$$

Degree of accuracy of  $I^n$  is the maximum integer  $K$  s.t.  $I(p) = I^n(p) \quad \forall p \in P^K$

For the midpoint quadrature rule



$$I^0(p) = I(p) \quad \forall p \in P^1$$

at least one case where

$k > n$   
because  
they are equal & lie left, passing through the midpoint

Theo

Given  $\{q_i\}_{i=0}^n$ , degree of accuracy is  $< 2(n+1)$

Proof : construct  $w(x) = \prod_{i=0}^n (x - q_i)$ ,  $w \in P^{n+1}$

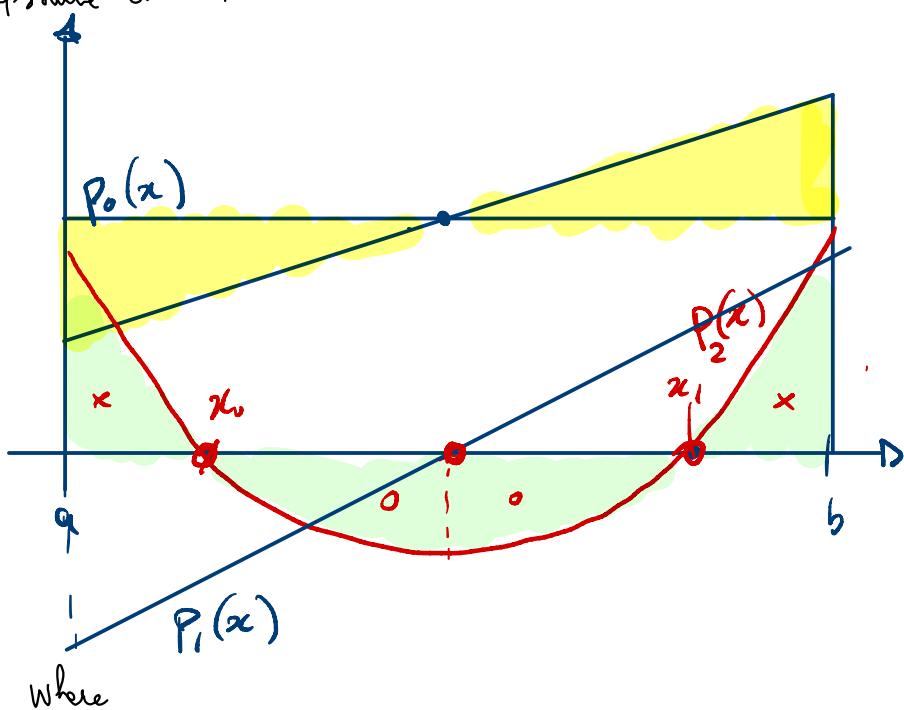
$w^2(x) > 0 \quad \forall x \neq q_i, w(q_i) = 0 \quad i = 0, \dots, n$

$\Rightarrow I(w^2) > 0 \quad I^b(w^2) = 0 \quad \text{because by construction.}$

$\Rightarrow \exists p = w^2 \in P^{2n+2} \text{ s.t. } I(p) + I^b(p)$

Can the degree of accuracy be  $k = 2n+1$ ? Yes

Assume current situation



Areas marked with  $x$  and with  $o$  must be equal

equivalent to saying  $\equiv$

average of  $P_2(x)$  is zero

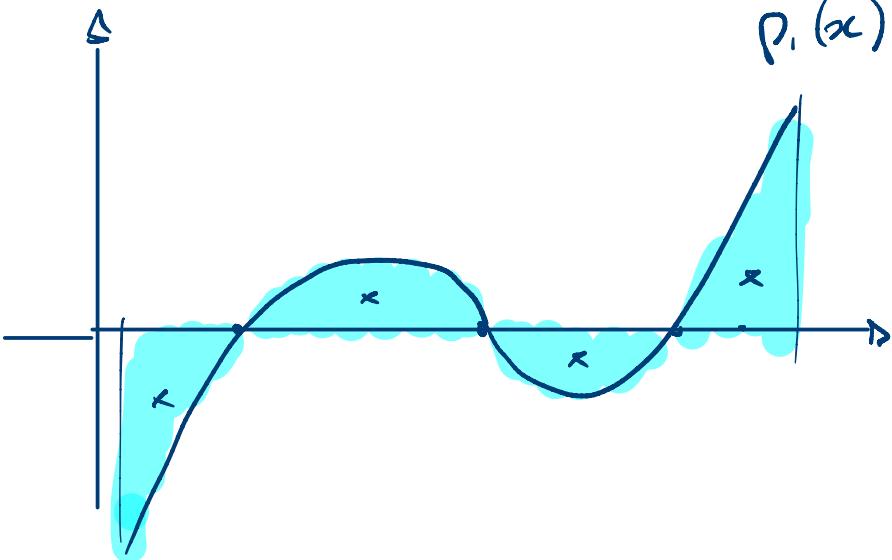
where

$$P_0(x) = 1 \quad \xrightarrow{\frac{a+b}{2}}$$

$$P_1(x) = (x - 0.5)$$

$$P_2(x) = (x - x_0)(x - x_1) \quad \text{s.t. } \int P_2(x) = 0$$

Now take the pt  
 $p_1(x) \cdot p_2(x) \in P^3(x)$



the areas marked  
with  $x$  are  
equal

Theo

Let  $u \in P^{n+m}$

$m \leq n+1$

then  $I^n(u) = I(u)$

$\Leftrightarrow$

$I^n$  has degree of accuracy  $k = n+m$

$\Leftrightarrow$

$$\int_a^b \omega(x) p = 0 \quad \forall p \in P^{m-1}$$

can be chosen as  $\perp$  to every  $p$

$$\omega = \prod_{i=0}^n (x - q_i) \in P^{n+1}$$

we are asking  $\omega$  to be Legende polynomial

Proof

Any polynomial  $p \in P^{n+m}$  can be

written

as :

$$p^{n+m}$$

$$\omega(x)$$

$$P^{m-1}$$

$$P_m$$

$$p(x) = \underbrace{\omega(x)}_{\text{by}} \underbrace{\prod(x)}_{\text{P}_m} + q(x)$$

basically  
division + rest

Ruffini's theorem

Given  $g(x) \in \mathbb{P}^{m-1} \subseteq \mathbb{P}^n$  ( $m \leq n+1$ )

Then  $I(g) = I^n(g)$

$$I(p) = I(\omega\pi) + I(g) \quad \text{by linearity of interpolation}$$

$$I^n(p) = I^n(\omega\pi) + I^n(g) \\ = I(g)$$

$$\text{If } I(\omega\pi) = 0 = \int_0^1 \omega\pi dx \quad \text{then } I(p) = I^n(p)$$

Now we choose  $\{g_i\}_{i=0}^n$  as the roots of Legendre basis

of order  $n+1$ :  $\int_a^b \omega_{n+1} l_j dx = 0 \quad j < n+1$

Legendre basis of  
order  $n+1$

The roots of  $\omega_{n+1}$  are called Gauss quadrature points. And  $I^n$  in this case is a Gauss quadrature formula, and it has degree of accuracy  $K = 2n+1$

General Strategy : find  $\{q_i\}, \{w_i\}$  s.t.

$$\sum_{j=0}^n v_i(q_j) w_j = \int_a^b v_i(x) dx \quad \text{for } \{v_i\}_{i=0}^{2n+1} \text{ known}$$

you impose that the interpolatory quadrature formula is EXACT for a fixed set of  
orthogonal functions

$2(n+1)$  unknowns, in  $2n+2$  equations.

NON LINEAR SYSTEM OF  $2n+2$  EQUATIONS.

if you can solve this, you win.

for  $v_i = x^i$  you obtain  
the GAUSS QUADRATURE  
FORMULA