

Numerical Methods for Computational Science and Engineering

Fall Semester 2017 (HS17)

Prof. Rima Alaifari, SAM, ETH Zurich

8. Iterative Methods for Nonlinear Systems of Equations

So far: learned direct methods for solving
linear systems of equations

Many models in real applications involve
nonlinear systems of equations

In general: these systems can't be solved directly

(not exactly)!

→ iterative methods for finding approximations
to the solution instead.

Example: liquid in spherical tank

r ... radius of the tank

g ... rate of constant liquid flow

Full tank ($h_0 = 2r$) at time $t=0$.

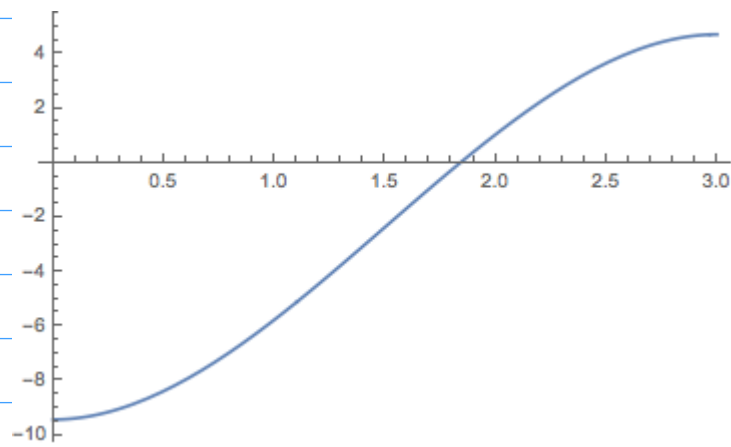
Height h of fluid at any time t :

$$-\frac{1}{3}\pi h^3 + \pi r h^2 + (g t - \frac{4}{3}\pi r^3) = 0$$

Task: Given any time t , determine height h .

Define $f_t(h) := -\frac{1}{3}\pi h^3 + \pi r h^2 + (g t - \frac{4}{3}\pi r^3)$

Find root h^* , i.e. solve for $f(h^*) = 0$.



$$t = \frac{1}{3} \cdot \underbrace{\frac{4}{3} \pi r^3}_S$$

$$(r = 1.5m)$$

Typical examples of nonlinear equations:

- thermodynamic models [involve equations of state for real gases]
- Colebrook equation for friction factor [= pressure drop in oil or gas pipeline]

How to solve a nonlinear equation

$$f(x) = 0 \quad ?$$

One question before that:

When is $f(x) = 0$ solvable?

Take $f(x) = e^{-\pi x^2}$ or $f(x) = \text{sign}(x)$

Both do not have roots.

First simple criterion: intermediate value theorem!

Theorem [IVT]: If $f: [a, b] \rightarrow \mathbb{R}$ is continuous and for some $t_l, t_r \in [a, b]$:

$$f(t_l) < u < f(t_r)$$

Then, there exists $z \in (t_l, t_r)$ s.t.

$$f(z) = u.$$

\swarrow f is real-val.

For root finding of $f \in C^0([a, b], \mathbb{R})$

If for some $t_l, t_r \in [a, b]$: $f(t_l) < 0$ & $f(t_r) > 0$

$\Rightarrow f$ has root in (t_l, t_r) .

→ Idea for our first algorithm for root finding:

Bisection algorithm: (for finding root x^*)

While $n \leq n_{max}$

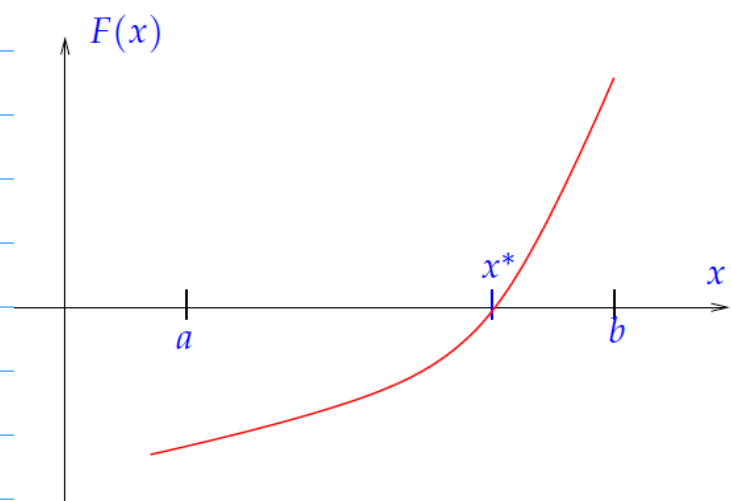
Compute $b \leftarrow \frac{t_l + t_r}{2}$

If $|f(b)| < TOL1$ or $|t_l - t_r| < TOL2$, return $x^* = b$.

$n \leftarrow n + 1$

If $\text{sign}(f(t_l)) == \text{sign}(f(b))$, then $t_l \leftarrow b$

Else $t_r \leftarrow b$



If f is cont: we always search in a region where a root exists

Region size is halved in every iteration step → convergence

Example [height of fluid in spherical tank]

Find height at time $t^* = \frac{1}{3} \cdot \frac{4}{3} \pi r^3$

t_{end} (no fluid in tank for first time)

$$f(h) = -\frac{1}{3} \pi h^3 + \pi r h^2 - \frac{8}{9} \pi r^3$$

We know: $f(0) < 0$

$$f(2r) > 0$$

→ apply bisection method for $t_l = 0$ and $t_r = 2r$.

For iterates $(x^{(k)})_{k \in \mathbb{N}}$ of approximate solutions
define iteration error $e^{(k)} = x^{(k)} - x^*$
 \nwarrow solution

Rate of convergence $x^{(k)} \rightarrow x^*$?

For bisection method:

$$|e^{(1)}| \leq \frac{1}{2} |a-b|$$

$$|e^{(2)}| \leq \frac{1}{2^2} |a-b|$$

$$|x^{(k)} - x^*| \leq 2^{-k} |a-b| \quad |e^{(k)}| \rightarrow 0 \text{ as } k \rightarrow \infty$$

"linear-type" convergence

error reduced by a fixed factor (here: $\frac{1}{2}$)
in each step.

Bisection: \oplus robustness, global convergence

\ominus rather slow convergence

no extension to higher dimensions

\downarrow
[n different quantities have
to be zero simultaneously]

Can we find methods that

- extend to higher dimensions
- guarantee faster convergence (under additional assumptions)

Fixed Point Iterations:

Root finding of f , i.e. $f(x) = 0$

\Leftrightarrow finding fixed point of $\Phi(x) := f(x) + x$

$$x^* \text{ is a fixed point (FP) of } \Phi \Leftrightarrow \Phi(x^*) = x^*$$

$$\Leftrightarrow f(x^*) = 0$$

For bisection: only needed continuity of f

Now: Assume that Φ is Lipschitz cont. on $[a, b]$:

$$\exists L > 0 \quad \forall x, y \in [a, b] \quad : \quad |\Phi(x) - \Phi(y)| \leq L \cdot |x - y|$$

notice order

Further assume $L < 1$ (Φ is a contractive mapping)

Idea of FPI:

- Start with initial guess $x^{(0)}$
- Iterate $x^{(k)} = \Phi(x^{(k-1)})$ [FPI]

$L < 1$ guarantees: If Φ has a fixed point, i.e. there exists x^* s.t. $\Phi(x^*) = x^*$, FPI will converge to x^* .

Derivation: We need to show $|e^{(k)}| \rightarrow 0$ as $k \rightarrow \infty$.

$$|e^{(k)}| = |x^{(k)} - x^*| = \underbrace{|\Phi(x^{(k-1)})|}_{\text{def. of } x^{(k)}} - \underbrace{|\Phi(x^*)|}_{x^* \text{ is FP}}$$

$$\leq L \cdot |x^{(k-1)} - x^*| = \underline{L \cdot |e^{(k-1)}|}$$

$$\Rightarrow |e^{(k)}| \leq \underbrace{L^k}_{\rightarrow 0} \cdot |e^{(0)}| \rightarrow 0 \text{ as } k \rightarrow \infty. \quad \square$$

$L < 1!!!$

Remarks:

- It would suffice to have Φ Lipschitz with $L < 1$ on $[x^* - \delta, x^* + \delta]$ and $x^{(0)} \in [x^* - \delta, x^* + \delta]$

- If $\Phi \in C^1([a,b])$ and $|\Phi'(x^*)| < 1$

Convergence of FPI in a neighborhood of x^* :

There are $\delta > 0, \varepsilon > 0$ s.t. $|\Phi'(x)| < 1 - \varepsilon$

for all $x \in [x^* - \delta, x^* + \delta] =: I^*$

Then, for all $x, y \in I^* \exists \theta \in [x, y]$ s.t.

$$|\Phi(x) - \Phi(y)| = |\Phi'(\theta)| \cdot |x - y| \quad [\text{mean value theorem}]$$

$$< (1 - \varepsilon) |x - y|$$

$\Rightarrow \Phi$ is Lipschitz with $L = 1 - \varepsilon < 1$ on I^* .

- Convergence rate is (at least) linear:

Definition 8.1.9. Linear convergence

A sequence $x^{(k)}, k = 0, 1, 2, \dots$, in \mathbb{R}^n converges linearly to $x^* \in \mathbb{R}^n$,

$$\exists 0 < L < 1: \underbrace{\|x^{(k+1)} - x^*\|}_{\|e^{(k+1)}\|} \leq L \underbrace{\|x^{(k)} - x^*\|}_{\|e^{(k)}\|} \quad \forall k \in \mathbb{N}_0.$$

[Note: bisection was not linear, only of "linear-type"

because there $|e^{(k)}| \leq L^k |a - b|$

but $|e^{(k)}| > L |e^{(k-1)}|$ was possible]

- If $\Phi'(x^*) = 0$ and Φ is C^2 in a neighborhood of x^* :

Taylor expansion of Φ at $x^{(k)}$ around x^* :

$$\Phi(x^{(k)}) = \Phi(x^*) + \underbrace{e^{(k)} \cdot \Phi'(x^*)}_{=0} + \frac{1}{2} (e^{(k)})^2 \Phi''(x^*) + \mathcal{O}((e^{(k)})^3)$$

$$|e^{(k)}| = |x^{(k)} - x^*| = |\Phi(x^{(k-1)}) - \Phi(x^*)|$$

$$= \frac{1}{2} (e^{(k-1)})^2 |\Phi''(x^*)| + \mathcal{O}((e^{(k-1)})^3)$$

if $|x^{(k-1)} - x^*| \leq \delta < 1$ [guaranteed by conv.]:

$$|e^{(k)}| \leq \frac{1}{2} |e^{(k-1)}|^2 (|\Phi''(x^*)| + \delta)$$

$$|e^{(k)}| \leq C \cdot |e^{(k-1)}|^2$$

↑
quadratic convergence

Definition 8.1.17. Order of convergence → [?, Sect. 17.2], [?, Def. 5.14], [?, Def. 6.1]

A convergent sequence $x^{(k)}$, $k = 0, 1, 2, \dots$, in \mathbb{R}^n with limit $x^* \in \mathbb{R}^n$ converges with order p , if

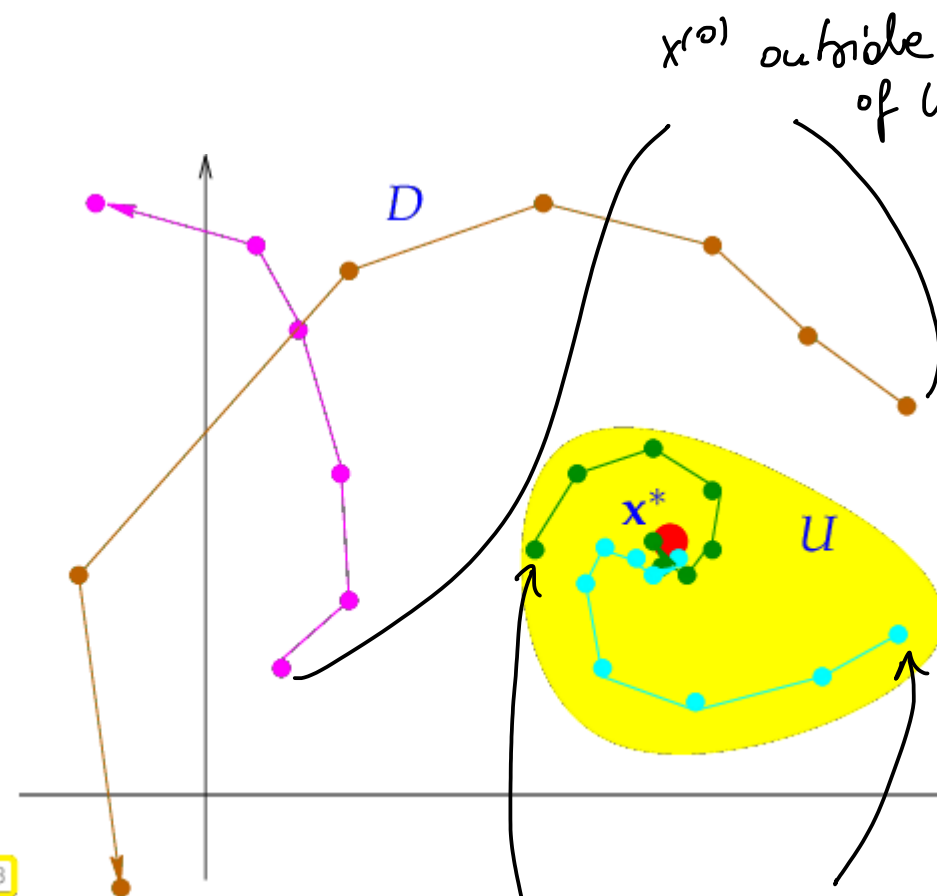
$$\exists C > 0: \|x^{(k+1)} - x^*\| \leq C \|x^{(k)} - x^*\|^p \quad \forall k \in \mathbb{N}_0, \quad (8.1.18)$$

and, in addition, $C < 1$ in the case $p = 1$ (linear convergence → Def. 8.1.9).

higher p = faster convergence (fewer iterations to reach same accuracy)

Note: If Φ is not globally Lipschitz (on $[a, b]$) but only locally (e.g. $\Phi \in C^1$, $|\Phi'(x^*)| < 1$) then FPI is only locally convergent (i.e. need $x^{(0)}$ suff. close to x^*)

Fig. 283



(illustration in 2D)

$x^{(0)}$ close in region where convergence is guaranteed

Algorithm for root-finding with quadratic convergence

Assumption needed: $f \in C^1$
↑
gives first order Taylor appr.

Intuition: Approximation around $x^{(k)}$:

$$f(x) \approx f(x^{(k)}) + (x - x^{(k)}) \cdot f'(x^{(k)})$$

To find x^* with $f(x^*) = 0$ take next iterate $x^{(k+1)}$

$$\text{s.t. } f(x^{(k)}) + (x^{(k+1)} - x^{(k)}) f'(x^{(k)}) = 0$$

⇒ Newton iteration:

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

Quadratic convergence when $f \in C^2$?

Reformulate as FPI: $\Phi(x) := x - \frac{f(x)}{f'(x)}$

$$\left[\begin{array}{l} \Phi(x^{(k)}) = x^{(k+1)} \\ \Phi(x^{(k)}) = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})} \end{array} \right]$$

Newton's method on $f \Leftrightarrow$ FPI on $\Phi(x)$

$$\Phi'(x) = 1 - \frac{(f'(x))^2 - f''(x) \cdot f(x)}{(f'(x))^2}$$

$$= \frac{f''(x) f(x)}{(f'(x))^2}$$

If $f'(x^*) \neq 0$ then

↑
simple root

$$\underline{\underline{\Phi'(x^*) = 0}}$$

↓
∃ neighborhood around x^*
 $|\Phi'| < 1$

For $x^{(0)}$ in a neighborhood I^* of x^*

$$[s.t. \forall x \in I^* \quad |\Phi'(x)| \leq 1 - \varepsilon]$$

Newton's method converges **quadratically** if

$$f'(x^*) \neq 0.$$

(due to $\Phi'(x^*) = 0$)

Note: $x^{(0)} \in I^*$ guarantees $f'(x^{(k)}) \neq 0$ for all the iterates. This has to be guaranteed

• For quadratic convergence of Newton's method

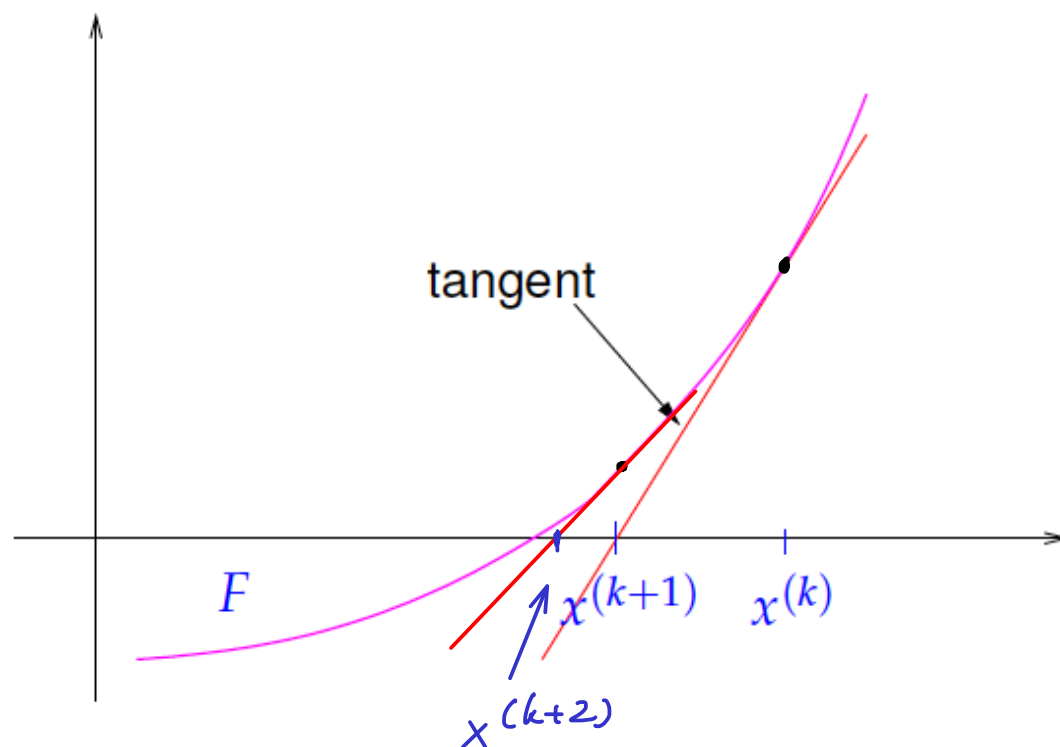
$$f \in C^2(I^*) \text{ suffices instead of } \Phi \in C^2(I^*)$$

In summary: we need a neighborhood I^* of x^* s.t.

• I^* suff. small

• $f'(x) \neq 0$ on I^*

• $f \in C^2(I^*)$



Example:

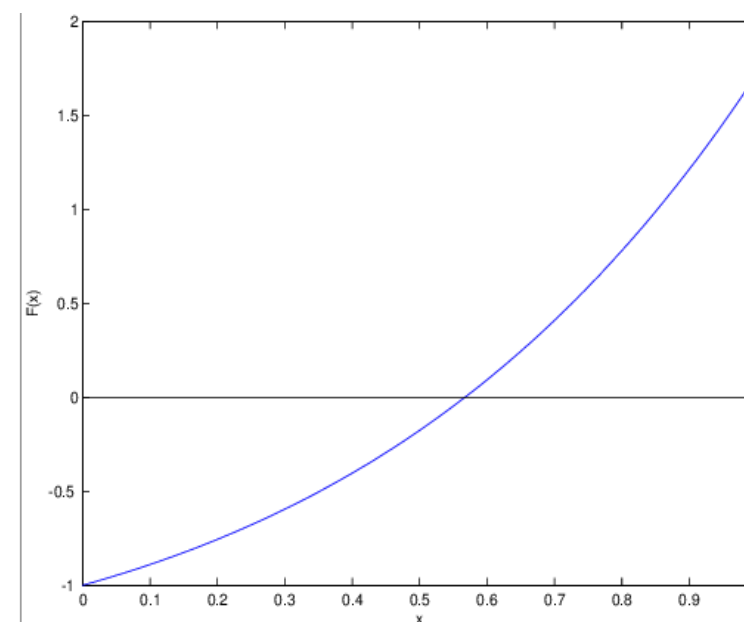
$$F(x) = xe^x - 1, \quad x \in [0, 1].$$

Different fixed point forms:

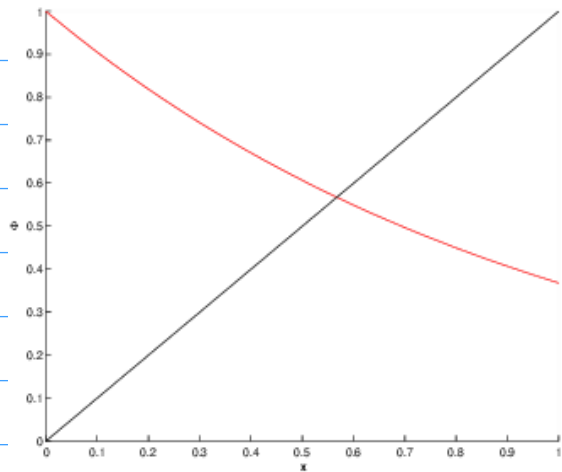
$$\Phi_1(x) = e^{-x},$$

$$\Phi_2(x) = \frac{1+x}{1+e^x},$$

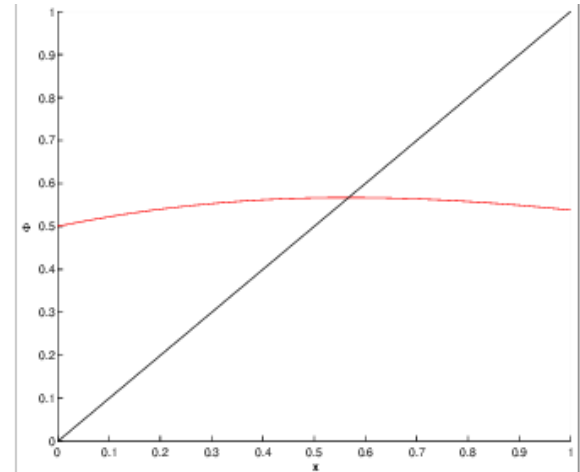
$$\Phi_3(x) = x + 1 - xe^x.$$



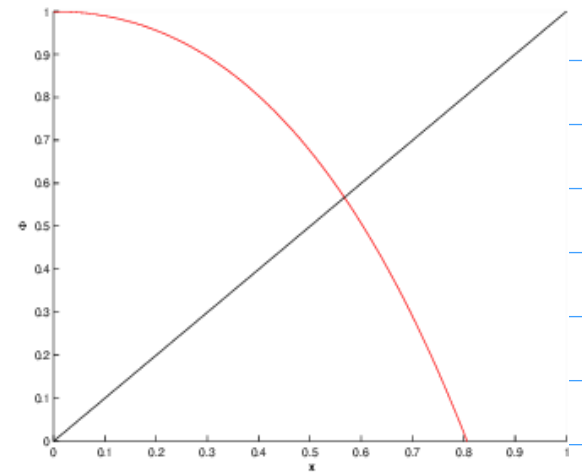
$$x^* = e^{-x^*} \Leftrightarrow x^* e^{x^*} = 1$$



function Φ_1



function Φ_2



function Φ_3

FPI with $x^{(0)} = 0.5$:

$$\Phi_1(x^*) = x^*$$

$$\Phi_2(x^*) = x^*$$

$$\Phi_3(x^*) = x^*$$

k	$ x_1^{(k+1)} - x^* $	$ x_2^{(k+1)} - x^* $	$ x_3^{(k+1)} - x^* $
0	0.067143290409784	0.067143290409784	0.067143290409784
1	0.039387369302849	0.000832287212566	0.108496074240152
2	0.021904078517179	0.000000125374922	0.219330611898582
3	0.012559804468284	0.0000000000000003	0.288178118764323
4	0.007078662470882	0.0000000000000000	0.723649245792953
5	0.004028858567431	0.0000000000000000	0.410183132337935
6	0.002280343429460	0.0000000000000000	1.186907542305364
7	0.001294757160282	0.0000000000000000	0.146569797006362
8	0.000733837662863	0.0000000000000000	0.310516641279937
9	0.000416343852458	0.0000000000000000	0.357777386500765
10	0.000236077474313	0.0000000000000000	0.974565695952037

roughly: FPI for Φ_1 : linearly conv.

FPI for Φ_2 : quadr. conv.

Φ_3 : no conv.

Why?

$$\Phi_1'(x) = -e^{-x}$$

$$\Rightarrow |\Phi_1'(x)| < 1 \text{ for } x \in [\delta, 1] = I_\delta^* \quad \forall \delta > 0$$

\rightarrow conv. of FPI with Φ_1 in I^*

$$\Phi_2'(x) = \frac{1 - xe^x}{(1 + e^x)^2}$$

$$|\Phi_2'(x)| < \frac{|1-e|}{4} < \frac{1}{2}$$

\rightarrow conv. of FPI of Φ_2 on $[0, 1]$

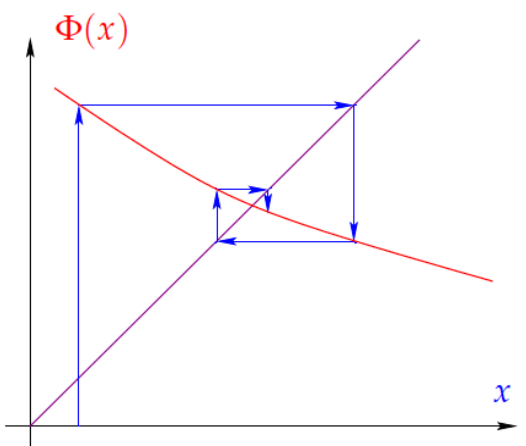
Furthermore $\Phi_2'(x^*) = 0 \rightarrow$ quadr convergence

$$\Phi_3'(x) = 1 - e^x - xe^x \quad e^{-x^*} = x^*$$

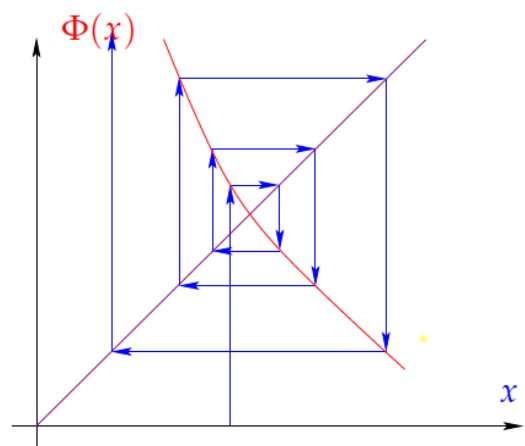
$$\Phi_3'(x^*) = -e^{x^*} = -\frac{1}{x^*}$$

$$x^* \in (0, 1) : |\Phi_3'(x^*)| > 1$$

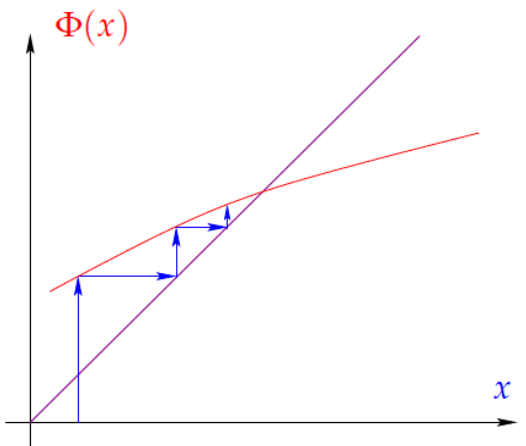
⇒ FPI not contractive around x^*



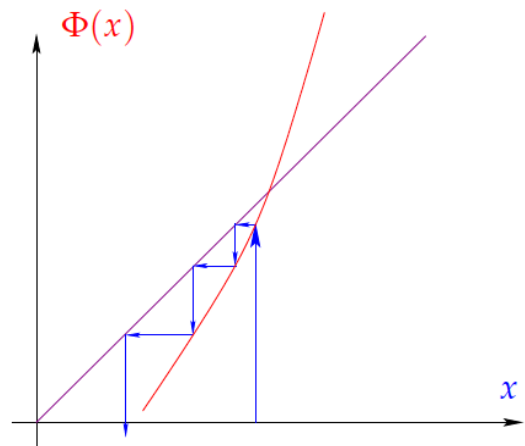
$-1 < \Phi'(x^*) \leq 0 \rightarrow$ convergence



$\Phi'(x^*) < -1 \rightarrow$ divergence



$0 \leq \Phi'(x^*) < 1 \rightarrow$ convergence



$1 < \Phi'(x^*) \rightarrow$ divergence

Remark on Newton's method

Requires computation of $f'(x^{(k)})$ for each iteration!

↑
can be costly

Alternative: Secant method

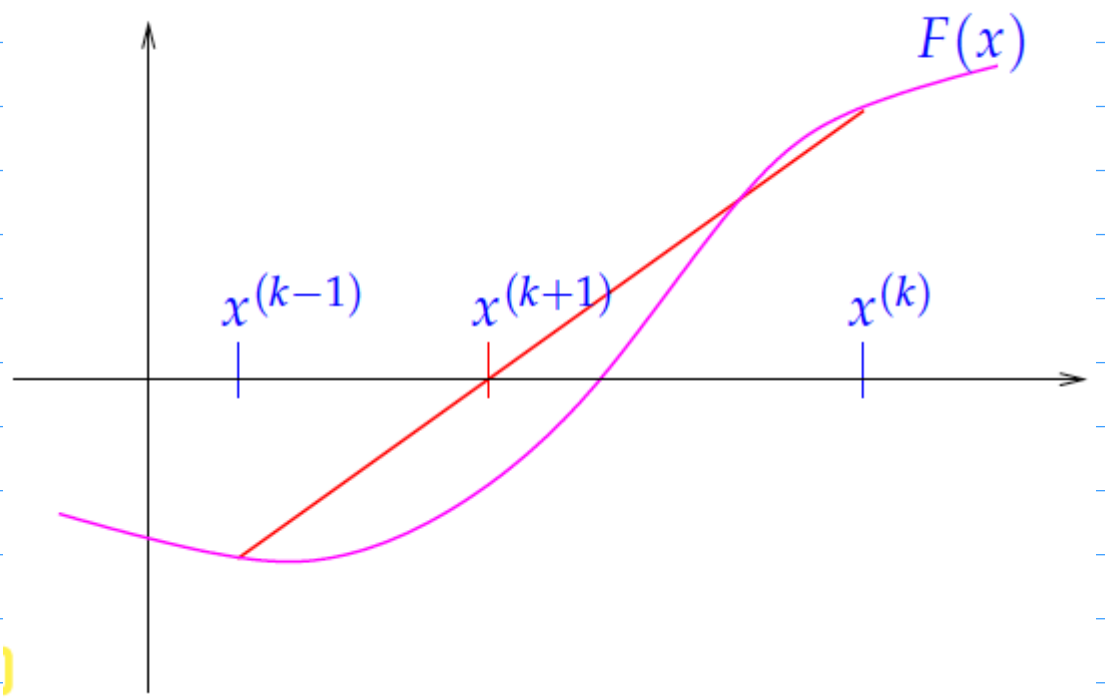
replace $f'(x^{(k)})$ by approximation

$$\frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$$

in Newton iteration to obtain

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)}) (x^{(k)} - x^{(k-1)})}{f(x^{(k)}) - f(x^{(k-1)})}$$

Remark: Secant method is a 2-point method
 (computing $x^{(k+1)}$ involves $x^{(k)}, x^{(k-1)}$)



Definition: Stationary m-point iterative method
 $x^{(k)}$ depends on m most recent iterates

$x^{(k-1)}, \dots, x^{(k-m)}$ iteration function

$$x^{(k)} = \Phi_F(x^{(k-1)}, \dots, x^{(k-m)})$$

for solving $F(x) = 0$.

Convergence of secant method:

- again local
 - need $f'(x^*) \neq 0$ [simple root]
 - f locally C^2
 - rate is superlinear but not quadratic
- } as for Newton's method

order $\rho = \frac{1+\sqrt{5}}{2} \approx 1.618$

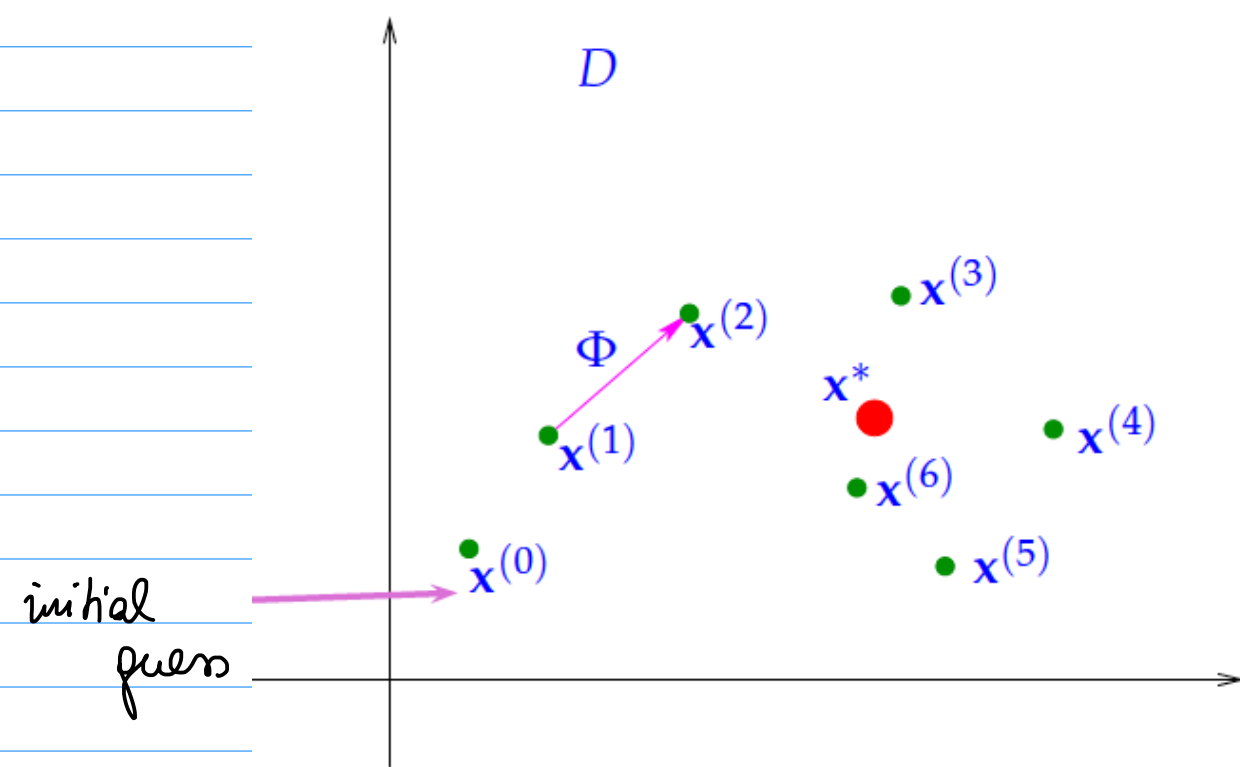
Note: m-point method requires m initial guesses

$$x^{(0)}, \dots, x^{(m-1)}$$

Nonlinear systems of equations

$F: \mathbb{D} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ (nonlinear system of eqns
n eqns, n unknowns)

Find x^* s.t. $F(x^*) = 0$.



$$x^{(k)} = \Phi_F(x^{(k-1)}, \dots, x^{(k-m)})$$

Aspects of iterative methods:

- Convergence: $(x^{(k)})_{k \in \mathbb{N}}$ convergent, $\lim_{k \rightarrow \infty} x^{(k)} = x^*$
- Consistency: $\Phi_F(x^*, \dots, x^*) = x^* \Leftrightarrow F(x^*) = 0$
- Rate of convergence $\|x^{(k)} - x^*\| \rightarrow 0$
with which order?

Note: $\|\cdot\|$ can be any norm for \mathbb{R}^n

in \mathbb{R}^n (finite dim. vector space): all norms
are equivalent

Definition 8.1.11. Equivalence of norms

Two norms $\|\cdot\|_a$ and $\|\cdot\|_b$ on a vector space V are equivalent if

$$\exists \underline{C}, \bar{C} > 0: \underline{C}\|v\|_a \leq \|v\|_b \leq \bar{C}\|v\|_a \quad \forall v \in V.$$

This implies: convergence in \mathbb{R}^n independent of choice of norm

But in general: convergence rate depends on chosen norm

Local vs. global convergence:

Definition 8.1.8. Local and global convergence \rightarrow [?, Def. 17.1]

As stationary m -point iterative method **converges locally** to $x^* \in \mathbb{R}^n$, if there is a neighborhood $U \subset D$ of x^* , such that

$$x^{(0)}, \dots, x^{(m-1)} \in U \Rightarrow x^{(k)} \text{ well defined } \wedge \lim_{k \rightarrow \infty} x^{(k)} = x^*$$

where $(x^{(k)})_{k \in \mathbb{N}_0}$ is the (infinite) sequence of iterates.
 If $U = D$, the iterative method is **globally convergent**.

Fixed Point Iterations in \mathbb{R}^n

Definition 8.2.1.

A fixed point iteration $x^{(k+1)} = \Phi(x^{(k)})$ is consistent with $F(x) = 0$ if for $x \in U \cap D$

$$F(x) = 0 \iff \Phi(x) = x.$$

Definition 8.2.6. [Contractive mapping]

$\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ is contractive (w.r.t. norm $\|\cdot\|$ on \mathbb{R}^n) if

$$\exists L < 1 \quad \|\Phi(x) - \Phi(y)\| \leq L \cdot \|x - y\| \quad \forall x, y \in U$$

① Contractivity of $\Phi \Rightarrow$ if $\Phi(x^*) = x^*$ then
 FPI will converge to x^* .

$$\underbrace{\|x^{(k+1)} - x^*\|}_{\|e^{(k+1)}\|} = \|\Phi(x^{(k)}) - \Phi(x^*)\| \leq \underbrace{L}_{\leq 1} \cdot \underbrace{\|x^{(k)} - x^*\|}_{\|e^{(k)}\|}$$

$$\|e^{(k+1)}\| \leq L^k \|e^{(0)}\|$$

② Convergence is at least linear.

③ If Φ is contractive $\Rightarrow \Phi$ has at most one FP.

Why? Suppose 2 FPs x_1^*, x_2^* :

$$\|x_1^* - x_2^*\| \underset{\text{FP}}{=} \|\Phi(x_1^*) - \Phi(x_2^*)\| \underset{\text{contractive}}{\leq} L \cdot \|x_1^* - x_2^*\|$$

with $L < 1 \Rightarrow x_1^* = x_2^*$.

Existence of a FP?

Theorem 8.2.9. Banach's fixed point theorem

If $D \subset \mathbb{K}^n$ ($\mathbb{K} = \mathbb{R}, \mathbb{C}$) closed and bounded and $\Phi: D \rightarrow D$ satisfies

$$\exists L < 1: \|\Phi(x) - \Phi(y)\| \leq L\|x - y\| \quad \forall x, y \in D,$$

then there is a unique fixed point $x^* \in D, \Phi(x^*) = x^*$, which is the limit of the sequence of iterates $x^{(k+1)} := \Phi(x^{(k)})$ for any $x^{(0)} \in D$.

Convergence criteria for FPI for Φ differentiable
 & knowing $\Phi(x^*) = x^*$.

Lemma 8.2.10. Sufficient condition for local linear convergence of fixed point iteration \rightarrow
 [?, Thm. 17.2], [?, Cor. 5.12]

If $\Phi: U \subset \mathbb{R}^n \rightarrow \mathbb{R}^n, \Phi(x^*) = x^*, \Phi$ differentiable in x^* , and $\|D\Phi(x^*)\| < 1$, then the fixed point iteration

$$x^{(k+1)} := \Phi(x^{(k)}), \tag{8.2.2}$$

converges locally and at least linearly. matrix norm, Def. 1.5.76!

$$D\Phi(x) := \left[\frac{\partial \Phi_j}{\partial x_i}(x) \right]_{j,i=1}^n \in \mathbb{R}^{n,n}$$

\uparrow
Jacobian

Lemma 8.2.12. Sufficient condition for linear convergence of fixed point iteration

Let U be convex and $\Phi : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ be continuously differentiable with

$$L := \sup_{x \in U} \|D\Phi(x)\| < 1.$$

If $\Phi(x^*) = x^*$ for some interior point $x^* \in U$, then the fixed point iteration $x^{(k+1)} = \Phi(x^{(k)})$ converges to x^* at least linearly with rate L .

\Rightarrow Locally contractive $\Phi \Rightarrow$ iteration converges locally around FP (at least linear conv.)

Termination criteria for contractive FPT:

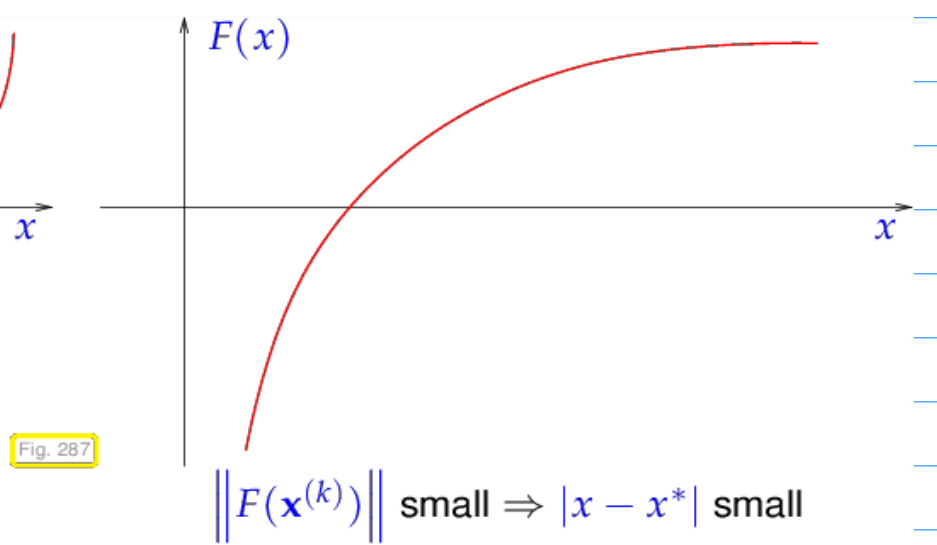
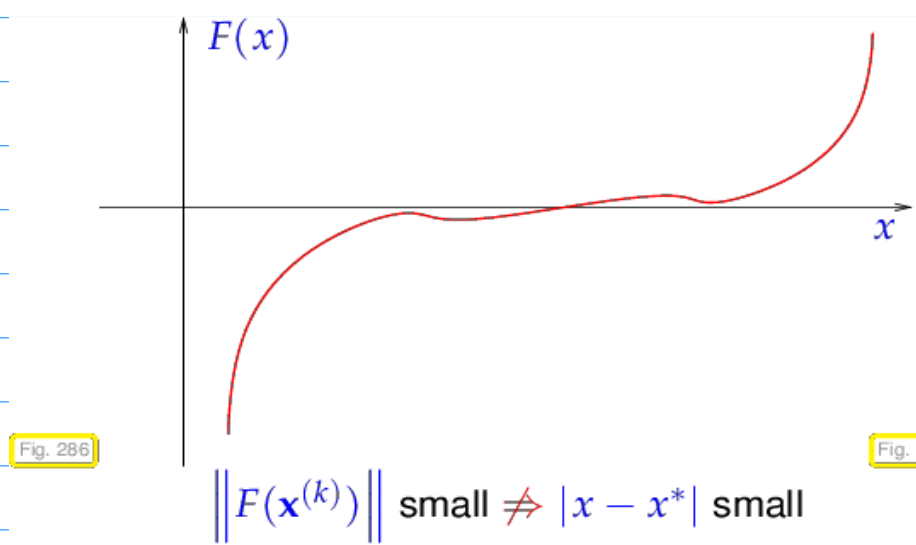
When to stop iterating (for finding x^* s.t. $F(x^*)=0$)?

① residual based
Stop when $\|F(x^{(k)})\| \leq \tau$

② Correction based
Stop when $\|x^{(k+1)} - x^{(k)}\| \leq \tau$
or $\|x^{(k+1)} - x^{(k)}\| \leq \tau_{rel} \|x^{(k+1)}\|$

Recall discussion about cond. number:

$$\underbrace{\|F(x^{(k)}) - \underbrace{F(x^*)}_{=0}\|}_{\text{compute}} \text{ small} \not\Rightarrow \underbrace{\|x^{(k)} - x^*\|}_{\text{can't compute}} \text{ small}$$



Ultimate goal: guarantee of the form $\|x^{(k)} - x^*\| \leq \tau$

If iteration is linearly convergent:

$$\begin{aligned} \|x^{(k)} - x^*\| &\leq \|x^{(k+1)} - x^{(k)}\| + \|x^{(k+1)} - x^*\| \\ &\leq \|x^{(k+1)} - x^{(k)}\| + L \|x^{(k)} - x^*\| \end{aligned}$$

$$\Rightarrow (1-L) \|x^{(k)} - x^*\| \leq \|x^{(k+1)} - x^{(k)}\|$$

$$\Rightarrow \underbrace{\|x^{(k+1)} - x^*\|}_{\text{not computable}} \leq L \|x^{(k)} - x^*\|$$

$$\leq \frac{L}{1-L} \underbrace{\|x^{(k+1)} - x^{(k)}\|}_{\text{computable!}}$$

Suggests to ask for

$$\frac{L}{1-L} \|x^{(k+1)} - x^{(k)}\| \leq \tau$$

as stopping criterion. It guarantees

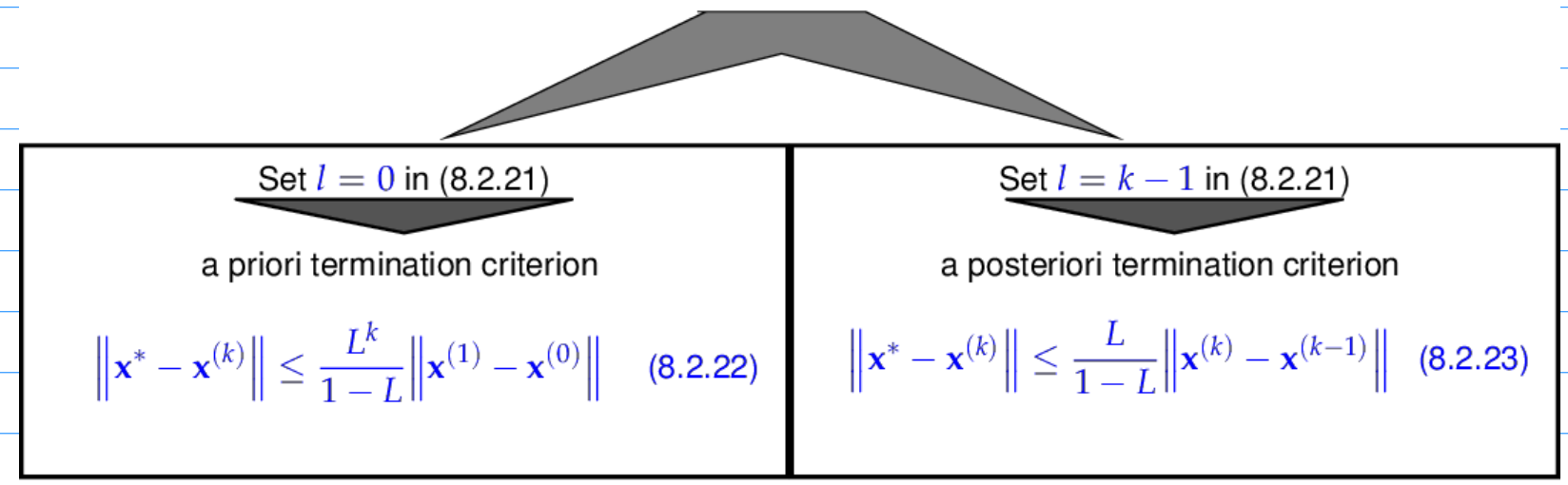
$$\|x^{(k+1)} - x^*\| \leq \tau.$$

Note: Estimating L can be difficult.

But: pessimistic estimate $\tilde{L} > L$ is still reliable!

More generally:

$$\|x^* - x^{(k)}\| \leq \frac{L^{k-l}}{1-L} \|x^{(l+1)} - x^{(l)}\|. \tag{8.2.21}$$



8.4. Newton's method (in higher dimensions)

Extend idea from 1D:

First order approximation (linearization)

$$F(x) \approx F(x^{(k)}) + \underbrace{DF(x^{(k)})}_{\in \mathbb{R}^{n,n} \text{ Jacobian of } F \text{ at } x^{(k)}} (x - x^{(k)}) =: \tilde{F}_k(x)$$

► **Newton iteration:** (generalizes (8.3.4) to $n > 1$)

$x^{(k+1)} := x^{(k)} - DF(x^{(k)})^{-1}F(x^{(k)})$, [if $DF(x^{(k)})$ regular] (8.4.1)

Terminology: $-DF(x^{(k)})^{-1}F(x^{(k)}) =$ Newton correction

Before: $x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$
needed $f'(x^{(k)}) \neq 0$

Now: Invertibility of Jacobian

To compute Newton correction: solve LSE

$DF(x^{(k)})y = -F(x^{(k)})$

Convergence of Newton's method:

If $F(x^*) = 0$ and $DF(x^*)$ is regular, then it is locally quadratically convergent.

Exact theorem: Thm 8.4.45 in lecture notes

~> hardly ever possible to verify in practice

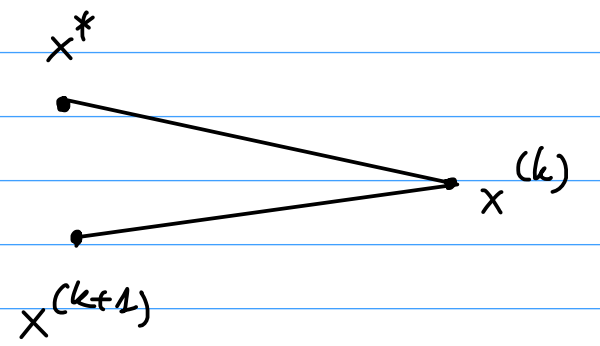
Note: If $DF(x^*)$ is singular, Newton is only locally linearly convergent.

A few remaining questions:

- 1., Stopping criterion for Newton's method?
- 2., Larger region of convergence possible? [At the cost of loosing quadratic convergence]
- 3., Newton correction is costly [Solve LSE with different system matrix in each iteration]
 - > Any remedy?

Ad 1. Quadratic convergence

$$\|x^{(k+1)} - x^*\| \ll \|x^{(k)} - x^*\|$$



Roughly: $\|x^{(k)} - x^*\| \approx \|x^{(k+1)} - x^{(k)}\|$

$$\|x^{(k+1)} - x^{(k)}\| = \underbrace{\|DF(x^{(k)})^{-1} F(x^{(k)})\|}_{\text{computable stopping criterion}} \leq \tau \|x^{(k)}\|$$

guarantee \downarrow

$$\|x^{(k)} - x^*\| \leq \tau \|x^{(k)}\|$$

BUT: If $x^{(k)}$ was a good approximation, we would have computed new

Newton correction $DF(x^{(k)})^{-1} F(x^{(k)})$

but not used in iteration.

Idea: "Cheaper" stopping criterion:

$$\|DF(x^{(k-1)})^{-1} F(x^{(k)})\| \leq \tau \|x^{(k)}\|$$

simplified Newton correction

Motivation: • Due to fast convergence

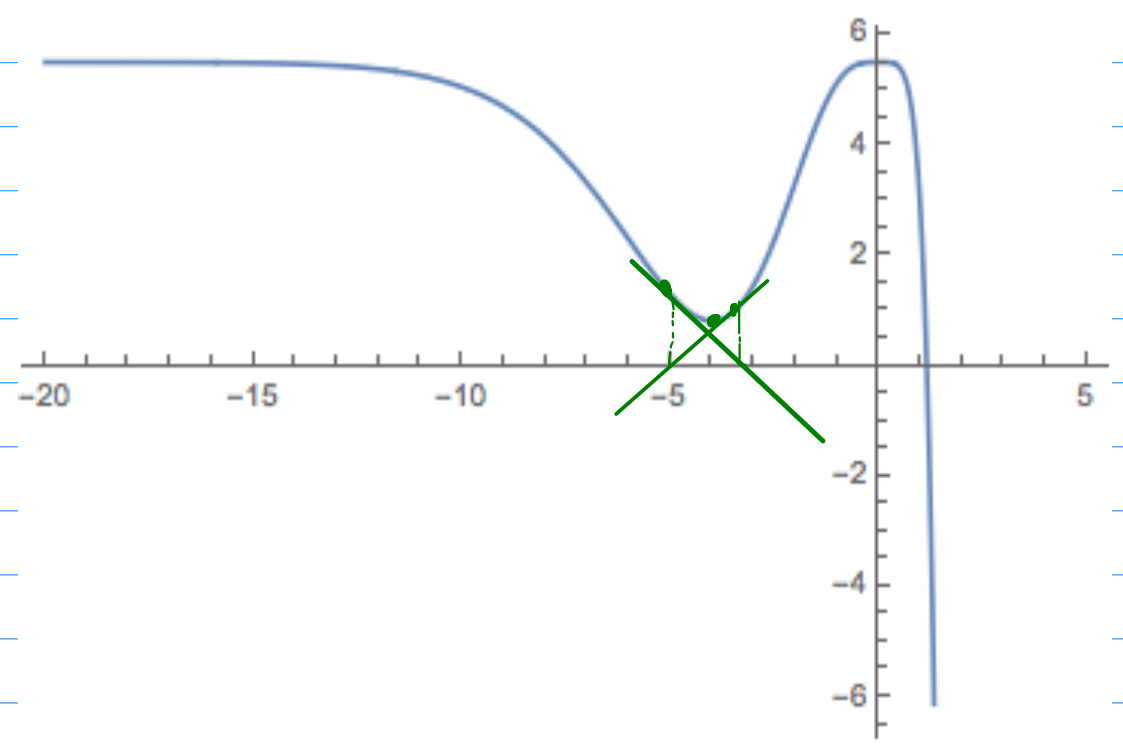
$$DF(x^{(k)}) \approx DF(x^{(k-1)})$$

during last steps of iteration

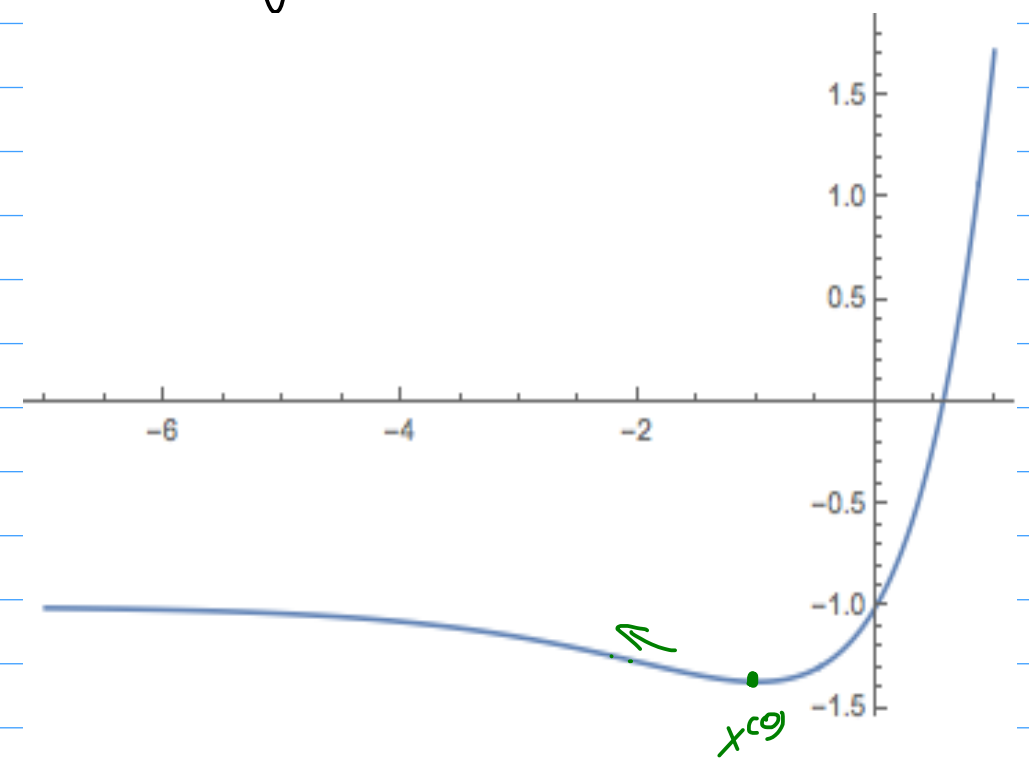
• We can reuse LU factorization of $DF(x^{(k-1)})$

Ad 2. Examples of failures of Newton's method

① Local min./max.



② Asymptotes



$$F(x) = xe^x - 1$$

$$F'(x) = xe^x + e^x$$

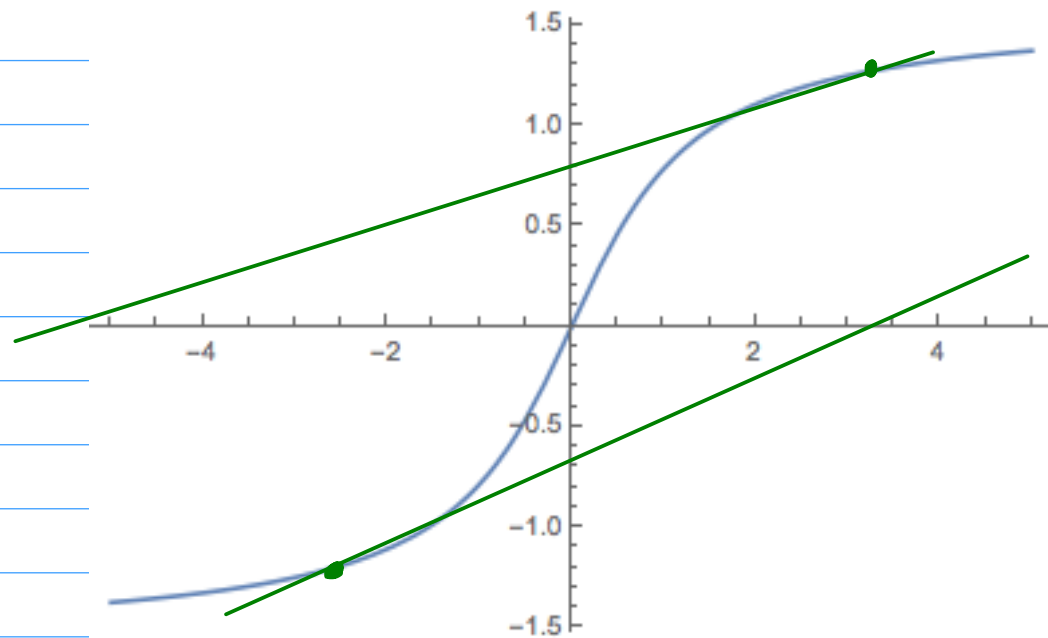
$$F'(-1) = 0$$

$$x^{(0)} < -1 \Rightarrow x^{(k)} \rightarrow -\infty$$

$$x^{(0)} > -1 \Rightarrow x^{(k)} \rightarrow x^*$$

③ Overshooting

$$F(x) = \arctan(x) \quad F(0) = 0$$



A remedy for overshooting:

8.4.4 Damped Newton method

Idea: In each iteration step, check whether distance $\|x^{(k+1)} - x^{(k)}\|$ is decreasing i.e. roughly whether at least

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

If not, don't take a **full** Newton step!

→ Damp the Newton correction

► we observe "overshooting" of Newton correction

Idea: **damping** of Newton correction:

$$\text{With } \lambda^{(k)} > 0: \quad x^{(k+1)} := x^{(k)} - \lambda^{(k)} \mathbf{D}F(x^{(k)})^{-1} F(x^{(k)}).$$

Terminology: $\lambda^{(k)}$ = damping factor

How to choose $\lambda^{(k)}$?

Strategy: largest possible $\lambda^{(k)}$ so that distance between iterates is decreasing

Affine invariant damping strategy

Choice of damping factor: affine invariant natural monotonicity test [?, Ch. 3]:

choose "maximal" $0 < \lambda^{(k)} \leq 1$: $\|\Delta \bar{x}(\lambda^{(k)})\| \leq (1 - \frac{\lambda^{(k)}}{2}) \|\Delta x^{(k)}\|_2$ (8.4.57)

where $\Delta x^{(k)} := DF(x^{(k)})^{-1} F(x^{(k)}) \rightarrow$ current Newton correction,
 $\Delta \bar{x}(\lambda^{(k)}) := DF(x^{(k)})^{-1} F(x^{(k)}) \leftarrow \lambda^{(k)} \Delta x^{(k)} \rightarrow$ tentative simplified Newton correction.



Approximation for $x^{(k+2)} - \tilde{x}^{(k+1)}$

where $\tilde{x}^{(k+1)} = x^{(k)} - \lambda^{(k)} \Delta x^{(k)}$

Check whether $\|\Delta \bar{x}(\lambda^{(k)})\|$ is strictly smaller than $\|x^{(k+1)} - x^{(k)}\|$ where $x^{(k+1)} = x^{(k)} - \Delta x^{(k)}$.

In practice: $\lambda^{(k)} = 1$ and check NMT, repeatedly take $\lambda^{(k)} \leftarrow \frac{\lambda^{(k)}}{2}$ until NMT passes for the first time.

Ad 3. Cheaper (approximate) Newton corrections?

Secant method in 1D:
 $f'(x^{(k)}) \approx \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$

$n > 1$?

Approximation $J_k \in \mathbb{R}^{n \times n} \approx \mathcal{D}F(x^{(k)})$ s.t.

$$J_k (x^{(k)} - x^{(k-1)}) = F(x^{(k)}) - F(x^{(k-1)}) \quad (*)$$

From Newton's method for $x^{(k)}$:

$$x^{(k)} = x^{(k-1)} - J_{k-1}^{-1} F(x^{(k-1)})$$

$$\Leftrightarrow J_{k-1} (x^{(k)} - x^{(k-1)}) = -F(x^{(k-1)}) \quad (**)$$

(*) - (**):

$$(J_k - J_{k-1}) (x^{(k)} - x^{(k-1)}) = F(x^{(k)})$$

underdetermined

Possible cheap choice: outer product in \mathbb{R}^n

$$J_k - J_{k-1} = \underbrace{\frac{F(x^{(k)}) (x^{(k)} - x^{(k-1)})^\top}{\|x^{(k)} - x^{(k-1)}\|_2^2}}_{\text{rank 1 matrix}}$$

Given initial J_0 (take $J_0 = \mathcal{D}F(x^{(0)})$),

get J_k by rank-1 updates:

$$J_k = J_{k-1} + \frac{F(x^{(k)}) (x^{(k)} - x^{(k-1)})^\top}{\|x^{(k)} - x^{(k-1)}\|_2^2}$$

Final form of Broyden's quasi-Newton method for solving $F(x) = 0$:

$$x^{(k+1)} := x^{(k)} + \Delta x^{(k)}, \quad \Delta x^{(k)} := -J_k^{-1} F(x^{(k)}),$$

$$J_{k+1} := J_k + \frac{F(x^{(k+1)}) (\Delta x^{(k)})^\top}{\|\Delta x^{(k)}\|_2^2}.$$

(8.4.66)

Note: Can use Sherman-Morrison-Woodbury formula to calculate J_k^{-1} from J_{k-1}^{-1} .

Remark: In general, iterative methods for nonlinear systems should have convergence monitor [i.e. simple check at each iteration whether convergence to be expected or not.

Example: NMT for damped Newton: if repeated failure: stop & report error