

A regular, diagonally dominant \Rightarrow partial pivoting according to (2.3.4) selects i -th row in i -th step.

Remark 2.6.3 (Telling MATLAB about matrix properties).

MATLAB-\ assumes generic matrix, cannot detect special properties of (fully populated) matrix (e.g. symmetric, s.p.d., triangular).

Use `y = linsolve(A,b,opts)`

opts \in { LT \leftrightarrow A lower triangular matrix
 UT \leftrightarrow A upper triangular matrix
 UHESS \leftrightarrow A upper Hessenberg matrix
 SYM \leftrightarrow A Hermitian matrix
 POSDEF \leftrightarrow A positive definite matrix }

\triangle

2.7 Essential Skills Learned in Chapter 2

You should know:

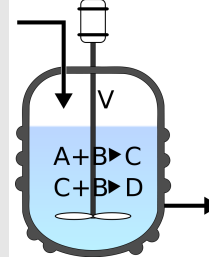
- how Gaussian elimination works in its block and LU-perspective
- particular types of matrices
- importance of pivoting
- when and how to solve directly linear systems in Matlab via `\` or `lu`
- when to use a direct solver and when not
- what is numerical stability and what to expect in case of Gaussian elimination
- what is residual, its size and impact in Gaussian elimination
- what is the condition number of a matrix and its importance
- what is a well-conditioned and an ill-conditioned problem
- what are sparse matrices and their importance
- how to create and use sparse matrices in matlab
- that several sparse matrix formats exist
- the impact of the `lu`-factorization or of pivoting on sparsity
- the impact of the `lu`-factorization on banded matrices
- when and why one should use the Cholesky factorization

3

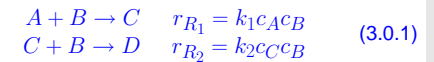
Iterative Methods for Non-Linear Systems of Equations

Example 3.0.1 (Continuous stirred-tank reactor (CSTR) in steady state).

Continuous stirred-tank reactor \triangleright



constant fluid volume V :



Assume spatially uniformity of each species:

$$\frac{d}{dt}(Vc_A) = v(c_{A,i} - c_A) + V(-k_1 c_A c_B) \quad (3.0.2)$$

$$\frac{d}{dt}(Vc_B) = v(c_{B,i} - c_B) + V(-k_1 c_A c_B - k_2 c_C c_B) \quad (3.0.3)$$

$$\frac{d}{dt}(Vc_C) = v(c_{C,i} - c_C) + V(k_1 c_A c_B - k_2 c_C c_B) \quad (3.0.4)$$

$$\frac{d}{dt}(Vc_D) = v(c_{D,i} - c_D) + V(k_2 c_C c_B) \quad (3.0.5)$$

v is the volumetric flow rate of the feed and outlet streams

V is the fixed reactor volume

c_j is the concentration of species j in the reactor

$c_{j,i}$ is the concentration of species j in the inlet stream

k_1, k_2 are the rate constants of each chemical reaction

steady state \Rightarrow time derivatives are zero

Non-linear system of equations

$$\begin{aligned} v(c_{A,i} - u_1) + V(-k_1 u_1 u_2) &= 0, \\ v(c_{B,i} - u_2) + V(-k_1 u_1 u_2 - k_2 u_3 u_2) &= 0, \\ v(c_{C,i} - u_3) + V(k_1 u_1 u_2 - k_2 u_3 u_2) &= 0, \\ v(c_{D,i} - u_4) + V(k_2 u_3 u_2) &= 0. \end{aligned} \quad (3.0.6)$$

4 equations \leftrightarrow 4 unknowns u_1, u_2, u_3, u_4

Formally:

$$(3.0.6) \iff F(\mathbf{u}) = 0$$



A **non-linear system of equations** is a concept almost *too abstract to be useful*, because it covers an extremely wide variety of problems. Nevertheless in this chapter we will mainly look at “generic” methods for such systems. This means that every method discussed may take a good deal of fine-tuning before it will really perform satisfactorily for a given non-linear system of equations.

Given:

$$\text{function } F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n, \quad n \in \mathbb{N}$$



Possible meaning:

Availability of a **procedure** function $y=F(x)$ evaluating F

Sought:

solution of **non-linear equation**

$$F(\mathbf{x}) = 0$$

Note: $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n \iff$ “same number of equations and unknowns”

In general no existence & uniqueness of solutions

3.1 Iterative methods

Remark 3.1.1 (Necessity of iterative approximation).

Gaussian elimination (\rightarrow Sect. 2.1) provides an algorithm that, if carried out in exact arithmetic, computes the solution of a linear system of equations with a **finite** number of elementary operations. However, linear systems of equations represent an exceptional case, because it is hardly ever possible to solve general systems of non-linear equations using only finitely many elementary operations. Certainly this is the case whenever irrational numbers are involved.



An **iterative method** for (approximately) solving the non-linear equation $F(\mathbf{x}) = 0$ is an algorithm generating a sequence $(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$ of **approximate solutions**.

Initial guess \rightarrow

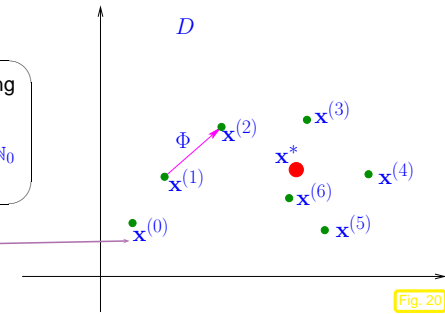


Fig. 20

Fundamental concepts: **convergence** \Rightarrow **speed of convergence**
consistency

- iterate $\mathbf{x}^{(k)}$ depends on F and (one or several) $\mathbf{x}^{(n)}$, $n < k$, e.g.,

$$\mathbf{x}^{(k)} = \underbrace{\Phi_F(\mathbf{x}^{(k-1)}, \dots, \mathbf{x}^{(k-m)})}_{\text{iteration function for } m\text{-point method}} \quad (3.1.1)$$

3.1
p. 137

- $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(m-1)}$ = **initial guess(es)** (*ger.:* Anfangsnäherung)

3.1
p. 13

Definition 3.1.1 (Convergence of iterative methods).

An iterative method **converges** (for fixed initial guess(es)) $\iff \mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ and $F(\mathbf{x}^*) = 0$.

Definition 3.1.2 (Consistency of iterative methods).

An iterative method is **consistent** with $F(\mathbf{x}) = 0$

$$\iff \Phi_F(\mathbf{x}^*, \dots, \mathbf{x}^*) = \mathbf{x}^* \iff F(\mathbf{x}^*) = 0$$

3.1
p. 138

Terminology:

error of iterates $\mathbf{x}^{(k)}$ is defined as: $\mathbf{e}^{(k)} := \mathbf{x}^{(k)} - \mathbf{x}^*$

3.1
p. 14

Definition 3.1.3 (Local and global convergence).

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

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

for the sequences $(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$ of iterates.

If $U = D$, the iterative method is **globally convergent**.

local convergence



(Only initial guesses "sufficiently close" to \mathbf{x}^* guarantee convergence.)

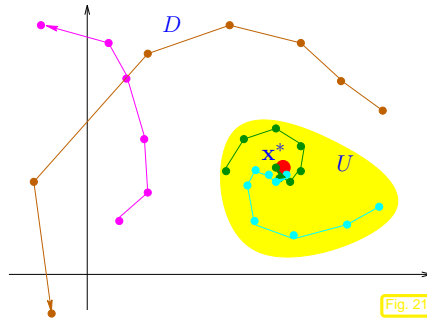


Fig. 21

Goal: Find iterative methods that converge (locally) to a solution of $F(\mathbf{x}) = 0$.

Two general questions: How to measure the speed of convergence?
When to terminate the iteration?

3.1.1 Speed of convergence

Here and in the sequel, $\|\cdot\|$ designates a generic vector norm, see Def. 2.4.1. Any occurring matrix norm is induced by this vector norm, see Def. 2.4.2.

It is important to be aware which statements depend on the choice of norm and which do not!

"Speed of convergence" \leftrightarrow decrease of norm (see Def. 2.4.1) of iteration error

Definition 3.1.4 (Linear convergence).

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

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

Terminology: least upper bound for L gives the **rate of convergence**

Remark 3.1.2 (Impact of choice of norm).

Fact of convergence of iteration is **independent** of choice of norm
 Fact of linear convergence **depends** on choice of norm
 Rate of linear convergence **depends** on choice of norm

Recall: equivalence of all norms on finite dimensional vector space \mathbb{K}^n :

Definition 3.1.5 (Equivalence of norms).

Two norms $\|\cdot\|_1$ and $\|\cdot\|_2$ on a vector space V are equivalent if

$$\exists \underline{C}, \bar{C} > 0: \underline{C} \|v\|_1 \leq \|v\|_2 \leq \bar{C} \|v\|_1 \quad \forall v \in V.$$

3.1
p. 141

Theorem 3.1.6 (Equivalence of all norms on finite dimensional vector spaces).
 If $\dim V < \infty$ all norms (\rightarrow Def. 2.4.1) on V are equivalent (\rightarrow Def. 3.1.5).

3.1
p. 14

Remark 3.1.3 (Seeing linear convergence).

norms of iteration errors
 \updownarrow
 \sim straight line in **lin-log** plot

$$\|e^{(k)}\| \leq L^k \|e^{(0)}\|, \quad \log(\|e^{(k)}\|) \leq k \log(L) + \log(\|e^{(0)}\|).$$

(• Any "faster" convergence also qualifies as linear!)

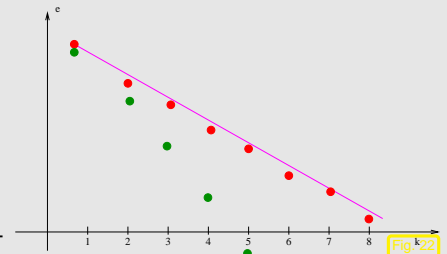


Fig. 22

3.1
p. 142

Let us abbreviate the error norm in step k by $\epsilon_k := \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$. In the case of linear convergence

3.1
p. 14

(see Def. 3.1.4) assume (with $0 < L < 1$)

$$\epsilon_{k+1} \approx L\epsilon_k \Rightarrow \log \epsilon_{k+1} \approx \log L + \log \epsilon_k \Rightarrow \log \epsilon_{k+1} \approx k \log L + \log \epsilon_0. \quad (3.1.2)$$

We conclude that $\log L < 0$ describes slope of graph in lin-log error chart.

Example 3.1.4 (Linearly convergent iteration).

Iteration ($n = 1$):

$$x^{(k+1)} = x^{(k)} + \frac{\cos x^{(k)} + 1}{\sin x^{(k)}}.$$

Code 3.1.5: simple fixed point iteration

```

1 y = [ ];
2 for i = 1:15
3   x = x + (cos(x)+1)/sin(x);
4   y = [y,x];
5 end
6 err = y - x;
7 rate = err(2:15) ./ err(1:14);

```

Note: $x^{(15)}$ replaces the exact solution x^* in the computation of the rate of convergence.

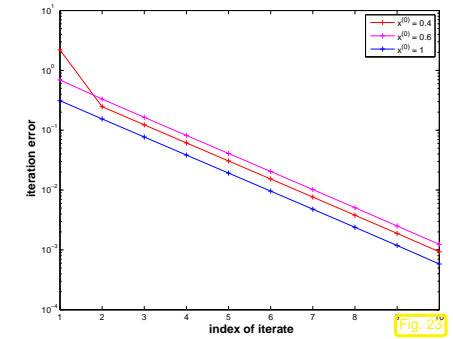
k	$x^{(0)} = 0.4$		$x^{(0)} = 0.6$		$x^{(0)} = 1$	
	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$
2	3.3887	0.1128	3.4727	0.4791	2.9873	0.4959
3	3.2645	0.4974	3.3056	0.4953	3.0646	0.4989
4	3.2030	0.4992	3.2234	0.4988	3.1031	0.4996
5	3.1723	0.4996	3.1825	0.4995	3.1224	0.4997
6	3.1569	0.4995	3.1620	0.4994	3.1320	0.4995
7	3.1493	0.4990	3.1518	0.4990	3.1368	0.4990
8	3.1454	0.4980	3.1467	0.4980	3.1392	0.4980

Rate of convergence ≈ 0.5

Linear convergence as in Def. 3.1.4



error graphs = straight lines in lin-log scale
→ Rem. 3.1.3



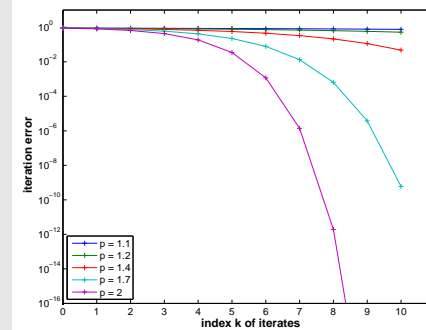
Definition 3.1.7 (Order of convergence).

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

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

with $C < 1$ for $p = 1$ (linear convergence → Def. 3.1.4)

3.1
p. 145



◁ Qualitative error graphs for convergence of order p (lin-log scale)

In the case of convergence of order p ($p > 1$) (see Def. 3.1.7):

$$\begin{aligned} \epsilon_{k+1} \approx C\epsilon_k^p &\Rightarrow \log \epsilon_{k+1} = \log C + p \log \epsilon_k \Rightarrow \log \epsilon_{k+1} = \log C \sum_{l=0}^k p^l + p^{k+1} \log \epsilon_0 \\ &\Rightarrow \log \epsilon_{k+1} = -\frac{\log C}{p-1} + \left(\frac{\log C}{p-1} + \log \epsilon_0\right) p^{k+1}. \end{aligned}$$

In this case, the error graph is a concave power curve (for sufficiently small ϵ_0 !)

3.1
p. 146

Example 3.1.6 (quadratic convergence). (= convergence of order 2)

3.1
p. 14

3.1
p. 14

Iteration for computing \sqrt{a} , $a > 0$:

$$x^{(k+1)} = \frac{1}{2}\left(x^{(k)} + \frac{a}{x^{(k)}}\right) \Rightarrow |x^{(k+1)} - \sqrt{a}| = \frac{1}{2x^{(k)}}|x^{(k)} - \sqrt{a}|^2. \quad (3.1.3)$$

By the arithmetic-geometric mean inequality (AGM) $\sqrt{ab} \leq \frac{1}{2}(a+b)$ we conclude: $x^{(k)} > \sqrt{a}$ for $k \geq 1$.

⇒ sequence from (3.1.3) converges with order 2 to \sqrt{a}

Note: $x^{(k+1)} < x^{(k)}$ for all $k \geq 2$ ⇒ $(x^{(k)})_{k \in \mathbb{N}_0}$ converges as a decreasing sequence that is bounded from below (→ analysis course)

How to guess the order of convergence in a numerical experiment?

Abbreviate $\epsilon_k := \|x^{(k)} - x^*\|$ and then

$$\epsilon_{k+1} \approx C \epsilon_k^p \Rightarrow \log \epsilon_{k+1} \approx \log C + p \log \epsilon_k \Rightarrow \frac{\log \epsilon_{k+1} - \log \epsilon_k}{\log \epsilon_k - \log \epsilon_{k-1}} \approx p.$$

Numerical experiment: iterates for $a = 2$:

k	$x^{(k)}$	$e^{(k)} := x^{(k)} - \sqrt{2}$	$\log \frac{ e^{(k)} }{ e^{(k-1)} } : \log \frac{ e^{(k-1)} }{ e^{(k-2)} }$
0	2.0000000000000000	0.58578643762690485	
1	1.5000000000000000	0.08578643762690485	
2	1.4166666666666665	0.00245310429357137	1.850
3	1.41421568627450966	0.00000212390141452	1.984
4	1.41421356237468987	0.0000000000159472	2.000
5	1.41421356237309492	0.0000000000000022	0.630

Note the **doubling** of the number of significant digits in each step ! [impact of roundoff !]

The doubling of the number of significant digits for the iterates holds true for any convergent second-order iteration:

Indeed, denoting the relative error in step k by δ_k , we have:

$$\begin{aligned} x^{(k)} &= x^*(1 + \delta_k) \Rightarrow x^{(k)} - x^* = \delta_k x^* \\ \Rightarrow |x^* \delta_{k+1}| &= |x^{(k+1)} - x^*| \leq C|x^{(k)} - x^*|^2 = C|x^* \delta_k|^2 \\ &\Rightarrow |\delta_{k+1}| \leq C|x^*| \delta_k^2. \end{aligned} \quad (3.1.4)$$

Note: $\delta_k \approx 10^{-\ell}$ means that $x^{(k)}$ has ℓ significant digits.

Also note that if $C \approx 1$, then $\delta_k = 10^{-\ell}$ and (3.1.6) implies $\delta_{k+1} \approx 10^{-2\ell}$.



3.1.2 Termination criteria

Usually (even without roundoff errors) the iteration will never arrive at an/the exact solution x^* after finitely many steps. Thus, we can only hope to compute an *approximate* solution by accepting $x^{(K)}$ as result for some $K \in \mathbb{N}_0$. Termination criteria (*ger.*: Abbruchbedingungen) are used to determine a suitable value for K .

For the sake of efficiency: ▷ stop iteration when iteration error is just “small enough”

“small enough” depends on concrete setting:

Usual goal: $\|x^{(K)} - x^*\| \leq \tau$, $\tau \hat{=}$ prescribed **tolerance**.

$$\text{Ideal: } K = \operatorname{argmin}\{k \in \mathbb{N}_0 : \|x^{(k)} - x^*\| < \tau\}. \quad (3.1.5)$$

3.1
p. 149

3.1
p. 15

① **A priori termination:** stop iteration after fixed number of steps (possibly depending on $x^{(0)}$).



Drawback: hardly ever possible !

Alternative:

A posteriori termination criteria

use already computed iterates to decide when to stop

②

Reliable termination: stop iteration $\{x^{(k)}\}_{k \in \mathbb{N}_0}$ with limit x^* , when

$$\|x^{(k)} - x^*\| \leq \tau, \quad \tau \hat{=} \text{prescribed tolerance} > 0. \quad (3.1.6)$$



x^* not known !

Invoking additional properties of either the non-linear system of equations $F(x) = 0$ or the iteration it is sometimes possible to tell that for sure $\|x^{(k)} - x^*\| \leq \tau$ for all $k \geq K$, though this K may be (significantly) larger than the optimal termination index from (3.1.5), see Rem. 3.1.8.

3.1
p. 150

3.1
p. 15

③ use that M is finite! (→ Sect. ??)

➤ possible to wait until (convergent) iteration becomes stationary

☹ possibly grossly inefficient!
(always computes “up to machine precision”)

Code 3.1.7: stationary iteration in M , → Ex. 3.1.6

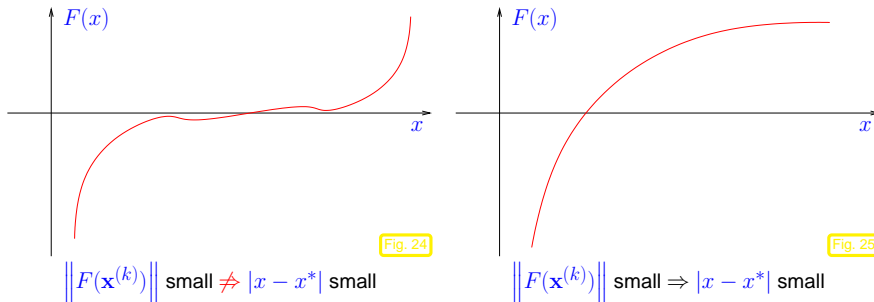
```

1 function x = sqrtit(a)
2   x_old = -1; x = a;
3   while (x_old ~= x)
4     x_old = x;
5     x = 0.5*(x+a/x);
6   end
  
```

④ Residual based termination: stop convergent iteration $\{x^{(k)}\}_{k \in \mathbb{N}_0}$, when

$$\|F(x^{(k)})\| \leq \tau, \quad \tau \hat{=} \text{prescribed tolerance} > 0.$$

☹ no guaranteed accuracy



Sometimes extra knowledge about the type/speed of convergence allows to achieve **reliable termination** in the sense that (3.1.6) can be guaranteed though the number of iterations might be (slightly) too large.

Remark 3.1.8 (A posteriori termination criterion for linearly convergent iterations).

Known: iteration linearly convergent with rate of convergence $0 < L < 1$:

Derivation of a posteriori termination criterion for linearly convergent iterations with rate of convergence $0 < L < 1$:

$$\|x^{(k)} - x^*\| \stackrel{\Delta\text{-inequ.}}{\leq} \|x^{(k+1)} - x^{(k)}\| + \|x^{(k+1)} - x^*\| \leq \|x^{(k+1)} - x^{(k)}\| + L \|x^{(k)} - x^*\|.$$

Iterates satisfy: $\|x^{(k+1)} - x^*\| \leq \frac{L}{1-L} \|x^{(k+1)} - x^{(k)}\|$. (3.1.7)

This suggests that we take the right hand side of (3.1.7) as a posteriori error bound.



Example 3.1.9 (A posteriori error bound for linearly convergent iteration).

Iteration of Example 3.1.4:

$$x^{(k+1)} = x^{(k)} + \frac{\cos x^{(k)} + 1}{\sin x^{(k)}} \Rightarrow x^{(k)} \rightarrow \pi \text{ for } x^{(0)} \text{ close to } \pi.$$

3.1

p. 153

Observed rate of convergence: $L = 1/2$

Error and error bound for $x^{(0)} = 0.4$:

k	$ x^{(k)} - \pi $	$\frac{L}{1-L} x^{(k)} - x^{(k-1)} $	slack of bound
1	2.191562221997101	4.933154875586894	2.741592653589793
2	0.247139097781070	1.944423124216031	1.697284026434961
3	0.122936737876834	0.124202359904236	0.001265622027401
4	0.061390835206217	0.061545902670618	0.000155067464401
5	0.030685773472263	0.030705061733954	0.000019288261691
6	0.015341682696235	0.015344090776028	0.000002408079792
7	0.007670690889185	0.007670991807050	0.000000300917864
8	0.003835326638666	0.003835364250520	0.000000037611854
9	0.001917660968637	0.001917665670029	0.000000004701392
10	0.000958830190489	0.000958830778147	0.000000000587658
11	0.000479415058549	0.000479415131941	0.00000000073392
12	0.000239707524646	0.000239707533903	0.00000000009257
13	0.000119853761949	0.000119853762696	0.00000000000747
14	0.000059926881308	0.000059926880641	0.00000000000667
15	0.000029963440745	0.000029963440563	0.00000000000181

3.1

p. 154

Hence: the a posteriori error bound is highly accurate in this case!

3.1

p. 15



Note: If L not known then using $\tilde{L} > L$ in error bound is playing safe.

3.2 Fixed Point Iterations

Non-linear system of equations $F(\mathbf{x}) = 0$, $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$,

A **fixed point iteration** is defined by **iteration function** $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$:

iteration function $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$
 initial guess $\mathbf{x}^{(0)} \in U$ \rightarrow iterates $(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$: $\mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})$
 \rightarrow 1-point method, cf. (3.1.1)

Sequence of iterates need not be well defined: $\mathbf{x}^{(k)} \notin U$ possible !

3.2.1 Consistent fixed point iterations

Definition 3.2.1 (Consistency of fixed point iterations, cf. Def. 3.1.2).
 A fixed point iteration $\mathbf{x}^{(k+1)} = \Phi(\mathbf{x}^{(k)})$ is **consistent** with $F(\mathbf{x}) = 0$, if

$$F(\mathbf{x}) = 0 \text{ and } x \in U \cap D \Leftrightarrow \Phi(\mathbf{x}) = \mathbf{x} .$$

Note: Φ continuous & fixed point iteration (locally) convergent to \mathbf{x}^* then \mathbf{x}^* is fixed point of iteration function Φ .

General construction of fixed point iterations that is consistent with $F(\mathbf{x}) = 0$:
 rewrite $F(\mathbf{x}) = 0 \Leftrightarrow \Phi(\mathbf{x}) = \mathbf{x}$ and then

use the **fixed point iteration** $\mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})$. (3.2.1)

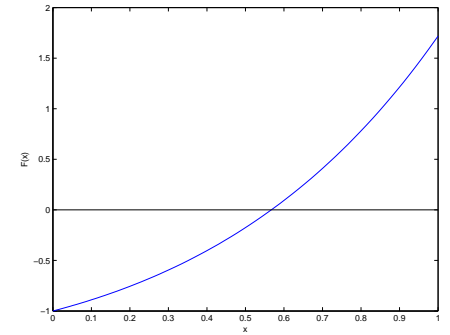
Note: there are *many ways* to transform $F(\mathbf{x}) = 0$ into a fixed point form !

Example 3.2.1 (Options for fixed point iterations).

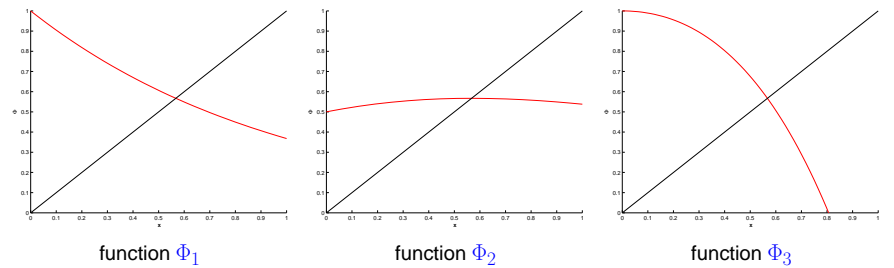
$$F(x) = xe^x - 1, \quad x \in [0, 1] .$$

Different fixed point forms:

$$\begin{aligned} \Phi_1(x) &= e^{-x} , \\ \Phi_2(x) &= \frac{1+x}{1+e^x} , \\ \Phi_3(x) &= x+1-xe^x . \end{aligned}$$



3.2
p. 157



3.2
p. 15

k	$x^{(k+1)} := \Phi_1(x^{(k)})$	$x^{(k+1)} := \Phi_2(x^{(k)})$	$x^{(k+1)} := \Phi_3(x^{(k)})$
0	0.5000000000000000	0.5000000000000000	0.5000000000000000
1	0.606530659712633	0.566311003197218	0.675639364649936
2	0.545239211892605	0.567143165034862	0.347812678511202
3	0.579703094878068	0.567143290409781	0.855321409174107
4	0.560064627938902	0.567143290409784	-0.156505955383169
5	0.571172148977215	0.567143290409784	0.977326422747719
6	0.564862946980323	0.567143290409784	-0.619764251895580
7	0.568438047570066	0.567143290409784	0.713713087416146
8	0.566409452746921	0.567143290409784	0.256626649129847
9	0.567559634262242	0.567143290409784	0.924920676910549
10	0.566907212935471	0.567143290409784	-0.407422405542253

3.2
p. 158

3.2
p. 16

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	0.000000000000003	0.288178118764323
4	0.007078662470882	0.000000000000000	0.723649245792953
5	0.004028858567431	0.000000000000000	0.410183132337935
6	0.002280343429460	0.000000000000000	1.186907542305364
7	0.001294757160282	0.000000000000000	0.146569797006362
8	0.000733837662863	0.000000000000000	0.310516641279937
9	0.000416343852458	0.000000000000000	0.357777386500765
10	0.000236077474313	0.000000000000000	0.974565695952037

Observed: linear convergence of $x_1^{(k)}$, quadratic convergence of $x_2^{(k)}$,
no convergence (erratic behavior) of $x_3^{(k)}$, $x_i^{(0)} = 0.5$.

Question: can we explain/forecast the behaviour of the iteration?

3.2.2 Convergence of fixed point iterations

In this section we will try to find easily verifiable conditions that ensure convergence (of a certain order) of fixed point iterations. It will turn out that these conditions are surprisingly simple and general.

Definition 3.2.2 (Contractive mapping).

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

$$\exists L < 1: \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in U. \quad (3.2.2)$$

A simple consideration: if $\Phi(\mathbf{x}^*) = \mathbf{x}^*$ (fixed point), then a fixed point iteration induced by a contractive mapping Φ satisfies

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = \|\Phi(\mathbf{x}^{(k)}) - \Phi(\mathbf{x}^*)\| \stackrel{(3.2.2)}{\leq} L \|\mathbf{x}^{(k)} - \mathbf{x}^*\|,$$

that is, the iteration **converges** (at least) **linearly** (\rightarrow Def. 3.1.4).

Remark 3.2.2 (Banach's fixed point theorem). \rightarrow [47, Satz 6.5.2]

A key theorem in calculus (also functional analysis):

Theorem 3.2.3 (Banach's fixed point theorem).

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

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

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

Proof. Proof based on 1-point iteration $\mathbf{x}^{(k)} = \Phi(\mathbf{x}^{(k-1)})$, $\mathbf{x}^{(0)} \in D$:

$$\begin{aligned} \|\mathbf{x}^{(k+N)} - \mathbf{x}^{(k)}\| &\leq \sum_{j=k}^{k+N-1} \|\mathbf{x}^{(j+1)} - \mathbf{x}^{(j)}\| \leq \sum_{j=k}^{k+N-1} L^j \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \\ &\leq \frac{L^k}{1-L} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \xrightarrow{k \rightarrow \infty} 0. \end{aligned}$$

$(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$ Cauchy sequence \rightarrow convergent $\mathbf{x}^{(k)} \xrightarrow{k \rightarrow \infty} \mathbf{x}^*$.

Continuity of $\Phi \rightarrow \Phi(\mathbf{x}^*) = \mathbf{x}^*$. Uniqueness of fixed point is evident. \square

3.2

p. 161

A simple criterion for a differentiable Φ to be contractive:

Lemma 3.2.4 (Sufficient condition for local linear convergence of fixed point iteration).

If $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$, $\Phi(\mathbf{x}^*) = \mathbf{x}^*$, Φ differentiable in \mathbf{x}^* , and $\|D\Phi(\mathbf{x}^*)\| < 1$, then the fixed point iteration (3.2.1) converges locally and at least linearly.

matrix norm, Def. 2.4.2!

notation: $D\Phi(\mathbf{x}) \hat{=}$ **Jacobian** (ger.: Jacobi-Matrix) of Φ at $\mathbf{x} \in D$
 \rightarrow [47, Sect. 7.6]

Example 3.2.3 (Fixed point iteration in 1D).

1D setting ($n = 1$): $\Phi : \mathbb{R} \mapsto \mathbb{R}$ continuously differentiable, $\Phi(x^*) = x^*$

fixed point iteration: $x^{(k+1)} = \Phi(x^{(k)})$

"Visualization" of the statement of Lemma 3.2.4: The iteration converges *locally*, if Φ is flat in a neighborhood of x^* , it will diverge, if Φ is steep there.

3.2

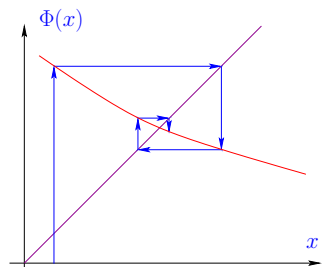
p. 162

3.2

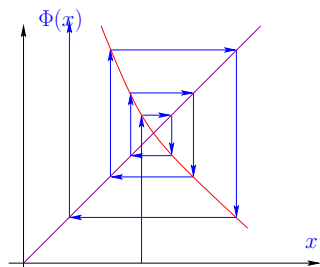
p. 16

3.2

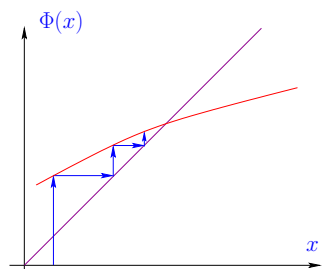
p. 16



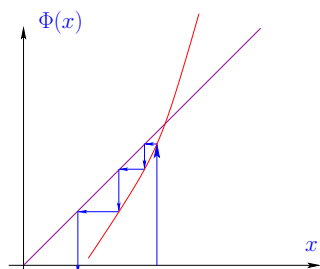
$-1 < \Phi'(x^*) \leq 0 \rightarrow$ convergence



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



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



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

Proof. (of Lemma 3.2.4) By definition of derivative

$$\|\Phi(\mathbf{y}) - \Phi(\mathbf{x}^*) - D\Phi(\mathbf{x}^*)(\mathbf{y} - \mathbf{x}^*)\| \leq \psi(\|\mathbf{y} - \mathbf{x}^*\|) \|\mathbf{y} - \mathbf{x}^*\|,$$

with $\psi : \mathbb{R}_0^+ \mapsto \mathbb{R}_0^+$ satisfying $\lim_{t \rightarrow 0} \psi(t) = 0$.

Choose $\delta > 0$ such that

$$L := \psi(t) + \|D\Phi(\mathbf{x}^*)\| \leq \frac{1}{2}(1 + \|D\Phi(\mathbf{x}^*)\|) < 1 \quad \forall 0 \leq t < \delta.$$

By inverse triangle inequality we obtain for fixed point iteration

$$\begin{aligned} \|\Phi(\mathbf{x}) - \mathbf{x}^*\| - \|D\Phi(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*)\| &\leq \psi(\|\mathbf{x} - \mathbf{x}^*\|) \|\mathbf{x} - \mathbf{x}^*\| \\ \Rightarrow \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| &\leq (\psi(t) + \|D\Phi(\mathbf{x}^*)\|) \|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq L \|\mathbf{x}^{(k)} - \mathbf{x}^*\|, \end{aligned}$$

if $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| < \delta$. □

Contractivity also guarantees the *uniqueness* of a fixed point, see the next Lemma.

Recalling the Banach fixed point theorem Thm. 3.2.3 we see that under some additional (usually mild) assumptions, it also ensures the *existence* of a fixed point.

Lemma 3.2.5 (Sufficient condition for local linear convergence of fixed point iteration (II)).

Let U be convex and $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ continuously differentiable with $L := \sup_{\mathbf{x} \in U} \|D\Phi(\mathbf{x})\| < 1$.

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

Recall: $U \subset \mathbb{R}^n$ convex $\Leftrightarrow (t\mathbf{x} + (1-t)\mathbf{y}) \in U$ for all $\mathbf{x}, \mathbf{y} \in U, 0 \leq t \leq 1$

Proof. (of Lemma 3.2.5) By the mean value theorem

$$\begin{aligned} \Phi(\mathbf{x}) - \Phi(\mathbf{y}) &= \int_0^1 D\Phi(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}) \, d\tau \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(\Phi). \\ \Rightarrow \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\| &\leq L \|\mathbf{y} - \mathbf{x}\|. \end{aligned}$$

There is $\delta > 0$: $B := \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \delta\} \subset \text{dom}(\Phi)$. Φ is contractive on B with unique fixed point \mathbf{x}^* .

Remark 3.2.4.

If $0 < \|D\Phi(\mathbf{x}^*)\| < 1$, $\mathbf{x}^{(k)} \approx \mathbf{x}^*$ then the asymptotic rate of linear convergence is $L = \|D\Phi(\mathbf{x}^*)\|$

3.2
p. 165

3.2
p. 16

Example 3.2.5 (Multidimensional fixed point iteration).

$$\begin{aligned} \text{System of equations} & \quad \text{in} & \quad \text{fixed point form:} \\ \begin{cases} x_1 - c(\cos x_1 - \sin x_2) = 0 \\ (x_1 - x_2) - c \sin x_2 = 0 \end{cases} & \Rightarrow & \begin{cases} c(\cos x_1 - \sin x_2) = x_1 \\ c(\cos x_1 - 2 \sin x_2) = x_2 \end{cases} \\ \text{Define: } \Phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = c \begin{pmatrix} \cos x_1 - \sin x_2 \\ \cos x_1 - 2 \sin x_2 \end{pmatrix} & \Rightarrow & D\Phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = -c \begin{pmatrix} \sin x_1 & \cos x_2 \\ \sin x_1 & 2 \cos x_2 \end{pmatrix}. \end{aligned}$$

Choose appropriate norm: $\|\cdot\| = \infty$ -norm $\|\cdot\|_\infty$ (\rightarrow Example 2.4.1);

$$\text{if } c < \frac{1}{3} \Rightarrow \|D\Phi(\mathbf{x})\|_\infty < 1 \quad \forall \mathbf{x} \in \mathbb{R}^2,$$

\Rightarrow (at least) linear convergence of the fixed point iteration.

The existence of a fixed point is also guaranteed, because Φ maps into the closed set $[-3, 3]^2$. Thus, the Banach fixed point theorem, Thm. 3.2.3, can be applied.

What about higher order convergence (\rightarrow Def. 3.1.7) ?

3.2
p. 166

Refined convergence analysis for $n = 1$ (scalar case, $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$):

3.2
p. 16

Theorem 3.2.6 (Taylor's formula). → [47, Sect. 5.5]

If $\Phi : U \subset \mathbb{R} \mapsto \mathbb{R}$, U interval, is $m + 1$ times continuously differentiable, $x \in U$

$$\Phi(y) - \Phi(x) = \sum_{k=1}^m \frac{1}{k!} \Phi^{(k)}(x)(y-x)^k + O(|y-x|^{m+1}) \quad \forall y \in U. \quad (3.2.3)$$

Apply Taylor expansion (3.2.3) to iteration function Φ :

If $\Phi(x^*) = x^*$ and $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$ is "sufficiently smooth"

$$x^{(k+1)} - x^* = \Phi(x^{(k)}) - \Phi(x^*) = \sum_{l=1}^m \frac{1}{l!} \Phi^{(l)}(x^*)(x^{(k)} - x^*)^l + O(|x^{(k)} - x^*|^{m+1}). \quad (3.2.4)$$

Lemma 3.2.7 (Higher order local convergence of fixed point iterations).

If $\Phi : U \subset \mathbb{R} \mapsto \mathbb{R}$ is $m + 1$ times continuously differentiable, $\Phi(x^*) = x^*$ for some x^* in the interior of U , and $\Phi^{(l)}(x^*) = 0$ for $l = 1, \dots, m$, $m \geq 1$, then the fixed point iteration (3.2.1) converges locally to x^* with **order** $\geq m + 1$ (→ Def. 3.1.7).

Proof. For neighborhood \mathcal{U} of x^*

$$(3.2.4) \Rightarrow \exists C > 0: |\Phi(y) - \Phi(x^*)| \leq C |y - x^*|^{m+1} \quad \forall y \in \mathcal{U}.$$

$$\delta^m C < 1/2: |x^{(0)} - x^*| < \delta \Rightarrow |x^{(k)} - x^*| < 2^{-k} \delta \quad \triangleright \text{local convergence.}$$

Then appeal to (3.2.4) □

Example 3.2.1 continued:

$$\Phi_2'(x) = \frac{1 - xe^x}{(1 + e^x)^2} = 0 \quad , \text{ if } xe^x - 1 = 0 \quad \text{hence quadratic convergence ! .}$$

Example 3.2.1 continued: Since $x^* e^{x^*} - 1 = 0$

$$\Phi_1'(x) = -e^{-x} \Rightarrow \Phi_1'(x^*) = -x^* \approx -0.56 \quad \text{hence local linear convergence .}$$

$$\Phi_3'(x) = 1 - xe^x - e^x \Rightarrow \Phi_3'(x^*) = -\frac{1}{x^*} \approx -1.79 \quad \text{hence no convergence .}$$

Remark 3.2.6 (Termination criterion for contractive fixed point iteration).

Recap of Rem. 3.1.8:

Termination criterion for contractive fixed point iteration, c.f. (3.2.3), with contraction factor $0 \leq L < 1$:

$$\begin{aligned} \|\mathbf{x}^{(k+m)} - \mathbf{x}^{(k)}\| &\stackrel{\Delta\text{-ineq.}}{\leq} \sum_{j=k}^{k+m-1} \|\mathbf{x}^{(j+1)} - \mathbf{x}^{(j)}\| \leq \sum_{j=k}^{k+m-1} L^{j-k} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \\ &= \frac{1 - L^m}{1 - L} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \frac{1 - L^m}{1 - L} L^{k-l} \|\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)}\|. \end{aligned}$$

hence for $m \rightarrow \infty$, with $\mathbf{x}^* := \lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$:

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{L^{k-l}}{1 - L} \|\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)}\|. \quad (3.2.5)$$



Set $l = 0$ in (3.2.5)

a priori termination criterion

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{L^k}{1 - L} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \quad (3.2.6)$$

Set $l = k - 1$ in (3.2.5)

a posteriori termination criterion

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{L}{1 - L} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \quad (3.2.7)$$

△

3.3 Zero Finding

Now, focus on scalar case $n = 1$:

$$F : I \subset \mathbb{R} \mapsto \mathbb{R} \text{ continuous, } I \text{ interval}$$

Sought:

$$x^* \in I: \quad F(x^*) = 0$$

3.3.1 Bisection

Idea: use ordering of real numbers & intermediate value theorem

Input: $a, b \in I$ such that $F(a)F(b) < 0$
(different signs !)

$\Rightarrow \exists x^* \in]\min\{a, b\}, \max\{a, b\}[$:
 $F(x^*) = 0$.

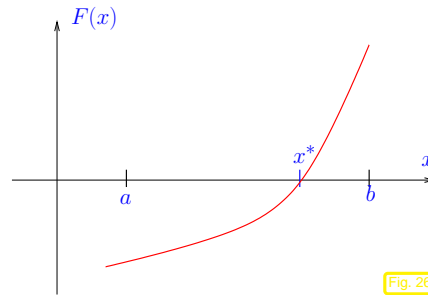


Fig. 26

Algorithm 3.3.1 (Bisection method).

MATLAB-CODE: bisection method

```
function x = bisect(F, a, b, tol)
% Searching zero by bisection
if (a > b), t=a; a=b; b=t; end;
fa = F(a); fb = F(b);
if (fa*fb > 0)
    error('f(a), f(b) same sign'); end;
if (fa > 0), v=-1; else v = 1; end
x = 0.5*(b+a);
while((b-a > tol) & ((a < x) & (x < b)))
    if (v*F(x) > 0), b=x; else a=x; end;
    x = 0.5*(a+b);
end
```

$\text{tol} \hat{=}$ absolute tolerance

Handle to MATLAB function providing F .

Avoid infinite loop, if $\text{tol} <$ resolution of \mathbb{M} at zero x^* ("M-based termination criterion").

This is an example for an algorithm that (in the case of $\text{tol}=0$) uses the properties of machine arithmetic to define an a posteriori termination criterion, see Sect. 3.1.2. The iteration will terminate, when, e.g., $a + \frac{1}{2}(b-a) = a$, which, by the Ass. ?? can only happen, when

$$\left| \frac{1}{2}(b-a) \right| \leq \text{eps} \cdot |a|.$$

Since the exact zero is located between a and b , this condition implies a relative error $\leq \text{eps}$ of the computed zero.



Advantages:

- "foolproof"
- requires only F evaluations



Drawbacks: Merely linear convergence: $|x^{(k)} - x^*| \leq 2^{-k}|b-a|$
 $\blacktriangleright \log_2 \left(\frac{|b-a|}{\text{tol}} \right)$ steps necessary

Remark 3.3.2. MATLAB function `fzero` is based on bisection approach.



3.3.2 Model function methods

$\hat{=}$ class of iterative methods for finding zeroes of F :



Idea: Given: approximate zeroes $x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}$

- 1 replace F with model function \tilde{F}
(using function values/derivative values in $x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}$)
- 2 $x^{(k+1)} :=$ zero of \tilde{F}
(has to be readily available \leftrightarrow analytic formula)

Distinguish (see (3.1.1)):

one-point methods: $x^{(k+1)} = \Phi_F(x^{(k)}), k \in \mathbb{N}$ (e.g., fixed point iteration \rightarrow Sect. 3.2)
 multi-point methods: $x^{(k+1)} = \Phi_F(x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}), k \in \mathbb{N}, m = 2, 3, \dots$

3.3.2.1 Newton method in scalar case

Assume: $F : I \mapsto \mathbb{R}$ continuously differentiable

model function := tangent at F in $x^{(k)}$:

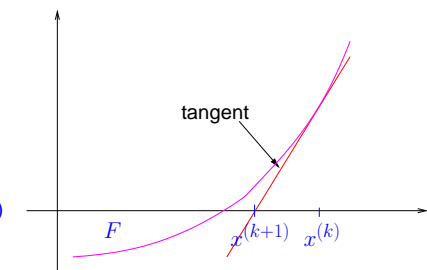
$$\tilde{F}(x) := F(x^{(k)}) + F'(x^{(k)})(x - x^{(k)})$$

take $x^{(k+1)} :=$ zero of tangent

We obtain **Newton iteration**

$$x^{(k+1)} := x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})}, \quad (3.3.1)$$

that requires $F'(x^{(k)}) \neq 0$.



Example 3.3.3 (Newton method in 1D). (→ Ex. 3.2.1)

Newton iterations for two different scalar non-linear equation with the same solution sets:

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

$$F(x) = x - e^{-x} \Rightarrow F'(x) = 1 + e^{-x} \Rightarrow x^{(k+1)} = x^{(k)} - \frac{x^{(k)} - e^{-x^{(k)}}}{1 + e^{-x^{(k)}}} = \frac{1 + x^{(k)}}{1 + e^{x^{(k)}}}$$

Ex. 3.2.1 shows quadratic convergence! (→ Def. 3.1.7) \diamond

Newton iteration (3.3.1) $\hat{=}$ fixed point iteration (→ Sect.3.2) with iteration function

$$\Phi(x) = x - \frac{F(x)}{F'(x)} \Rightarrow \Phi'(x) = \frac{F(x)F''(x)}{(F'(x))^2} \Rightarrow \Phi'(x^*) = 0, \text{ , if } F(x^*) = 0, F'(x^*) \neq 0.$$

From Lemma 3.2.7:

Newton method locally quadratically convergent (→ Def. 3.1.7) to zero x^* , if $F'(x^*) \neq 0$

3.3.2.2 Special one-point methods

Idea underlying other one-point methods: non-linear local approximation

Useful, if *a priori knowledge* about the structure of F (e.g. about F being a rational function, see below) is available. This is often the case, because many problems of 1D zero finding are posed for functions given in analytic form with a few parameters.

Prerequisite: Smoothness of F : $F \in C^m(I)$ for some $m > 1$

Example 3.3.4 (Halley's iteration).

Given $x^{(k)} \in I$, next iterate := zero of model function: $h(x^{(k+1)}) = 0$, where

$$h(x) := \frac{a}{x+b} + c \text{ (rational function) such that } F^{(j)}(x^{(k)}) = h^{(j)}(x^{(k)}), \quad j = 0, 1, 2.$$

$$\frac{a}{x^{(k)} + b} + c = F(x^{(k)}), \quad -\frac{a}{(x^{(k)} + b)^2} = F'(x^{(k)}), \quad \frac{2a}{(x^{(k)} + b)^3} = F''(x^{(k)}).$$

$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})} \cdot \frac{1}{1 - \frac{1}{2} \frac{F(x^{(k)})F''(x^{(k)})}{(F'(x^{(k)})^2)}}.$$

$$F(x) = \frac{1}{(x+1)^2} + \frac{1}{(x+0.1)^2} - 1, \quad x > 0: \quad x^{(0)} = 0$$

k	$x^{(k)}$	$F(x^{(k)})$	$x^{(k)} - x^{(k-1)}$	$x^{(k)} - x^*$
				0.0
				0.000
				0.00000000

3.3
p. 177

3.3.1

3.3
p. 17

k	$x^{(k)}$	$F(x^{(k)})$	$x^{(k)} - x^{(k-1)}$	$x^{(k)} - x^*$
1	0.04995004995005	44.38117504792020	-0.04995004995005	-0.99625244494443
2	0.12455117953073	19.62288236082625	-0.07460112958068	-0.92165131536375
3	0.23476467495811	8.57909346342925	-0.11021349542738	-0.81143781993637
4	0.39254785728080	3.63763326452917	-0.15778318232269	-0.65365463761368
5	0.60067545233191	1.42717892023773	-0.20812759505112	-0.44552704256257
6	0.82714994286833	0.46286007749125	-0.22647449053641	-0.21905255202615
7	0.99028203077844	0.09369191826377	-0.16313208791011	-0.05592046411604
8	1.04242438221432	0.00592723560279	-0.05214235143588	-0.00377811268016
9	1.04618505691071	0.00002723158211	-0.00376067469639	-0.00001743798377
10	1.04620249452271	0.00000000058056	-0.00001743761199	-0.00000000037178

Note that Halley's iteration is superior in this case, since F is a rational function.

! Newton method converges more slowly, but also needs less effort per step (→ Sect. 3.3.3) \diamond

In the previous example Newton's method performed rather poorly. Often its convergence can be boosted by converting the non-linear equation to an equivalent one (that is, one with the same solutions) for another function g , which is "closer to a linear function":

3.3
p. 178

3.3
p. 18

Assume $F \approx \widehat{F}$, where \widehat{F} is invertible with an inverse \widehat{F}^{-1} that can be evaluated with little effort.

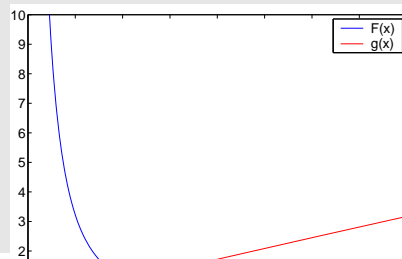
$$\blacktriangleright g(x) := \widehat{F}^{-1}(F(x)) \approx x.$$

Then apply Newton's method to $g(x)$, using the formula for the derivative of the inverse of a function

$$\frac{d}{dy}(\widehat{F}^{-1})(y) = \frac{1}{\widehat{F}'(\widehat{F}^{-1}(y))} \Rightarrow g'(x) = \frac{1}{\widehat{F}'(g(x))} \cdot F'(x).$$

Example 3.3.5 (Adapted Newton method).

As in Ex. 3.3.4: $F(x) = \frac{1}{(x+1)^2} + \frac{1}{(x+0.1)^2} - 1, \quad x > 0:$



Observation:

$$F(x) + 1 \approx 2x^{-2} \text{ for } x \gg 1$$

and so $g(x) := \frac{1}{\sqrt{F(x)+1}}$ "almost" linear for $x \gg 1$

Idea: instead of $F(x) \stackrel{!}{=} 0$ tackle $g(x) \stackrel{!}{=} 1$ with Newton's method (3.3.1).

$$\begin{aligned} x^{(k+1)} &= x^{(k)} - \frac{g(x^{(k)}) - 1}{g'(x^{(k)})} = x^{(k)} + \left(\frac{1}{\sqrt{F(x^{(k)})+1}} - 1 \right) \frac{2(F(x^{(k)})+1)^{3/2}}{F'(x^{(k)})} \\ &= x^{(k)} + \frac{2(F(x^{(k)})+1)(1 - \sqrt{F(x^{(k)})+1})}{F'(x^{(k)})}. \end{aligned}$$

Convergence recorded for $x^{(0)} = 0$:

k	$x^{(k)}$	$F(x^{(k)})$	$x^{(k)} - x^{(k-1)}$	$x^{(k)} - x^*$
1	0.91312431341979	0.24747993091128	0.91312431341979	-0.13307818147469
2	1.04517022155323	0.00161402574513	0.13204590813344	-0.00103227334125
3	1.04620244004116	0.00000008565847	0.00103221848793	-0.00000005485332
4	1.04620249489448	0.00000000000000	0.00000005485332	-0.00000000000000

◇

For zero finding there is wealth of iterative methods that offer higher order of convergence.

One idea: **consistent modification** of the Newton-Iteration:

$$\blacktriangleright \text{fixed point iteration: } \Phi(x) = x - \frac{F(x)}{F'(x)}H(x) \text{ with "proper" } H: I \mapsto \mathbb{R}.$$

Aim: find H such that the method is of p -th order; tool: Lemma 3.2.7.

Assume: F smooth "enough" and $\exists x^* \in I: F(x^*) = 0, F'(x^*) \neq 0$.

$$\begin{aligned} \Phi &= x - uH, \quad \Phi' = 1 - u'H - uH', \quad \Phi'' = -u''H - 2u'H - uH'', \\ \text{with } u &= \frac{F}{F'} \Rightarrow u' = 1 - \frac{FF''}{(F')^2}, \quad u'' = -\frac{F''}{F'} + 2\frac{F(F'')^2}{(F')^3} - \frac{FF'''}{(F')^2}. \end{aligned}$$

$$F(x^*) = 0 \blacktriangleright u(x^*) = 0, u'(x^*) = 1, u''(x^*) = -\frac{F''(x^*)}{F'(x^*)}.$$

$$\blacktriangleright \Phi'(x^*) = 1 - H(x^*), \quad \Phi''(x^*) = \frac{F''(x^*)}{F'(x^*)}H(x^*) - 2H'(x^*). \quad (3.3.2)$$

Lemma 3.2.7 \blacktriangleright **Necessary** conditions for local convergence of order p :

$$p = 2 \text{ (quadratic convergence): } H(x^*) = 1,$$

$$p = 3 \text{ (cubic convergence): } H(x^*) = 1 \wedge H'(x^*) = \frac{1}{2} \frac{F''(x^*)}{F'(x^*)}.$$

3.3
p. 181

3.3
p. 18

In particular: $H(x) = G(1 - u'(x))$ with "proper" G

$$\blacktriangleright \text{fixed point iteration } x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})}G\left(\frac{F(x^{(k)})F''(x^{(k)})}{(F'(x^{(k)}))^2}\right). \quad (3.3.3)$$

Lemma 3.3.1. If $F \in C^2(I)$, $F(x^*) = 0$, $F'(x^*) \neq 0$, $G \in C^2(U)$ in a neighbourhood U of 0, $G(0) = 1$, $G'(0) = \frac{1}{2}$, then the fixed point iteration (3.3.3) converge locally cubically to x^* .

Proof: Lemma 3.2.7, (3.3.2) and

$$H(x^*) = G(0), \quad H'(x^*) = -G'(0)u''(x^*) = G'(0)\frac{F''(x^*)}{F'(x^*)}.$$

Example 3.3.6 (Example of modified Newton method).

$$\bullet G(t) = \frac{1}{1 - \frac{1}{2}t} \Rightarrow \text{Halley's iteration } (\rightarrow \text{Ex. 3.3.4})$$

$$\bullet G(t) = \frac{2}{1 + \sqrt{1 - 2t}} \Rightarrow \text{Euler's iteration}$$

3.3
p. 182

3.3
p. 18

- $G(t) = 1 + \frac{1}{2}t \rightarrow$ quadratic inverse interpolation

k	$e^{(k)} := x^{(k)} - x^*$		
	Halley	Euler	Quad. Inv.
1	2.81548211105635	3.57571385244736	2.03843730027891
2	1.37597082614957	2.76924150041340	1.02137913293045
3	0.34002908011728	1.95675490333756	0.28835890388161
4	0.00951600547085	1.25252187565405	0.01497518178983
5	0.00000024995484	0.51609312477451	0.00000315361454
6		0.14709716035310	
7		0.00109463314926	
8		0.00000000107549	

Numerical experiment:

$$F(x) = xe^x - 1, \quad x^{(0)} = 5$$

3.3.2.3 Multi-point methods



Idea:

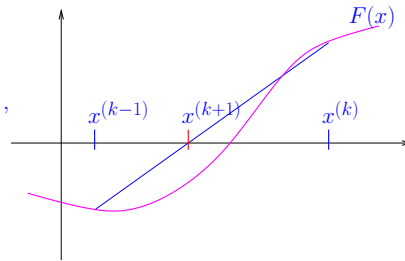
Replace F with **interpolating polynomial** producing interpolatory model function methods

Simplest representative: **secant method**

$x^{(k+1)}$ = zero of secant

$$s(x) = F(x^{(k)}) + \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}(x - x^{(k)}), \quad (3.3.4)$$

$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})(x^{(k)} - x^{(k-1)})}{F(x^{(k)}) - F(x^{(k-1)})}. \quad (3.3.5)$$



Code 3.3.7: secant method

```

1 function x = secant(x0,x1,F,tol)
2 fo = F(x0);
3 for i=1:MAXIT
4   fn = F(x1);
5   s = fn*(x1-x0)/(fn-fo);
6   x0 = x1; x1 = x1-s;
7   if (abs(s) < tol), x = x1; return;
8   end
9   fo = fn;
end

```

secant method

(MATLAB implementation)

- Only one function evaluation per step
- **no derivatives required!**

Remember: $F(x)$ may only be available as output of a (complicated) procedure. In this case it is difficult to find a procedure that evaluates $F'(x)$. Thus the significance of methods that do not involve evaluations of derivatives.

Example 3.3.8 (secant method). $F(x) = xe^x - 1, \quad x^{(0)} = 0, \quad x^{(1)} = 5.$

k	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log e^{(k+1)} - \log e^{(k)} }{\log e^{(k)} - \log e^{(k-1)} }$
2	0.00673794699909	-0.99321649977589	-0.56040534341070	
3	0.01342122983571	-0.98639742654892	-0.55372206057408	24.43308649757745
4	0.98017620833821	1.61209684919288	0.41303291792843	2.70802321457994
5	0.38040476787948	-0.44351476841567	-0.18673852253030	1.48753625853887
6	0.50981028847430	-0.15117846201565	-0.05733300193548	1.51452723840131
7	0.57673091089295	0.02670169957932	0.00958762048317	1.70075240166256
8	0.56668541543431	-0.00126473620459	-0.00045787497547	1.59458505614449
9	0.56713970649585	-0.00000990312376	-0.00000358391394	1.62641838319117
10	0.56714329175406	0.00000000371452	0.0000000134427	
11	0.56714329040978	-0.00000000000001	-0.00000000000000	

A startling observation: the method seems to have a **fractional** (!) order of convergence, see Def. 3.1.7.

3.3

p. 185

3.3

p. 18

Remark 3.3.9 (Fractional order of convergence of secant method).

Indeed, a fractional order of convergence can be proved for the secant method, see [28, Sect. 18.2]. Here is a crude outline of the reasoning:

Asymptotic convergence of secant method: error $e^{(k)} := x^{(k)} - x^*$

$$e^{(k+1)} = \Phi(x^* + e^{(k)}, x^* + e^{(k-1)}) - x^*, \quad \text{with } \Phi(x, y) := x - \frac{F(x)(x-y)}{F(x) - F(y)}. \quad (3.3.6)$$

Use MAPLE to find Taylor expansion (assuming F sufficiently smooth):

```

> Phi := (x,y) -> x-F(x)*(x-y)/(F(x)-F(y));
> F(s) := 0;
> e2 = normal(mtaylor(Phi(s+e1,s+e0)-s,[e0,e1],4));

```

> linearized error propagation formula:

$$e^{(k+1)} \doteq \frac{1}{2} \frac{F''(x^*)}{F'(x^*)} e^{(k)} e^{(k-1)} = C e^{(k)} e^{(k-1)}. \quad (3.3.7)$$

3.3

p. 186

3.3

p. 18

Try $e^{(k)} = K(e^{(k-1)})^p$ to determine the order of convergence (\rightarrow Def. 3.1.7):

$$\Rightarrow e^{(k+1)} = K^{p+1}(e^{(k-1)})^{p^2}$$

$$\Rightarrow (e^{(k-1)})^{p^2-p-1} = K^{-p}C \Rightarrow p^2 - p - 1 = 0 \Rightarrow p = \frac{1}{2}(1 \pm \sqrt{5}).$$

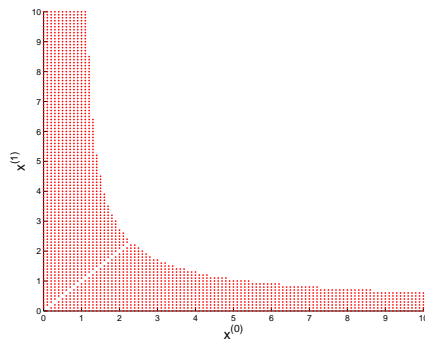
As $e^{(k)} \rightarrow 0$ for $k \rightarrow \infty$ we get the rate of convergence $p = \frac{1}{2}(1 + \sqrt{5}) \approx 1.62$ (see Ex. 3.3.8 !)

Example 3.3.10 (local convergence of secant method).

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

$\hat{=}$ secant method converges for a pair $(x^{(0)}, x^{(1)})$ of initial guesses. \triangleright

= local convergence \rightarrow Def. 3.1.3 \diamond



Another class of multi-point methods: *inverse interpolation*

Assume:

$$F : I \subset \mathbb{R} \mapsto \mathbb{R} \text{ one-to-one}$$

$$F(x^*) = 0 \Rightarrow F^{-1}(0) = x^*.$$

• Interpolate F^{-1} by polynomial p of degree d determined by

$$p(F(x^{(k-m)})) = x^{(k-m)}, \quad m = 0, \dots, d.$$

• New approximate zero $x^{(k+1)} := p(0)$



$$F(x^*) = 0 \Leftrightarrow F^{-1}(0) = x^*$$

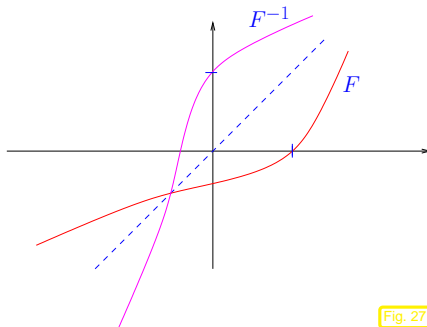


Fig. 27

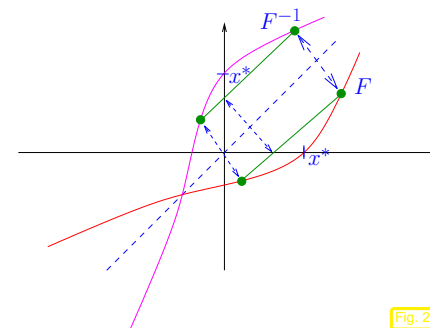


Fig. 28

Case $m = 1$ \triangleright secant method

Case $m = 2$: quadratic inverse interpolation, see [33, Sect. 4.5]

MAPLE code: `p := x -> a*x^2+b*x+c;`
`solve({p(f0)=x0,p(f1)=x1,p(f2)=x2},{a,b,c});`
`assign(%); p(0);`

$$x^{(k+1)} = \frac{F_0^2(F_1x_2 - F_2x_1) + F_1^2(F_2x_0 - F_0x_2) + F_2^2(F_0x_1 - F_1x_0)}{F_0^2(F_1 - F_2) + F_1^2(F_2 - F_0) + F_2^2(F_0 - F_1)}.$$

($F_0 := F(x^{(k-2)})$, $F_1 := F(x^{(k-1)})$, $F_2 := F(x^{(k)})$, $x_0 := x^{(k-2)}$, $x_1 := x^{(k-1)}$, $x_2 := x^{(k)}$)

Example 3.3.11 (quadratic inverse interpolation). $F(x) = xe^x - 1$, $x^{(0)} = 0$, $x^{(1)} = 2.5$, $x^{(2)} = 5$.

k	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log e^{(k+1)} - \log e^{(k)} }{\log e^{(k)} - \log e^{(k-1)} }$
3	0.08520390058175	-0.90721814294134	-0.48193938982803	
4	0.16009252622586	-0.81211229637354	-0.40705076418392	3.33791154378839
5	0.79879381816390	0.77560534067946	0.23165052775411	2.28740488912208
6	0.63094636752843	0.18579323999999	0.06380307711864	1.82494667289715
7	0.56107750991028	-0.01667806436181	-0.00606578049951	1.87323264214217
8	0.56706941033107	-0.00020413476766	-0.00007388007872	1.79832936980454
9	0.56714331707092	0.00000007367067	0.00000002666114	1.84841261527097
10	0.56714329040980	0.000000000000003	0.000000000000001	

Also in this case the numerical experiment hints at a fractional rate of convergence, as in the case of the secant method, see Rem. 3.3.9. \diamond

3.3.3 Note on Efficiency

Efficiency of an iterative method (for solving $F(\mathbf{x}) = 0$) \leftrightarrow computational effort to reach prescribed number of significant digits in result.

Abstract:

$W \hat{=}$ computational effort per step

(e.g. $W \approx \frac{\#\{\text{evaluations of } F\}}{\text{step}} + n \cdot \frac{\#\{\text{evaluations of } F'\}}{\text{step}} + \dots$)

Crucial: number of steps $k = k(\rho)$ to achieve relative reduction of error

$$\|e^{(k)}\| \leq \rho \|e^{(0)}\|, \quad \rho > 0 \text{ prescribed?} \quad (3.3.8)$$

Error recursion can be converted into expressions (3.3.9) and (3.3.10) that related the error norm $\|e^{(k)}\|$ to $\|e^{(0)}\|$ and lead to quantitative bounds for the number of steps to achieve (3.3.8) for iterative method of order $p, p \geq 1$ (\rightarrow Def. 3.1.7):

$$\exists C > 0: \|e^{(k)}\| \leq C \|e^{(k-1)}\|^p \quad \forall k \geq 1 \quad (C < 1 \text{ for } p = 1).$$

Assuming $C \|e^{(0)}\|^{p-1} < 1$ (guarantees convergence):

$$p = 1: \|e^{(k)}\| \leq C^k \|e^{(0)}\| \quad \text{requires } k \geq \frac{\log \rho}{\log C}, \quad (3.3.9)$$

$$p > 1: \|e^{(k)}\| \leq C^{\frac{k-1}{p-1}} \|e^{(0)}\|^{p^k} \quad \text{requires } p^k \geq 1 + \frac{\log \rho}{\log C/p-1 + \log(\|e^{(0)}\|)}$$

$$\Rightarrow k \geq \log\left(1 + \frac{\log \rho}{\log L_0}\right) / \log p, \quad (3.3.10)$$

$$L_0 := C^{1/p-1} \|e^{(0)}\| < 1.$$

If $\rho \ll 1$, then $\log\left(1 + \frac{\log \rho}{\log L_0}\right) \approx \log |\log \rho| - \log |\log L_0| \approx \log |\log \rho|$. This simplification will be made in the context of asymptotic considerations $\rho \rightarrow 0$ below.

Notice: $|\log \rho| \leftrightarrow$ No. of significant digits of $x^{(k)}$

Measure for efficiency:

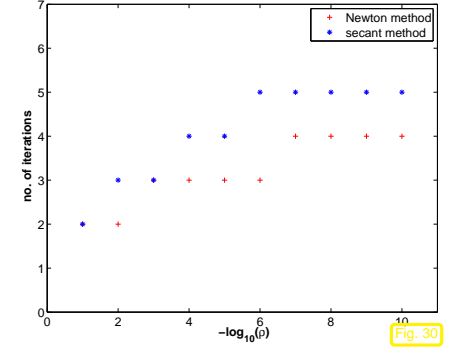
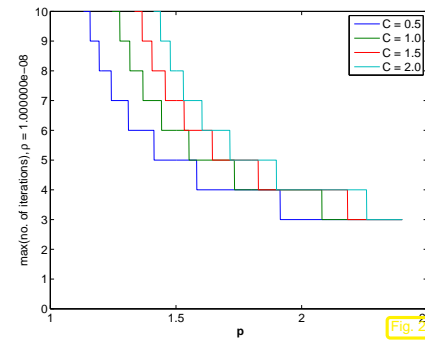
$$\text{Efficiency} := \frac{\text{no. of digits gained}}{\text{total work required}} = \frac{|\log \rho|}{k(\rho) \cdot W}$$

(3.3.11)

asymptotic efficiency w.r.t. $\rho \rightarrow 0 \rightarrow |\log \rho| \rightarrow \infty$:

$$\text{Efficiency}_{|\rho \rightarrow 0} = \begin{cases} -\frac{\log C}{W} & , \text{ if } p = 1, \\ \frac{\log p |\log \rho|}{W \log |\log \rho|} & , \text{ if } p > 1. \end{cases} \quad (3.3.12)$$

Example 3.3.12 (Efficiency of iterative methods).



3.3
p. 193

Evaluation (3.3.10) for $\|e^{(0)}\| = 0.1, \rho = 10^{-8}$

Newton's method \leftrightarrow secant method, $C = 1$,
initial error $\|e^{(0)}\| = 0.1$

3.3
p. 19

$$\frac{W_{\text{Newton}}}{p_{\text{Newton}}} = 2 \frac{W_{\text{secant}}}{p_{\text{secant}}} \Rightarrow \frac{\log p_{\text{Newton}}}{W_{\text{Newton}}} : \frac{\log p_{\text{secant}}}{W_{\text{secant}}} = 0.71.$$

secant method is more efficient than Newton's method!

\Rightarrow

\diamond

3.4 Newton's Method

Non-linear system of equations: for $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ find $\mathbf{x}^* \in D: F(\mathbf{x}^*) = 0$

Assume: $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ continuously differentiable

3.3
p. 194

3.4
p. 19

3.4.1 The Newton iteration



Idea (→ Sect. 3.3.2.1):

local linearization:

Given $\mathbf{x}^{(k)} \in D \ni \mathbf{x}^{(k+1)}$ as zero of affine linear model function

$$F(\mathbf{x}) \approx \tilde{F}(\mathbf{x}) := F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}),$$

$$DF(\mathbf{x}) \in \mathbb{R}^{n,n} = \text{Jacobian (ger.: Jacobi-Matrix)}, DF(\mathbf{x}) = \left(\frac{\partial F_j}{\partial x_k}(\mathbf{x}) \right)_{j,k=1}^n.$$

► **Newton iteration:** (\Leftrightarrow (3.3.1) for $n = 1$)

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - DF(\mathbf{x}^{(k)})^{-1}F(\mathbf{x}^{(k)}), \quad [\text{if } DF(\mathbf{x}^{(k)}) \text{ regular}] \quad (3.4.1)$$

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

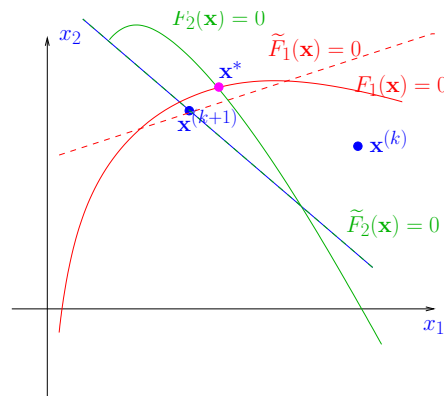
Illustration of idea of Newton's method for $n = 2$:

▷

Sought: intersection point \mathbf{x}^* of the curves

$$F_1(\mathbf{x}) = 0 \text{ and } F_2(\mathbf{x}) = 0.$$

Idea: $\mathbf{x}^{(k+1)}$ = the intersection of two straight lines (= zero sets of the components of the model function, cf. Ex. 2.4.11) that are approximations of the original curves



```
MATLAB-CODE: Newtonverfahren
function x = newton(x,F,DF,tol)
for i=1:MAXIT
    s = DF(x)\F(x);
    x = x-s;
    if (norm(s) < tol*norm(x))
        return; end;
end
```

MATLAB template for Newton method:

Solve linear system:

$$A \setminus b = A^{-1}b \rightarrow \text{Chapter 2}$$

F, DF : function handles

A posteriori termination criterion

Example 3.4.1 (Newton method in 2D).

$$F(\mathbf{x}) = \begin{pmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^2 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2 \quad \text{with solution } F \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0.$$

$$\text{Jacobian (analytic computation): } DF(\mathbf{x}) = \begin{pmatrix} \partial_{x_1} F_1(x) & \partial_{x_2} F_1(x) \\ \partial_{x_1} F_2(x) & \partial_{x_2} F_2(x) \end{pmatrix} = \begin{pmatrix} 2x_1 & -4x_2^3 \\ 1 & -3x_2^2 \end{pmatrix}$$

Realization of Newton iteration (3.4.1)

1. Solve LSE

$$\begin{pmatrix} 2x_1 & -4x_2^3 \\ 1 & -3x_2^2 \end{pmatrix} \Delta \mathbf{x}^{(k)} = - \begin{pmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^2 \end{pmatrix}$$

where $\mathbf{x}^{(k)} = (x_1, x_2)^T$.

2. Set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)}$

Code 3.4.2: Newton iteration in 2D

```
1 F = @(x) [ x(1)^2-x(2)^4; x(1)-x(2)^3 ];
2 DF = @(x) [ 2*x(1), -4*x(2)^3; 1, -3*x(2)^2 ];
3 x = [0.7; 0.7]; x_ast = [1;1]; tol = 1E-10;
4
5 res = [0,x',norm(x-x_ast)];
6 s = feval(DF,x)\feval(F,x); x = x-s;
7 res = [res; 1,x',norm(x-x_ast)]; k=2;
8 while (norm(s) > tol*norm(x))
9     s = DF(x)\F(x); x = x-s;
10    res = [res; k,x',norm(x-x_ast)];
11    k = k+1;
12 end
13
14 logdiff = diff(log(res(:,4)));
15 rates = logdiff(2:end)./logdiff(1:end-1);
```

k	$\mathbf{x}^{(k)}$	$\epsilon_k := \ \mathbf{x}^* - \mathbf{x}^{(k)}\ _2$
0	$(0.7, 0.7)^T$	4.24e-01
1	$(0.87850000000000, 1.064285714285714)^T$	1.37e-01
2	$(1.01815943274188, 1.00914882463936)^T$	2.03e-02
3	$(1.00023355916300, 1.00015913936075)^T$	2.83e-04
4	$(1.00000000583852, 1.00000002726552)^T$	2.79e-08
5	$(0.999999999999998, 1.000000000000000)^T$	2.11e-15
6	$(1, 1)^T$	

◇



New aspect for $n \gg 1$ (compared to $n = 1$ -dimensional case, section. 3.3.2.1):

Computation of the Newton correction is eventually costly!

Remark 3.4.3 (Affine invariance of Newton method).

An important property of the Newton iteration (3.4.1): **affine invariance** \rightarrow [12, Sect .1.2.2]

set $G(\mathbf{x}) := \mathbf{A}F(\mathbf{x})$ with regular $\mathbf{A} \in \mathbb{R}^{n,n}$ so that $F(\mathbf{x}^*) = 0 \Leftrightarrow G(\mathbf{x}^*) = 0$.

affine invariance: Newton iteration for $G(\mathbf{x}) = 0$ is the same for all $\mathbf{A} \in GL(n)$!

$$DG(\mathbf{x}) = \mathbf{A}DF(\mathbf{x}) \Rightarrow DG(\mathbf{x})^{-1}G(\mathbf{x}) = DF(\mathbf{x})^{-1}\mathbf{A}^{-1}\mathbf{A}F(\mathbf{x}) = DF(\mathbf{x})^{-1}F(\mathbf{x}).$$

Use affine invariance as guideline for

- convergence theory for Newton's method: assumptions and results should be affine invariant, too.
- modifying and extending Newton's method: resulting schemes should preserve affine invariance.

Remark 3.4.4 (Differentiation rules). \rightarrow Repetition: basic analysis

Statement of the Newton iteration (3.4.1) for $F : \mathbb{R}^n \mapsto \mathbb{R}^n$ given as analytic expression entails computing the Jacobian DF . To avoid cumbersome component-oriented considerations, it is useful to know the *rules of multidimensional differentiation*:

Let V, W be finite dimensional vector spaces, $F : D \subset V \mapsto W$ sufficiently smooth. The **differential** $DF(\mathbf{x})$ of F in $\mathbf{x} \in V$ is the *unique*

linear mapping $DF(\mathbf{x}) : V \mapsto W$,
such that $\|F(\mathbf{x} + \mathbf{h}) - F(\mathbf{x}) - DF(\mathbf{x})\mathbf{h}\| = o(\|\mathbf{h}\|) \quad \forall \mathbf{h}, \|\mathbf{h}\| < \delta$.

• For $F : V \mapsto W$ linear, i.e. $F(\mathbf{x}) = \mathbf{A}\mathbf{x}$, \mathbf{A} matrix $\rightarrow DF(\mathbf{x}) = \mathbf{A}$.

• **Chain rule:** $F : V \mapsto W, G : W \mapsto U$ sufficiently smooth

$$D(G \circ F)(\mathbf{x})\mathbf{h} = DG(F(\mathbf{x}))(DF(\mathbf{x})\mathbf{h}), \quad \mathbf{h} \in V, \mathbf{x} \in D. \quad (3.4.2)$$

• **Product rule:** $F : D \subset V \mapsto W, G : D \subset V \mapsto U$ sufficiently smooth, $b : W \times U \mapsto Z$ **bilinear**

$$T(\mathbf{x}) = b(F(\mathbf{x}), G(\mathbf{x})) \Rightarrow DT(\mathbf{x})\mathbf{h} = b(DF(\mathbf{x})\mathbf{h}, G(\mathbf{x})) + b(F(\mathbf{x}), DG(\mathbf{x})\mathbf{h}), \quad (3.4.3)$$

$$\mathbf{h} \in V, \mathbf{x} \in D.$$

For $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}$ the **gradient** $\text{grad } F : D \mapsto \mathbb{R}^n$, and the **Hessian matrix** $HF(\mathbf{x}) : D \mapsto \mathbb{R}^{n,n}$ are defined as

$$\text{grad } F(\mathbf{x})^T \mathbf{h} := DF(\mathbf{x})\mathbf{h}, \quad \mathbf{h}_1^T HF(\mathbf{x})\mathbf{h}_2 := D(DF(\mathbf{x})(\mathbf{h}_1))(\mathbf{h}_2), \quad \mathbf{h}, \mathbf{h}_1, \mathbf{h}_2 \in V.$$

3.4
p. 201

△

3.4
p. 20

Remark 3.4.5 (Simplified Newton method).

Simplified Newton Method: use the same $DF(\mathbf{x}^{(k)})$ for more steps
 \rightarrow (usually) merely linear convergence instead of quadratic convergence

△

Remark 3.4.6 (Numerical Differentiation for computation of Jacobian).

If $DF(\mathbf{x})$ is not available (e.g. when $F(\mathbf{x})$ is given only as a procedure):

$$\text{Numerical Differentiation: } \frac{\partial F_i}{\partial x_j}(\mathbf{x}) \approx \frac{F_i(\mathbf{x} + h\mathbf{e}_j) - F_i(\mathbf{x})}{h}.$$

Caution: impact of roundoff errors for small h !

△

Example 3.4.7 (Roundoff errors and difference quotients).

$$\text{Approximate derivative by difference quotient: } f'(x) \approx \frac{f(x+h) - f(x)}{h}.$$

3.4
p. 202

Calculus: better approximation for smaller $h > 0$, isn't it ?

3.4
p. 20

$$F(\mathbf{x}^*) = 0 \Rightarrow D\Phi(\mathbf{x}^*) = 0.$$

3.2.7

Local quadratic convergence of Newton's method, if $DF(\mathbf{x}^*)$ regular

Example 3.4.9 (Convergence of Newton's method).

Ex. 3.4.1 cont'd: record of iteration errors, see Code 3.4.1:

k	$\mathbf{x}^{(k)}$	$\epsilon_k := \ \mathbf{x}^* - \mathbf{x}^{(k)}\ _2$	$\frac{\log \epsilon_{k+1} - \log \epsilon_k}{\log \epsilon_k - \log \epsilon_{k-1}}$
0	$(0.7, 0.7)^T$	4.24e-01	
1	$(0.87850000000000, 1.064285714285714)^T$	1.37e-01	1.69
2	$(1.01815943274188, 1.00914882463936)^T$	2.03e-02	2.23
3	$(1.00023355916300, 1.00015913936075)^T$	2.83e-04	2.15
4	$(1.00000000583852, 1.00000002726552)^T$	2.79e-08	1.77
5	$(0.999999999999998, 1.000000000000000)^T$	2.11e-15	
6	$(1, 1)^T$		

◇

There is a sophisticated theory about the convergence of Newton's method. For example one can find the following theorem in [14, Thm. 4.10], [12, Sect. 2.1]:

Theorem 3.4.1 (Local quadratic convergence of Newton's method). If:

- (A) $D \subset \mathbb{R}^n$ open and convex,
- (B) $F : D \mapsto \mathbb{R}^n$ continuously differentiable,
- (C) $DF(\mathbf{x})$ regular $\forall \mathbf{x} \in D$,
- (D) $\exists L \geq 0: \left\| DF(\mathbf{x})^{-1}(DF(\mathbf{x} + \mathbf{v}) - DF(\mathbf{x})) \right\|_2 \leq L \|\mathbf{v}\|_2 \quad \forall \mathbf{v} \in \mathbb{R}^n, \mathbf{v} + \mathbf{x} \in D, \mathbf{x} \in D$,
- (E) $\exists \mathbf{x}^*: F(\mathbf{x}^*) = 0$ (existence of solution in D)
- (F) initial guess $\mathbf{x}^{(0)} \in D$ satisfies $\rho := \left\| \mathbf{x}^* - \mathbf{x}^{(0)} \right\|_2 < \frac{2}{L} \wedge B_\rho(\mathbf{x}^*) \subset D$.

then the Newton iteration (3.4.1) satisfies:

- (i) $\mathbf{x}^{(k)} \in B_\rho(\mathbf{x}^*) := \{\mathbf{y} \in \mathbb{R}^n, \|\mathbf{y} - \mathbf{x}^*\| < \rho\}$ for all $k \in \mathbb{N}$,
- (ii) $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$,
- (iii) $\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\|_2 \leq \frac{L}{2} \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|_2^2$ (local quadratic convergence).

notation: ball $B_\rho(\mathbf{z}) := \{\mathbf{x} \in \mathbb{R}^n: \|\mathbf{x} - \mathbf{z}\|_2 \leq \rho\}$

Terminology: (D) $\hat{=}$ affine invariant Lipschitz condition

Problem: Usually neither ω nor x^* are known !

► In general: a priori estimates as in Thm. 3.4.1 are of little practical relevance.

3.4.3 Termination of Newton iteration

A first viable idea:

Asymptotically due to quadratic convergence:

$$\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \ll \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| \Rightarrow \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| \approx \left\| \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \right\|. \quad (3.4.4)$$

3.4
p. 209

3.4
p. 21

► quit iterating as soon as $\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \right\| = \left\| DF(\mathbf{x}^{(k)})^{-1}F(\mathbf{x}^{(k)}) \right\| < \tau \left\| \mathbf{x}^{(k)} \right\|$, with $\tau =$ tolerance

→ uneconomical: one needless update, because $\mathbf{x}^{(k)}$ already accurate enough !

Remark 3.4.10. New aspect for $n \gg 1$: computation of Newton correction may be expensive !

Therefore we would like to use an a-posteriori termination criterion that dispenses with computing (and "inverting") another Jacobian $DF(\mathbf{x}^{(k)})$ just to tell us that $\mathbf{x}^{(k)}$ is already accurate enough.

Practical a-posteriori termination criterion for Newton's method:

$$DF(\mathbf{x}^{(k-1)}) \approx DF(\mathbf{x}^{(k)}): \text{ quit as soon as } \left\| DF(\mathbf{x}^{(k-1)})^{-1}F(\mathbf{x}^{(k)}) \right\| < \tau \left\| \mathbf{x}^{(k)} \right\|$$

affine invariant termination criterion

Justification: we expect $DF(\mathbf{x}^{(k-1)}) \approx DF(\mathbf{x}^{(k)})$, when Newton iteration has converged. Then appeal to (3.4.4).

3.4
p. 210

3.4
p. 21

If we used the residual based termination criterion

$$\|F(\mathbf{x}^{(k)})\| \leq \tau,$$

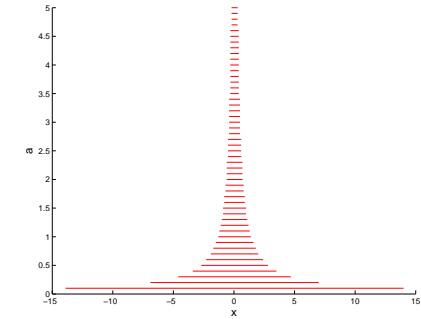
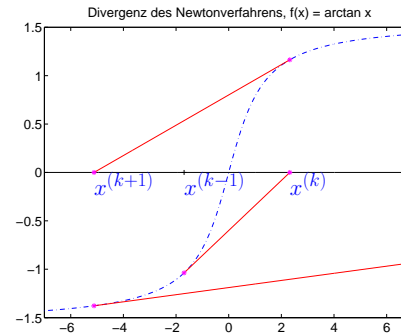
then the resulting algorithm would not be affine invariant, because for $F(\mathbf{x}) = 0$ and $\mathbf{A}F(\mathbf{x}) = 0$, $\mathbf{A} \in \mathbb{R}^{n,n}$ regular, the Newton iteration might terminate with different iterates.

Terminology: $\Delta \bar{\mathbf{x}}^{(k)} := DF(\mathbf{x}^{(k-1)})^{-1}F(\mathbf{x}^{(k)}) \hat{=} \text{simplified Newton correction}$

Reuse of LU-factorization (\rightarrow Rem. 2.2.6) of $DF(\mathbf{x}^{(k-1)}) \blacktriangleright \Delta \bar{\mathbf{x}}^{(k)}$ available with $O(n^2)$ operations

Summary: The Newton Method

- 😊 converges *asymptotically* very fast: doubling of number of significant digits in each step
- 😞 often a very small region of convergence, which requires an initial guess rather close to the solution.



red zone = $\{x^{(0)} \in \mathbb{R}, x^{(k)} \rightarrow 0\}$



we observe "overshooting" of Newton correction

Idea: **damping** of Newton correction:

$$\text{With } \lambda^{(k)} > 0: \mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \lambda^{(k)} DF(\mathbf{x}^{(k)})^{-1}F(\mathbf{x}^{(k)}). \quad (3.4.5)$$

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

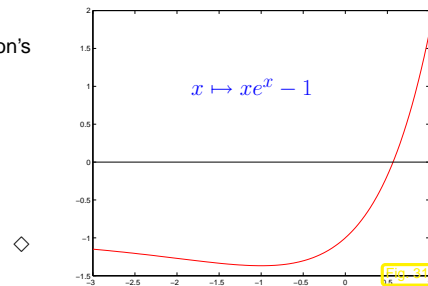
3.4.4 Damped Newton method

Example 3.4.11 (Local convergence of Newton's method).

$$F(x) = xe^x - 1 \Rightarrow F'(-1) = 0$$

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

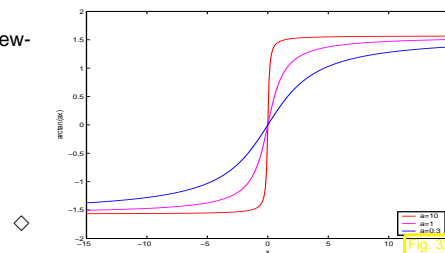
$$x^{(0)} > -1 \Rightarrow x^{(k)} \rightarrow x^*$$



Example 3.4.12 (Region of convergence of Newton method).

$$F(x) = \arctan(ax), \quad a > 0, x \in \mathbb{R}$$

with zero $x^* = 0$.



Choice of damping factor: affine invariant **natural monotonicity test** [12, Ch. 3]

$$\text{"maximal" } \lambda^{(k)} > 0: \|\Delta \bar{\mathbf{x}}(\lambda^{(k)})\| \leq \left(1 - \frac{\lambda^{(k)}}{2}\right) \|\Delta \mathbf{x}^{(k)}\|_2 \quad (3.4.6)$$

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

Heuristics: convergence \Leftrightarrow size of Newton correction decreases

Code 3.4.14: Damped Newton method (non-optimal implementation !)

```

1 function [x,cvg] = dampnewton(x,F,DF,tol)
2 [L,U] = lu(DF(x)); s = U\(\L\F(x));
3 xn = x-s; lambda = 1; cvg = 0;
4 f = F(xn); st = U\(\L\f);
5 while (norm(st) > tol*norm(xn))
6   while (norm(st) > (1-lambda/2)*norm(s))
7     lambda = lambda/2;
8     if (lambda < LMIN), cvg = -1; return; end
9     xn = x-lambda*s; f = F(xn); st = U\(\L\f);
10  end
11  x = xn; [L,U] = lu(DF(x)); s = U\(\L\f);
12  lambda = min(2*lambda,1);
13  xn = x-lambda*s; f = F(xn); st = U\(\L\f);
14 end
15 x = xn;
  
```

Reuse of LU-factorization, see Rem. 2.2.6

a-posteriori termination criterion (based on simplified Newton correction, cf. Sect. 3.4.3)

Natural monotonicity test (3.4.6)

Reduce damping factor λ

Policy: Reduce damping factor by a factor $q \in]0, 1[$ (usually $q = \frac{1}{2}$) until the affine invariant natural monotonicity test (3.4.6) passed.

Example 3.4.15 (Damped Newton method). (\rightarrow Ex. 3.4.12)

$F(x) = \arctan(x)$,

- $x^{(0)} = 20$
- $q = \frac{1}{2}$
- LMIN = 0.001

Observation: asymptotic quadratic convergence

k	$\lambda^{(k)}$	$x^{(k)}$	$F(x^{(k)})$
1	0.03125	0.94199967624205	0.75554074974604
2	0.06250	0.85287592931991	0.70616132170387
3	0.12500	0.70039827977515	0.61099321623952
4	0.25000	0.47271811131169	0.44158487422833
5	0.50000	0.20258686348037	0.19988168667351
6	1.00000	-0.00549825489514	-0.00549819949059
7	1.00000	0.00000011081045	0.00000011081045
8	1.00000	-0.000000000000001	-0.000000000000001

Example 3.4.16 (Failure of damped Newton method).

- As in Ex. 3.4.11:

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

- Initial guess for damped Newton method $x^{(0)} = -1.5$

k	$\lambda^{(k)}$	$x^{(k)}$	$F(x^{(k)})$
1	0.25000	-4.4908445351690	-1.0503476286303
2	0.06250	-6.1682249558799	-1.0129221310944
3	0.01562	-7.6300006580712	-1.0037055902301
4	0.00390	-8.8476436930246	-1.0012715832278
5	0.00195	-10.5815494437311	-1.0002685596314

Bailed out because of $\lambda < \text{LMIN}$!

Observation: Newton correction pointing in "wrong direction" so no convergence.

3.4.5 Quasi-Newton Method

What to do when $DF(x)$ is not available and numerical differentiation (see remark 3.4.6) is too expensive?



Idea: in one dimension ($n = 1$) apply the secant method (3.3.4) of section 3.3.2.3

$$F'(x^{(k)}) \approx \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}} \quad \text{"difference quotient"} \quad (3.4.7)$$



Generalisation for $n > 1$?

Idea: rewrite (3.4.7) as a **secant condition** for the approximation $J_k \approx DF(x^{(k)})$, $x^{(k)} \hat{=}$ iterate:



$$J_k(x^{(k)} - x^{(k-1)}) = F(x^{(k)}) - F(x^{(k-1)}) \quad (3.4.8)$$

BUT: many matrices J_k fulfill (3.4.8)

Hence: we need more conditions for $J_k \in \mathbb{R}^{n,n}$

3.4 p. 217

Idea: get

J_k by a **modification** of J_{k-1}

Broyden conditions: $J_k z = J_{k-1} z \quad \forall z: z \perp (x^{(k)} - x^{(k-1)})$. (3.4.9)



i.e.: $J_k := J_{k-1} + \frac{F(x^{(k)})(x^{(k)} - x^{(k-1)})^T}{\|x^{(k)} - x^{(k-1)}\|_2^2}$ (3.4.10)

Broydens Quasi-Newton Method for solving $F(x) = 0$:

$$x^{(k+1)} := x^{(k)} + \Delta x^{(k)}, \Delta x^{(k)} := -J_k^{-1} F(x^{(k)}), J_{k+1} := J_k + \frac{F(x^{(k+1)})(\Delta x^{(k)})^T}{\|\Delta x^{(k)}\|_2^2} \quad (3.4.11)$$

Initialize J_0 e.g. with the exact Jacobi matrix $DF(x^{(0)})$.

Remark 3.4.17 (Minimal property of Broydens rank 1 modification).

3.4 p. 218

Let $J \in \mathbb{R}^{n,n}$ fulfill (3.4.8) and $J_k, x^{(k)}$ from (3.4.11)

then $(I - J_k^{-1} J)(x^{(k+1)} - x^{(k)}) = -J_k^{-1} F(x^{(k+1)})$

3.4 p. 21

3.4 p. 22

and hence

$$\begin{aligned} \left\| \mathbf{I} - \mathbf{J}_k^{-1} \mathbf{J}_{k+1} \right\|_2 &= \left\| \frac{-\mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)}) \Delta \mathbf{x}^{(k)}}{\left\| \Delta \mathbf{x}^{(k)} \right\|_2^2} \right\|_2 = \left\| \left(\mathbf{I} - \mathbf{J}_k^{-1} \mathbf{J} \right) \frac{\Delta \mathbf{x}^{(k)} (\Delta \mathbf{x}^{(k)})^T}{\left\| \Delta \mathbf{x}^{(k)} \right\|_2^2} \right\|_2 \\ &\leq \left\| \mathbf{I} - \mathbf{J}_k^{-1} \mathbf{J} \right\|_2. \end{aligned}$$

In conclusion,

(3.4.10) gives the $\|\cdot\|_2$ -minimal relative correction of \mathbf{J}_{k-1} , such that the secant condition (3.4.8) holds.

△

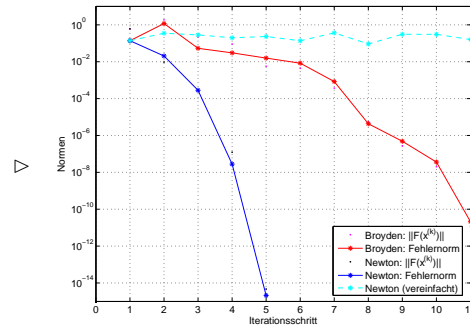
Example 3.4.18 (Broydens Quasi-Newton Method: Convergence).

- In the non-linear system of the example 3.4.1, $n = 2$ take $\mathbf{x}^{(0)} = (0.7, 0.7)^T$ and $\mathbf{J}_0 = DF(\mathbf{x}^{(0)})$

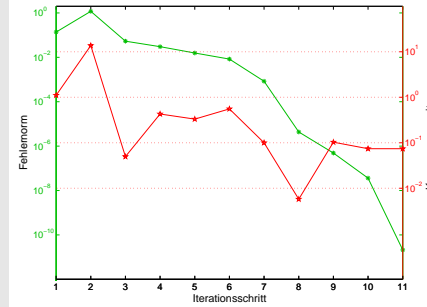
The numerical example shows that the method is:

slower than Newton method (3.4.1), but

better than simplified Newton method (see remark. 3.4.5)



▷



convergence monitor

=

quantity that displays difficulties in the convergence of an iteration

Here:

$$\mu := \frac{\left\| \mathbf{J}_{k-1}^{-1} F(\mathbf{x}^{(k)}) \right\|}{\left\| \Delta \mathbf{x}^{(k-1)} \right\|}$$

Heuristics: no convergence whenever $\mu > 1$

◇

Remark 3.4.19. Option: damped Broyden method (as for the Newton method, section 3.4.4)

△

3.4
p. 221

Implementation of (3.4.11): with Sherman-Morrison-Woodbury Update-Formula

$$\mathbf{J}_{k+1}^{-1} = \left(\mathbf{I} - \frac{\mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)}) (\Delta \mathbf{x}^{(k)})^T}{\left\| \Delta \mathbf{x}^{(k)} \right\|_2^2 + \Delta \mathbf{x}^{(k)} \cdot \mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)})} \right) \mathbf{J}_k^{-1} = \left(\mathbf{I} + \frac{\Delta \mathbf{x}^{(k+1)} (\Delta \mathbf{x}^{(k)})^T}{\left\| \Delta \mathbf{x}^{(k)} \right\|_2^2} \right) \mathbf{J}_k^{-1} \quad (3.4.12)$$

that makes sense in the case that

$$\left\| \mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)}) \right\|_2 < \left\| \Delta \mathbf{x}^{(k)} \right\|_2$$

"simplified Quasi-Newton correction"

3.4
p. 222

3.4
p. 22

3.4
p. 22

MATLAB-CODE: Broyden method (3.4.11)

```
function x = broyden(F,x,J,tol)
k = 1;
[L,U] = lu(J);
s = U\ (L\F(x)); sn = dot(s,s);
dx = [s]; dxn = [sn];
x = x - s; f = F(x);

while (sqrt(sn) > tol), k=k+1
w = U\ (L\f);
for l=2:k-1
w = w+dx(:,l)*(dx(:,l-1)'*w)...
/dxn(l-1);
end
if (norm(w)>=sn)
warning('Dubious step %d!',k);
end
z = s'*w; s = (1+z/(sn-z))*w; sn=s'*s;
dx = [dx,s]; dxn = [dxn,sn];
x = x - s; f = F(x);
end
```

unique LU-decomposition !
store $\Delta x^{(k)}, \|\Delta x^{(k)}\|_2^2$
(see (3.4.12))
solve two SLEs
Termination, see 3.4.2
construct $w := J_k^{-1}F(x^{(k)})$
(\rightarrow recursion (3.4.12))
convergence monitor
 $\frac{\|J_{k-1}^{-1}F(x^{(k)})\|}{\|\Delta x^{(k-1)}\|} < 1 ?$
correction $s = J_k^{-1}F(x^{(k)})$

Computational cost : $\bullet O(N^2 \cdot n)$ operations with vectors, (Level I)
 N steps
 \bullet 1 LU-decomposition of J , $N \times$ solutions of SLEs, see section 2.2
 $\bullet N$ evaluations of F !

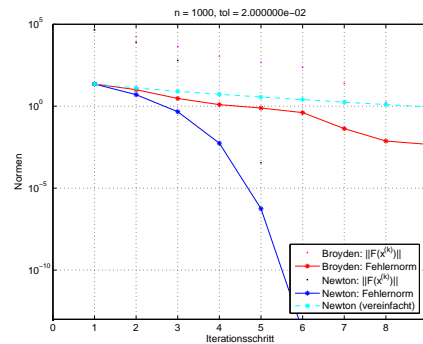
Memory cost : \bullet LU-factors of J + auxiliary vectors $\in \mathbb{R}^n$
 N steps
 $\bullet N$ vectors $x^{(k)} \in \mathbb{R}^n$

Example 3.4.20 (Broyden method for a large non-linear system).

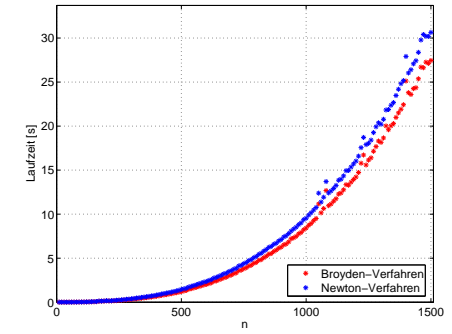
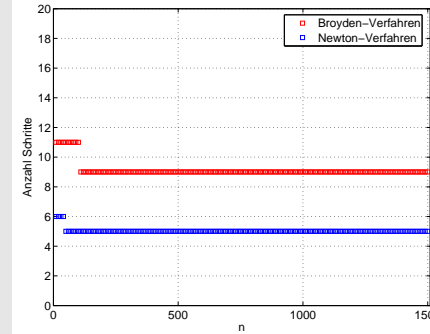
$$F(x) = \begin{cases} \mathbb{R}^n \mapsto \mathbb{R}^n \\ x \mapsto \text{diag}(x)Ax - b, \\ b = (1, 2, \dots, n) \in \mathbb{R}^n, \\ A = I + aa^T \in \mathbb{R}^{n,n}, \\ a = \frac{1}{\sqrt{1 \cdot b - 1}}(b - 1). \end{cases}$$

The interpretation of the results resemble the example 3.4.18 \triangleright

$h = 2/n; x_0 = (2:h:4-h)'$;



Efficiency comparison: Broyden method \longleftrightarrow Newton method:
(in case of dimension n use tolerance $\text{tol} = 2n \cdot 10^{-5}$, $h = 2/n$; $x_0 = (2:h:4-h)'$;)



In conclusion,
the Broyden method is worthwhile for dimensions $n \gg 1$ and low accuracy requirements.



3.5 Essential Skills Learned in Chapter 3

You should know:

- what is a linear convergent iteration, its rate and dependence of the choice of the norm
- what is the the order of convergence and how to recognize it from plots or from error data
- possible termination criteria and their risks
- how to use fixed-point iterations; convergence criteria
- bisection-method: pros and contras
- Newton-iteration: pros and contras
- the idea behind multi-point methods and an example
- how to use the Newton-method in several dimensions and how to reduce its computational effort (simplified Newton, quasi-Newton, Broyden method)