

# Numerical Methods for Computational Science and Engineering

Prof. R. Hiptmair, SAM, ETH Zurich

(with contributions from Prof. P. Arbenz and Dr. V. Gradinaru)

Autumn Term 2016

(C) Seminar für Angewandte Mathematik, ETH Zürich

URL: <http://www.sam.math.ethz.ch/~hiptmair/tmp/NumCSE/NumCSE16.pdf>

## VII. Iterative Methods for non-linear Systems of Equations

$$F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n : F(\underline{x}) = \underline{0}$$

↑      ↑  
Same number of unknowns  
and equations

▷ No general theory  
o

▷ vector of unknowns

## 2. 1. Iterative Methods

An iterative method for (approximately) solving the non-linear equation  $F(\underline{x}) = \underline{0}$  is an algorithm generating an arbitrarily long sequence  $(\underline{x}^{(k)})_k$  of approximate solutions.

$\underline{x}^{(k)} \doteq k\text{-th iterate}$

Initial guess

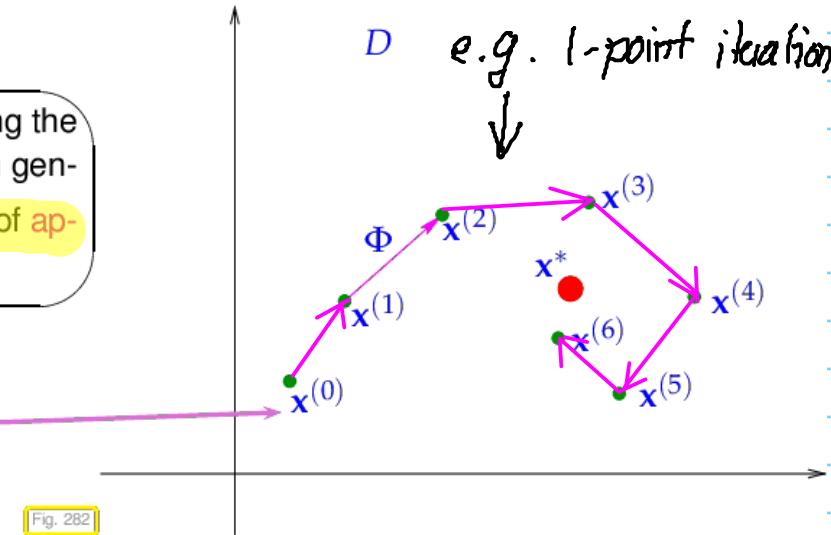


Fig. 282

Iteration error :  $e^{(k)} := \underline{x}^{(k)} - \underline{x}^*$ ,  $\underline{x}^* \doteq$  solution

More concrete :

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

↳ iteration function

Needed initial guesses  $\underline{x}^{(0)}, \dots, \underline{x}^{(m-1)}$

- Issues : • Convergence :  $\underline{x}^{(k)} \rightarrow \underline{x}^*$
- Consistency :  $F(\underline{x}^*) = \underline{0} \stackrel{?}{\iff} \Phi(\underline{x}^*, \dots, \underline{x}^*) = \underline{x}^*$

[  $\Phi$  cont. &  $\underline{x}^{(k)}$  cvg. & consistent  $\Rightarrow$  limit is solution :

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

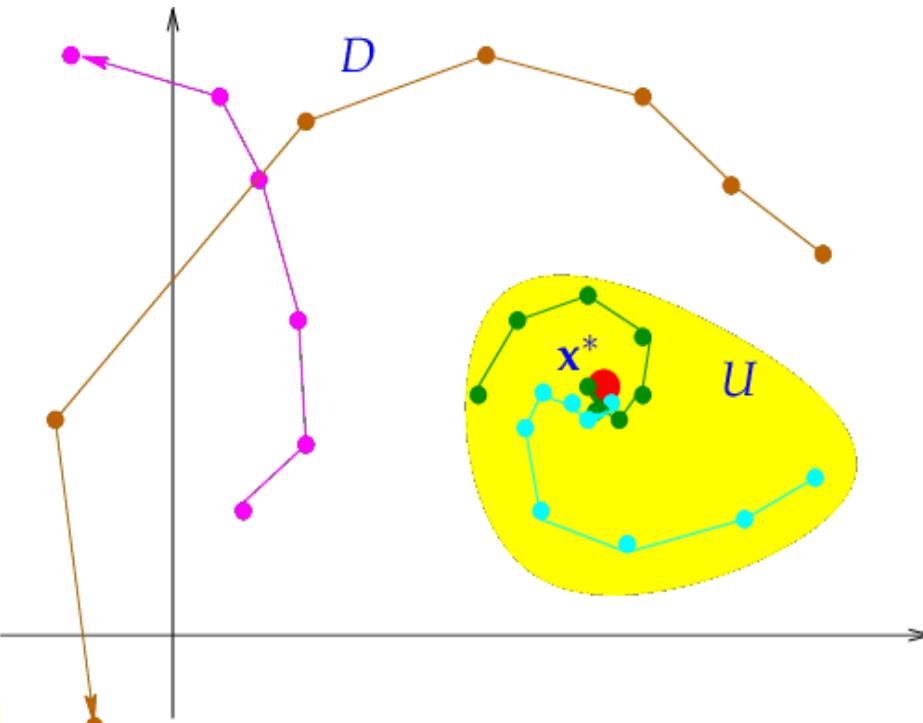
D e.g. 1-point iteration

2

$$\underline{x}^* = \lim_{k \rightarrow \infty} \phi(x^{(k-1)}, \dots, x^{(k-m)}) \\ \stackrel{\text{cont.}}{=} \phi_F(x^*, \dots, x^*)$$

- Speed of conv.: How fast  $\|\underline{x}^{(k)} - \underline{x}^*\| \rightarrow 0$  for some norm on  $\mathbb{R}^n$

Much depends on initial guess:



### Definition 8.1.8. Local and global convergence $\rightarrow [?, \text{Def. 17.1}]$

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

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

where  $(\underline{x}^{(k)})_{k \in \mathbb{N}_0}$  is the (infinite) sequence of iterates.

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

### 2.1.1. Speed of Convergence

#### Definition 8.1.9. Linear convergence

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

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

smallest  $L$ : rate

How to tell linear conv. in an experiment?

Assume sharpness  $E_k := \|\underline{x}^{(k)} - \underline{x}^*\| \approx L \|\underline{x}^{(k-1)} - \underline{x}^*\| = L E_{k-1}$

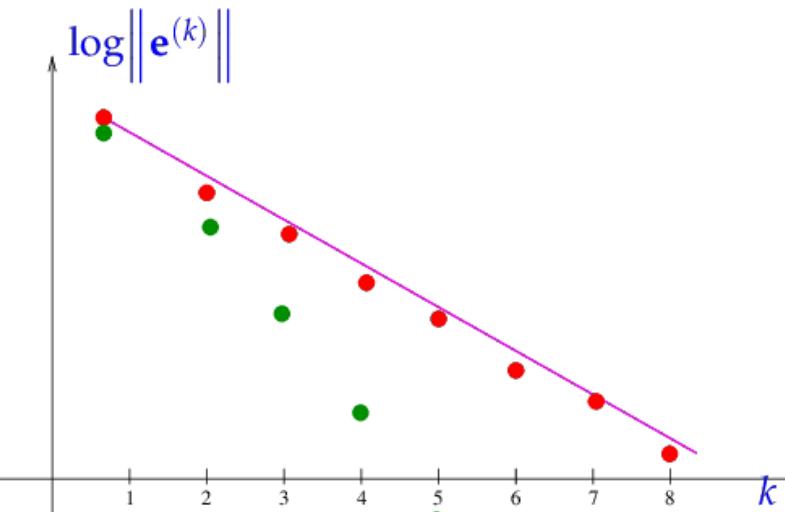
Visual evidence of linear conv.:

$$E_k \propto L E_{k-1}$$

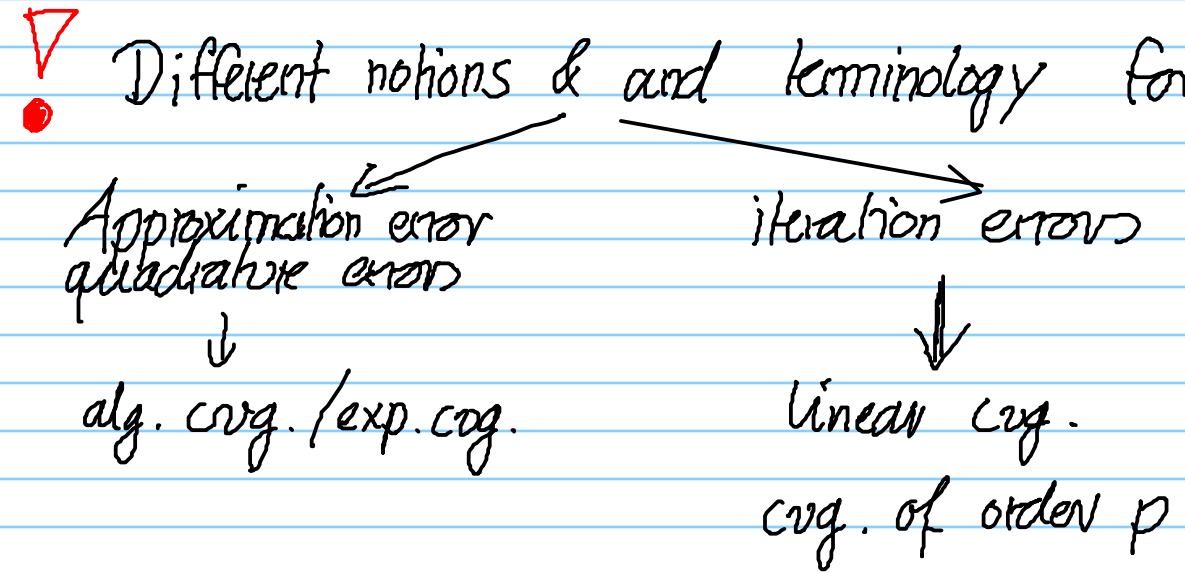
$$\log \cdot | \quad E_k \propto L^k E_0 \Rightarrow \log E_k \approx k \log L + \log E_0$$

③

$$\log \cdot | \quad \varepsilon_k \propto L^k \varepsilon_0 \Rightarrow \underbrace{\log \varepsilon_k \approx k \log L + \log \varepsilon_0}_{\text{Data points on straight line}}$$



Data points on straight line  
in  $k$  vs.  $\log \varepsilon_k$  plot.  
\* slope  $\sim$  rate



"Faster" than linear cog:

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).

How to tell cog of order  $p > 1$ ?

[assuming sharpness of estimate]

$$\cdot \log | \quad \varepsilon_k \propto C \varepsilon_{k-1}^p$$

$$\log \varepsilon_k \approx \log C + p \log \varepsilon_{k-1}$$

$$\log \varepsilon_{k-1} \approx \log C + p \log \varepsilon_{k-2}$$

$$p \approx \frac{(\log \varepsilon_k - \log \varepsilon_{k-1})}{(\log \varepsilon_{k-1} - \log \varepsilon_{k-2})}$$

study these quotients

(4)

$$\text{NLSE : } F(\underline{x}) = 0, F: \mathbb{R}^n \rightarrow \mathbb{R}^n$$

Iteration produces sequence  $(\underline{x}^{(k)})_{k \in \mathbb{N}_0}$

Linear Cvg :  $\|\underline{x}^{(k+1)} - \underline{x}^*\| \leq L \|\underline{x}^{(k)} - \underline{x}^*\|$  for  $0 \leq L < 1$

order-p Cvg. :  $\|\underline{x}^{(k+1)} - \underline{x}^*\| \leq C \|\underline{x}^{(k)} - \underline{x}^*\|^p, C > 0$   
 $(p \geq 1)$

Example :  $\sqrt{\cdot}$ -iteration ( $n = 1$ )

$$\underline{x}^{(k+1)} = \frac{1}{2} \left( \underline{x}^{(k)} + \frac{a}{\underline{x}^{(k)}} \right) \Rightarrow |\underline{x}^{(k+1)} - \sqrt{a}| = \frac{1}{2\underline{x}^{(k)}} |\underline{x}^{(k)} - \sqrt{a}|^2. \quad (8.1.21)$$

$a > 0$

$$1\text{-point iteration : } \phi(t) = \frac{1}{2} \left( t + \frac{a}{t} \right)$$



Convergence to  $\sqrt{a}$  with order 2

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

Note: Convergence of order  $p > 1$  guarantees local convergence

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

If  $\underline{x}^{(0)} \in M := \{ \underline{z} : C \|\underline{z} - \underline{x}^*\|^{p-1} \leq 1 \}$   
 $\Rightarrow \underline{x}^{(k)} \rightarrow \underline{x}^* \text{ for } k \rightarrow \infty$

$k$	$\underline{x}^{(k)}$	$e^{(k)} := \underline{x}^{(k)} - \sqrt{2}$	$\log \frac{ e^{(k)} }{ e^{(k-1)} } : \log \frac{ e^{(k-1)} }{ e^{(k-2)} }$
0	2.00000000000000000000	0.58578643762690485	
1	1.50000000000000000000	0.08578643762690485	
2	1.41666666666666652	0.00245310429357137	1.850
3	1.41421568627450966	0.00000212390141452	1.984
4	1.41421356237468987	0.00000000000159472	2.000
5	1.41421356237309492	0.00000000000000022	0.630

Note the doubling of the number of significant digits in each step!

[impact of roundoff !]

$$\text{Relative error : } \underline{x}^{(k)} = \underline{x}^* (1 + \delta_k)$$

$$\text{order-2 cvg : } \|\underline{x}^{(k+1)} - \underline{x}^*\| \leq C \|\underline{x}^{(k)} - \underline{x}^*\|^2$$

$$\delta_{k+1} \|\underline{x}^*\| \leq C \delta_k^2 \|\underline{x}^*\|^2$$

$$\text{"sharpness"} \quad \delta_{k+1} \approx C \|\underline{x}^*\| \delta_k^2$$

$$\text{IF } C \|\underline{x}^*\| \approx 1, \delta_k \approx 10^{-\ell} \Rightarrow \delta_{k+1} \approx 10^{-2\ell}$$

⑤

## 2.1.2 Termination

→ Stopping rules

Ideal:

STOP, if  $\|x^{(k)} - x^*\| \leq T_{rel} \|x^*\|$  : relative error

$\|x^{(k)} - x^*\| \leq T_{abs}$  : absolute error

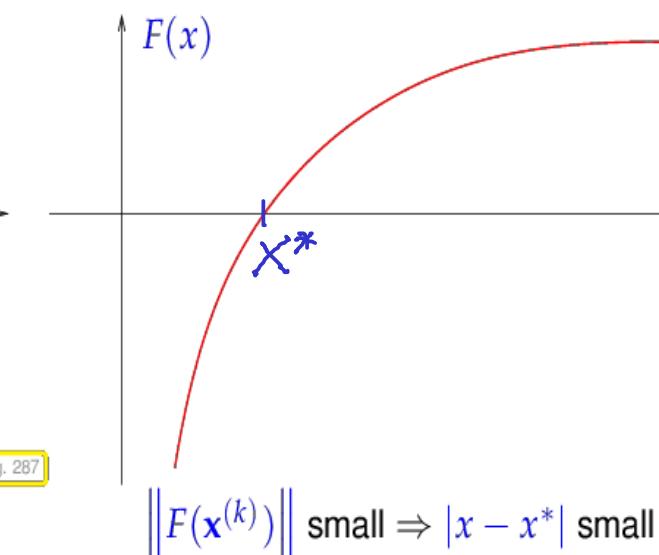
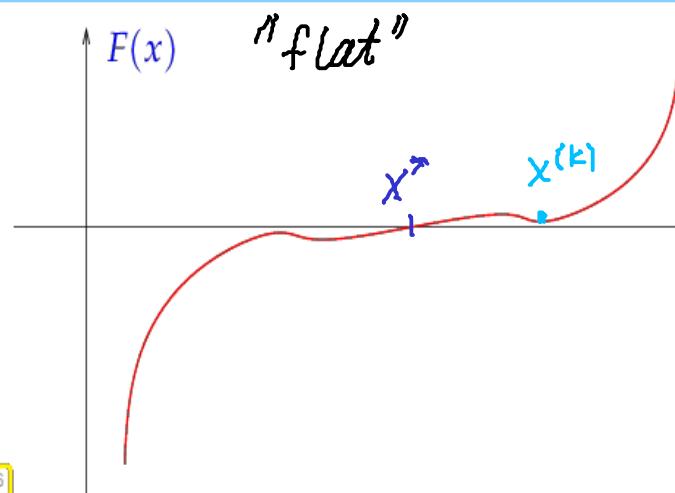
↑ not accessible

Practical:

### ① Residual based termination

STOP, if  $\|F(x^{(k)})\| \leq \varepsilon$

tells little about error



### ② Correction based termination

STOP, if  $\|x^{(k+1)} - x^{(k)}\| \leq \begin{cases} T_{rel} \|x^{(k+1)}\| \\ T_{abs} \end{cases}$

→ Still no guarantee, except for linearly org. iterations with known rate

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

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

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

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

► Iterates satisfy:

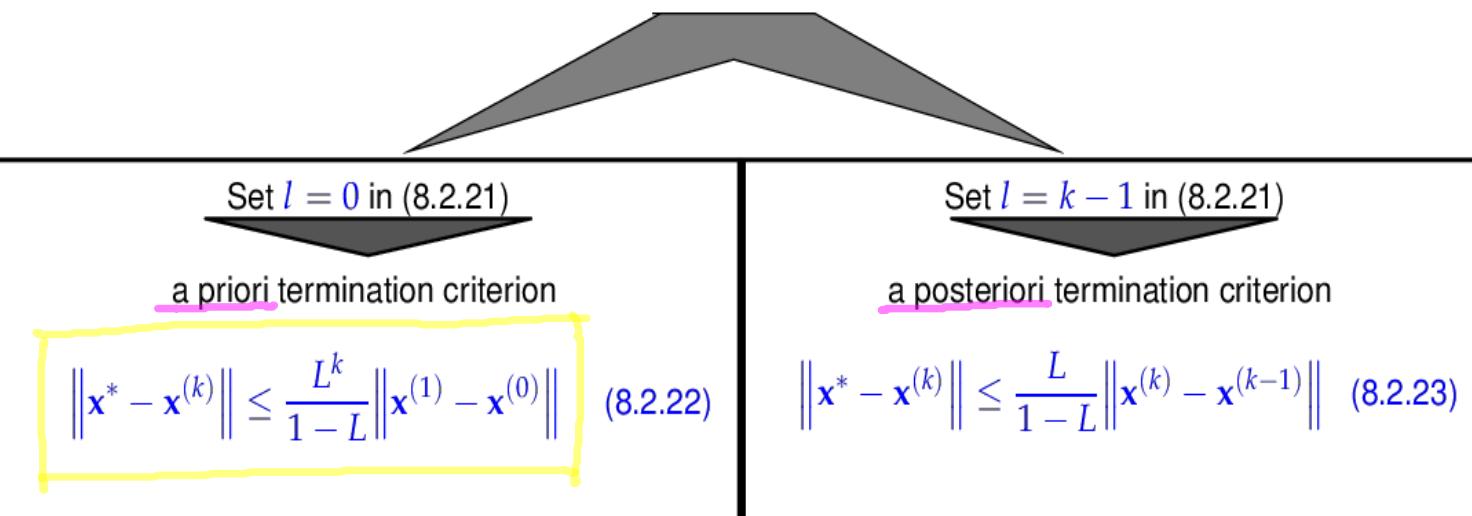
$$\|x^{(k+1)} - x^*\| \leq \frac{L}{1-L} \|x^{(k+1)} - x^{(k)}\|. \quad (8.1.29)$$

Upper bound  $\hat{\varepsilon} < 1$  for  $L$  sufficient to obtain reliable bound!

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

6

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



$$\begin{aligned} \Phi_2(x) = x &\Leftrightarrow 1+x = (1+e^x)x \Leftrightarrow 1 = xe^x \\ &\Leftrightarrow F(x) = 0 \end{aligned}$$

•

## 2.2. Fixed Point Iterations

$$\mathbf{x}^{(k+1)} = \Phi_F(\mathbf{x}^{(k)})$$

$\Delta \Phi: M \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  iteration function

consistent:  
with  $F(x) = 0$

$$\Phi_F(x^*) = x^* \Leftrightarrow F(x^*) = 0$$

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

Different fixed point forms:

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

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

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

[all consistent]  
 $x^{(0)} = 0.5$ :

$\Phi_2$ :

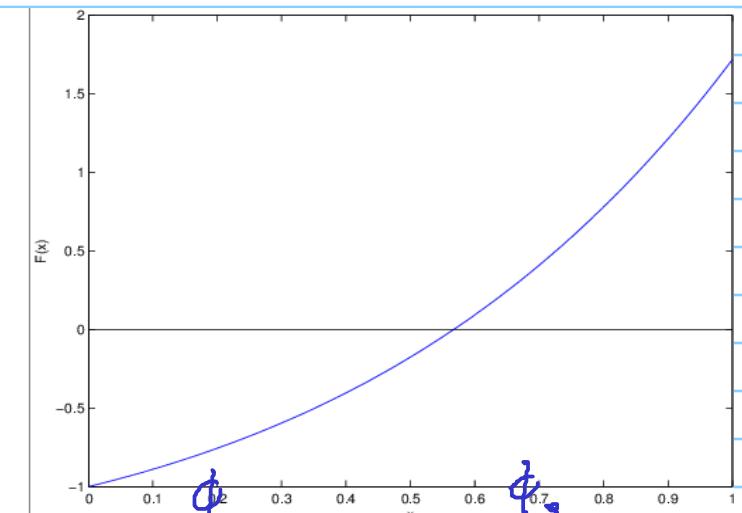
order 2-cvg

$\Phi_1$ :

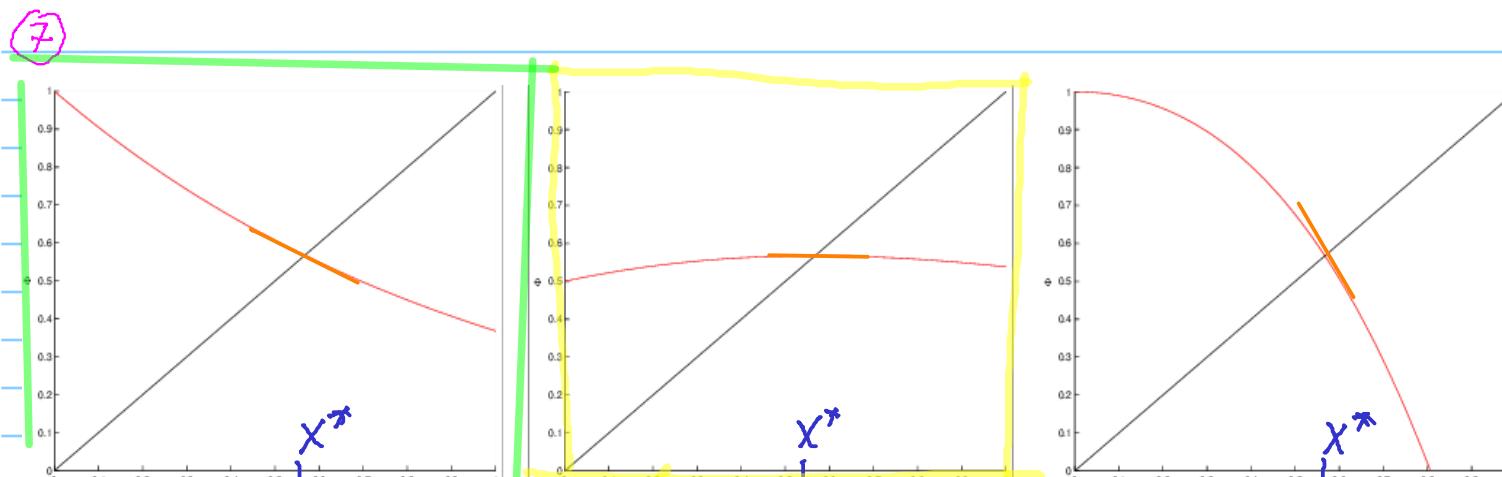
lin cvg.

$\Phi_3$ :

no cvg.



$k$	$ 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



$$|\Phi'(x^*)| < 1$$

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

$$|\Phi'(x^*)| > 1$$

► Slope of  $\Phi$  at  $x^*$  matters! Taylor argument

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

$$= \Phi'(x^*)(x^{(k)} - x^*) + \frac{1}{2}\Phi''(x^*)(x^{(k)} - x^*)^2 + \underbrace{\frac{1}{6}\Phi'''(\bar{x})(x^{(k)} - x^*)^3}_{\text{neglect}}$$

If  $|x^{(k)} - x^*| \ll 1$   $\longrightarrow$  neglect

• and  $\Phi'(x^*) = 0$   $\Rightarrow$   $|x^{(k+1)} - x^*| \approx |\frac{1}{2}\Phi''(x^*)| |x^{(k)} - x^*|^2$   
 $(\Phi'' \text{ flat})$   $\cong$  order  $\sim 2$  conv.

• and  $|\Phi'(x^*)| < 1$   $\Rightarrow$   $|x^{(k+1)} - x^*| \leq |\Phi'(x^*)| |x^{(k)} - x^*|$   
 $+ \text{small terms}$   
 $\cong$  Linear convergence with rate  $\approx |\Phi'(x^*)|$

## Generalization to n-dim. :

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

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

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

converges locally and at least linearly.

predictor for rate of lin. conv.

If  $\mathbf{D}\Phi(x^*) = 0 \Rightarrow$  local order  $\sim 2$  conv.

↑  
 matrix norm, Def. 1.5.76!

⑧

## 2.3. Zero Finding

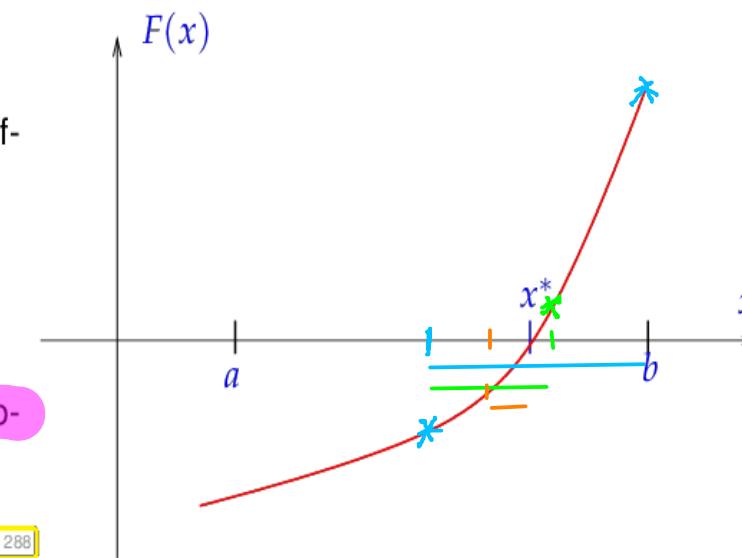
Find  $x^* \in I \subset \mathbb{R}$  :  $F(x^*) = 0$  for  $F: I \rightarrow \mathbb{R}$   
 [  $F$  continuous ]

### 2.3.1. Bisection

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

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

as we conclude from the intermediate value theorem.



Convergence: "linear type"

Intervals shrink by factor 2 in each step

$$|a^{(k)} - b^{(k)}| \leq 2^{-k} |a^{(0)} - b^{(0)}|$$

$$\Rightarrow |x^{(k)} - x^*| \leq 2^{-k} |a^{(0)} - b^{(0)}|$$

$\Rightarrow$  robust method

### C++11 code 8.3.2: Bisection method for solving $F(x) = 0$ on $[a, b]$

```

2 // Searching zero of F in [a,b] by bisection
3 template <typename Func>
4 double bisect(Func& F, double a, double b, double tol)
5 {
6     if (a > b) std::swap(a,b); // sort interval bounds
7     double fa = F(a), fb = F(b);
8     if (fa*fb > 0) throw "f(a) and f(b) have same sign";
9     int v=1; if (fa > 0) v=-1;
10    double x = 0.5*(b+a); // determine midpoint
11    // termination, relies on machine arithmetic if tol = 0
12    while (b-a > tol && ((a<x) && (x<b))) //
13    {
14        // sgn(f(x)) = sgn(f(b)), then use x as next right boundary
15        if (v*F(x) > 0) b=x;
16        // sgn(f(x)) = sgn(f(a)), then use x as next left boundary
17        else a=x;
18        x = 0.5*(a+b); // determine next midpoint
19    }
20    return x;
21 }
```

$tol = 0$  : Loop will stop, when

$$a \approx \frac{1}{2}(b-a) = a$$

$\Rightarrow$  happens, if  $|1/2(b-a)| \leq EPS \cdot |a|$

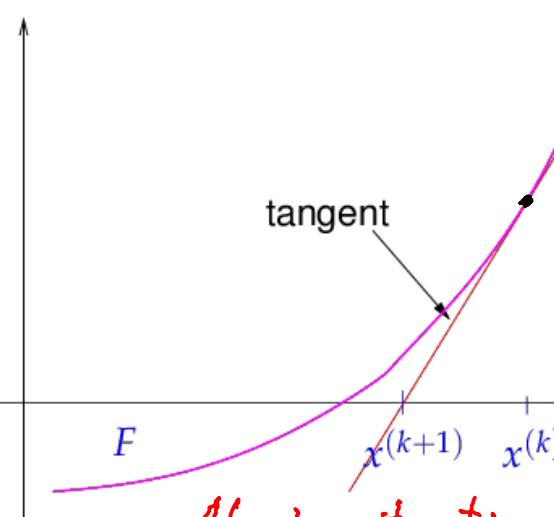
### ⑨ 2.3.2 Model Function Methods

Idea: Replace  $F$  with  $\tilde{F}: I \rightarrow \mathbb{R}$  (model function)

- depending on point values of  $F / F'$
- whose zeros are easy to compute

→ do this in every step of an iteration

#### 2.3.2.1. Newton's method



Newton iteration

Model function  $\tilde{F}$

≡ tangent at graph of  $F$   
in  $x^{(k)}$

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

$$x^{(k+1)} : \tilde{F}(x^{(k+1)}) = 0$$

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

[if  $F'(x^{(k)}) \neq 0$ ]

↪ ≡ fixed point iteration with  $\phi(x) = x - \frac{F(x)}{F'(x)}$

(Local) crg. of Newton's method:

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

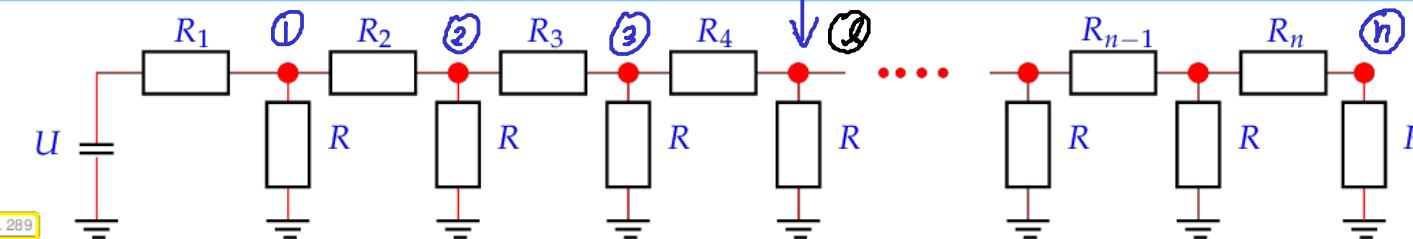
▷  $F(x^*) = 0 \Rightarrow \phi'(x^*) = 0$ , if  $F'(x^*) \neq 0$   
 $\Rightarrow$  local 2nd-order crg.

#### Convergence of Newton's method in 1D

Newton's method locally quadratically converges ( $\rightarrow$  Def. 8.1.17) to a zero  $x^*$  of  $F$ , if  $F'(x^*) \neq 0$

Example:  $F'$  by implicit differentiation

achieve prescribed voltage = 1 here



[Nodal analysis]  $\Rightarrow$  LSE for nodal voltages  $u_i$

$$(A + xI) \underline{u}(x) = \underline{b} \in \mathbb{R}^n$$

↑  
↓ n.v.  
s.p.d. tridiagonal  
matrix  $\in \mathbb{R}^{n \times n}$   
=  $\mathcal{Y}R$  [→ unknown]

(1D)

Task:  $x \in \mathbb{R} : \underline{w}^T \underline{\mu}(x) = 1, \underline{w} \in \mathbb{Q}_d$

⇒

Find zero of  $F(x) = \underline{w}^T \underline{\mu}(x) - 1 = \underline{w}^T (A + xI)^{-1} \underline{b} - 1$

$$F'(x) = \underline{w}^T \underline{\mu}'(x)$$

$$\underline{v}(x) : (A + xI) \cdot \underline{\mu}(x) = \underline{b} \quad | \frac{d}{dx}$$

$$I \cdot \underline{\mu}(x) + (A + xI) \cdot \underline{\mu}'(x) = 0$$

with product rule:  $(fg)' = f' \cdot g + f \cdot g'$

$$\Rightarrow \underline{\mu}'(x) = - (A + xI)^{-1} \underline{\mu}(x)$$

▷ Newton iteration:

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

$$= x^{(k)} - \frac{\underline{w}^T \underline{\mu}(x^{(k)}) - 1}{\underline{w}^T \underline{\mu}}$$

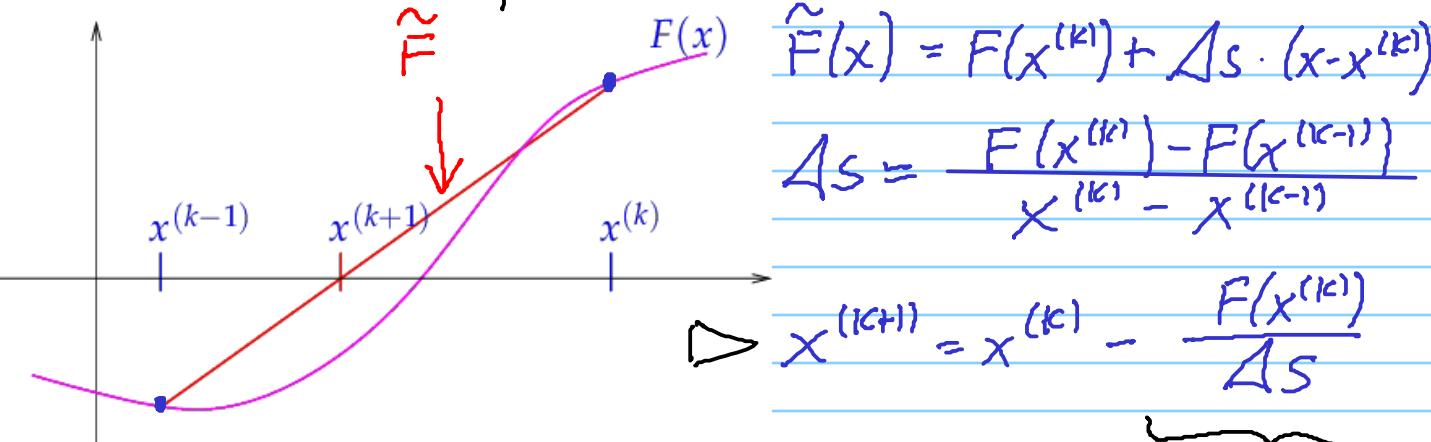
where:  $(A + x^{(k)}I) \underline{\mu}(x^{(k)}) = \underline{b}$

$$(A + x^{(k)}I) \underline{\mu} = - \underline{\mu}(x^{(k)}) \leftarrow \\ := \underline{\mu}'(x^{(k)})$$

### 2.3.2.2. Multi-point methods

m-point method:  $x^{(k+1)} = \phi(x^{(k)}, \dots, x^{(k-m+1)})$   
 ↳ m initial  $x^{(j)}$  required

Secant method: 2-point method



C++11 code 8.3.25: Secant method for 1D non-linear equation

```

2 // Secant method for solving F(x)=0 for F:D ⊂ ℝ → ℝ,
3 // initial guesses x0, x1,
4 // tolerances atol (absolute), rtol (relative)
5 template<typename Func>
6 double secant(double x0, double x1, Func&& F,
7                 double rtol, double atol, unsigned int maxIt)
8 {
9     double fo = F(x0);
10    for (unsigned int i=0; i<maxIt; ++i) {
11        double fn = F(x1);
12        double s = fn*(x1-x0)/(fn-fo); // secant correction
13        x0 = x1; x1 = x1-s;
14        // correction based termination (relative and absolute)
15        if (abs(s) < max(atol, rtol*min(abs(x0), abs(x1)))) {
16            return x1;
17            fo = fn;
18        }
19    }
20    return x1;
}

```

| F-evaluation per step

## (Empiric) convergence of secant method :

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

↓  
estimate for  
order of conv.

$k$	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log  e^{(k+1)}  - \log  e^{(k)} }{\log  e^{(k)}  - \log  e^{(k-1)} }$
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	0.00000000134427	
11	0.56714329040978	-0.00000000000001	-0.00000000000000	

→ order 1.6 convergence

Generalization: Inverse interpolation

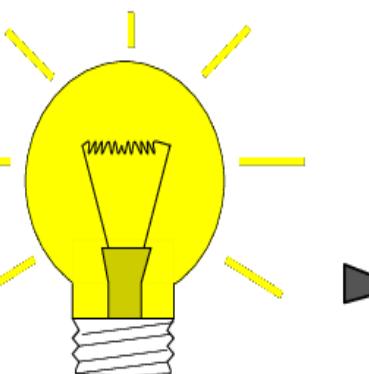
Assume:

$F : I \subset \mathbb{R} \mapsto \mathbb{R}$  one-to-one (monotone)

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

- Interpolate  $F^{-1}$  by polynomial  $p$  of degree  $m-1$  determined by

$$p(F(x^{(k-j)})) = x^{(k-j)}, \quad j = 0, \dots, m-1.$$

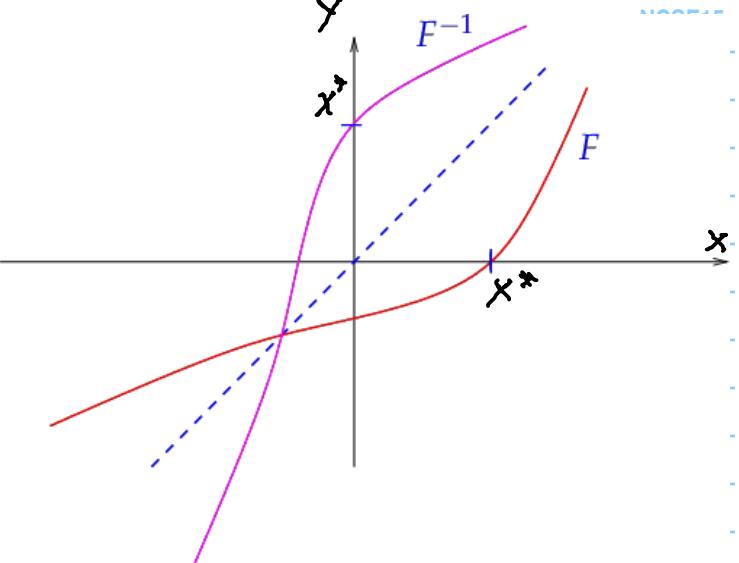


- New approximate zero  $x^{(k+1)} := p(0)$

The graph of  $F^{-1}$  can be obtained by reflecting the graph of  $F$  at the angular bisector. ▷

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

Fig. 292



Case  $m = 2$  (2-point method)

➢ secant method

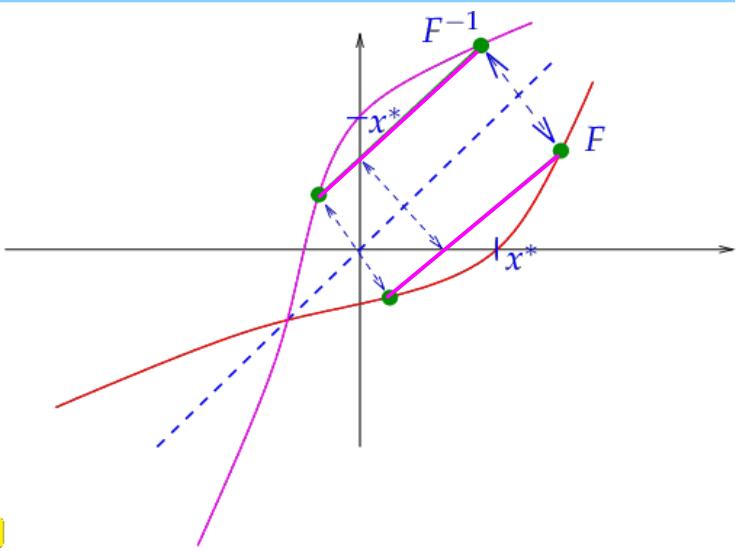
The interpolation polynomial is a line. In this case we do not get a new method, because the inverse function of a linear function (polynomial of degree 1) is again a polynomial of degree 1.

Fig. 293

$m = 3$ :

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

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



(12)

$$F(x) = xe^x - 1, \quad x^{(0)} = 0, x^{(1)} = 2.5, x^{(2)} = 5.$$

$k$	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log  e^{(k+1)}  - \log  e^{(k)} }{\log  e^{(k)}  - \log  e^{(k-1)} }$
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	0.00000000000003	0.00000000000001	

↓  
fractional order

### 2.3.3. Asymptotic efficiency

$$\frac{\text{gain}}{\text{work}} \leftarrow \frac{\text{prescribed error reduction}}{\# F\text{-eval} + \# F'\text{-eval}}$$

[no of digits]

gain:  $\boxed{\|e^{(k)}\| \leq \beta \|e^{(0)}\|}$  with  $\beta \leq 1$  ( $\times$ )

$K(\beta)$ : minimal no. of iterations required for ( $\times$ )

Efficiency:  $\frac{|\log \beta|}{K(\beta) \cdot W}$ ,  $W \equiv \text{work/iteration}$

• Linearly conjg. iteration:  $\|e^{(k)}\| \leq L^k \|e^{(0)}\|$ ,  $L < 1$

$\Rightarrow K(\beta) \geq \frac{\log \beta}{\log L} \Rightarrow \text{Eff.} = \frac{|\log L|}{W}$

• order  $p > 1$ :  $\|e^{(k)}\| \leq C \|e^{(k-1)}\|^p$

$$\|e^{(k)}\| \leq \underbrace{C^{1+p+p^2+\dots+p^{k-1}}}_{= C^{\frac{p^k-1}{p-1}}} \|e^{(0)}\|^{p^{k-1}} \cdot \|e^{(0)}\|$$

$\triangleright K(\beta) : \underbrace{\left( C^{\frac{1}{p-1}} \|e^{(0)}\| \right)^{p^{k-1}}}_{=: L_0} \leq \beta$

$$p^k \geq \frac{\log \beta}{\log L_0} + 1$$

$$k \geq \log \left( \frac{\log \beta}{\log L_0} + 1 \right) / \log p$$

$\beta \ll 1$ :  $k \cdot \log p \geq \log \log \beta - \log \log L_0$

$\therefore K(\beta) = \frac{\log \log \beta}{\log p}$

$$\text{Efficiency} = \frac{\log p}{W} \frac{\log \beta}{\log \log \beta}$$

$\hookrightarrow$  different for different methods

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

$$W_{\text{Newton}} = 2W_{\text{secant}}, \\ p_{\text{Newton}} = 2, p_{\text{secant}} = 1.62$$

(B)

► "Secant method is 30% more efficient than Newton"

## 2.4. Newton's Method (for $n \geq 1$ )

$$\underline{x}^* \in \mathbb{R}^n : F(\underline{x}^*) = 0, F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$\uparrow$   
differentiable

"Tangent plane" at  $\underline{x}^{(k)}$   $\hat{=}$  model function by linearization

$$\underline{x}^{(k+1)} : \hat{F}(\underline{x}) = F(\underline{x}^{(k)}) + DF(\underline{x}^{(k)})(\underline{x} - \underline{x}^{(k)}) \stackrel{!}{=} 0$$

$\uparrow$   
Jacobian  $\in \mathbb{R}^{n,n}$

Newton iteration

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - \underbrace{DF(\underline{x}^{(k)})^{-1} F(\underline{x}^{(k)})}_{\text{Newton correction}}$$

Newton correction  $\Delta \underline{x}^{(k)}$ : by solving a LSE

Affine invariance:

Newton iterations for  $F(\underline{x}) = 0$   
 $AF(\underline{x}) = 0, A \in \mathbb{R}^{n,n}$  regular

produce the same iterates, if started with the same value

$$[ G(\underline{x}) = AF(\underline{x}) \Rightarrow DG(\underline{x}) = ADF(\underline{x}) ]$$

N.I. for  $G(\underline{x}) = 0$ :

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - [ \cancel{ADF(\underline{x}^{(k)})} ]^{-1} [ \cancel{AF(\underline{x}^{(k)})} ]$$

### 2.4.2. Convergence of Newton's method

→ = Fixed point iteration with iteration function

$$\phi(\underline{x}) := \underline{x} - DF(\underline{x})^{-1} F(\underline{x}), \underline{x} \in D$$

[product rule]

$$D\phi(\underline{x}^*) = I - ? \cdot \underbrace{F(\underline{x}^*)}_{=0} - DF(\underline{x})^{-1} DF(\underline{x}) = 0$$

⇒ local quadratic conv.

Example: Quasi-linear system of equation

$$A(\underline{x})\underline{x} = \underline{b}, A: \mathbb{R}^n \rightarrow \mathbb{R}^{n,n}$$

(14)

$$\mathbf{A}(\underline{x})\underline{x} = \mathbf{b}, \quad \mathbf{A}(\underline{x}) = \begin{pmatrix} \gamma(\underline{x}) & 1 & & \\ 1 & \gamma(\underline{x}) & 1 & \\ & \ddots & \ddots & \ddots \\ & & \ddots & \ddots \\ & & & 1 & \gamma(\underline{x}) & 1 \\ & & & & 1 & \gamma(\underline{x}) \end{pmatrix} \in \mathbb{R}^{n \times n},$$

$$g(\underline{x}) = 3 + \|\underline{x}\|$$

$$\Leftrightarrow F(\underline{x}) = \underline{b} - \mathbf{A}(\underline{x})\underline{x} = \underline{0}$$

$$\mathbf{A}(\underline{x})\underline{x} = \mathbf{T}\underline{x} + \underline{x}\|\underline{x}\|_2, \quad \mathbf{T} := \begin{pmatrix} 3 & 1 & & & \\ 1 & 3 & 1 & & \\ & \ddots & 3 & \ddots & \\ & & \ddots & \ddots & \ddots \\ & & & 1 & 3 & 1 \\ & & & & 1 & 3 \end{pmatrix}.$$

$$\Rightarrow F(\underline{x}) = \underline{b} - \mathbf{T}\underline{x} - \|\underline{x}\|\cdot\underline{x}$$

$$\Rightarrow DF(\underline{x}) = -\mathbf{T} - D\{\underline{x} \rightarrow \|\underline{x}\|\cdot\underline{x}\}$$

$\nearrow =: g(\underline{x})$

- $DF(\underline{x})$  : (i) product rule  
(ii) direct computation

(8.4.20)

$$(g(\underline{x}))_i = \sqrt{x_1^2 + \dots + x_n^2} \cdot x_i$$

$$(Dg(\underline{x}))_{i,j} = \frac{\partial g_i(\underline{x})}{\partial x_j} = \begin{cases} * & j=i \\ \frac{x_j}{\sqrt{x_1^2 + \dots + x_n^2}} \cdot x_i, & j \neq i \end{cases}$$

$$\frac{\partial}{\partial x_i} (\underline{x}\|\underline{x}\|)_i = \underbrace{\sqrt{x_1^2 + \dots + x_n^2}}_{=\|\underline{x}\|} + x_i \frac{x_i}{\sqrt{x_1^2 + \dots + x_n^2}},$$

$$\frac{\partial}{\partial x_j} (\underline{x}\|\underline{x}\|)_i = x_i \frac{x_j}{\sqrt{x_1^2 + \dots + x_n^2}}, \quad j \neq i.$$

$$DF(\underline{x}) = \mathbf{T} + \|\underline{x}\|_2 \cdot \underline{x} \frac{\underline{x}^\top}{\|\underline{x}\|_2} = \left( \mathbf{A}(\underline{x}) + \frac{\underline{x}\underline{x}^\top}{\|\underline{x}\|_2} \right).$$

↑  
rank-1-modification of  
tridiagonal matrix  $A$

$\Rightarrow$  Use Sherman-Morrison-Woodbury formula for  
computing Newton correction: cost  $O(n)$

(15)

Example: Matrix inversion via Newton

$$A \in \mathbb{R}^{n,n} \text{ regular: } F(X) = A - X^{-1} = 0$$

$$F: \mathbb{R}^{n,n} \rightarrow \mathbb{R}^{n,n}$$

Newton iteration for  $F(X) = 0$

Trick: implicit differentiation

$$F(X) \cdot X = AX - I$$

[product rule]

"General derivative of  $\phi: V \rightarrow W$

$$\mathcal{D}\phi(v)h = \lim_{t \rightarrow 0} \frac{\phi(v+th) - \phi(v)}{t} \quad \vdash$$

↑ direction of differentiation

Product rule:  $F: D \subset V \mapsto W, G: D \subset V \mapsto U$  sufficiently smooth,  $b: W \times U \mapsto Z$  bilinear, ie., linear in each argument:

$$T(x) = b(F(x), G(x)) \Rightarrow \mathcal{D}T(x)h = b(\mathcal{D}F(x)h, G(x)) + b(F(x), \mathcal{D}G(x)h), \quad (8.4.9)$$

$h \in V, x \in D.$

$$(f \cdot g)' = f' \cdot g + f \cdot g'$$

$$\begin{array}{l} \frac{d}{dx} \square \\ \downarrow \end{array} \quad \begin{array}{l} F(X) \cdot X = AX - I \\ \mathcal{D}(F(X) \cdot X)H = \mathcal{D}(AX - I)H \end{array}$$

• General product rule with  
 $b \Leftrightarrow$  matrix multiplication.  
 •  $\mathcal{D}fX \rightarrow X \mathcal{D}H = H$

$H \in \mathbb{R}^{n,n}$

$$\mathcal{D}F(X)(H) \cdot X + F(X) \cdot H = A \cdot H$$

$$\triangleright \underbrace{\mathcal{D}F(X)(H)}_{\in \mathbb{R}^{n,n}} = (A \cdot H - F(X) \cdot H) X^{-1} = X^{-1} H X^{-1}$$

Newton iteration:  $\mathcal{D}F(X)$  is not a matrix!

$$X^{(k+1)} = X^{(k)} - S, \quad \mathcal{D}F(X^{(k)})(S) = F(X^{(k)})$$

$$(X^{(k)})^{-1} S (X^{(k)})^{-1} = A - (X^{(k)})^{-1}$$

$$\Leftrightarrow S = X^{(k)} A X^{(k)} - X^{(k)}$$

$$\triangleright X^{(k+1)} = 2X^{(k)} - X^{(k)} A X^{(k)}$$

## General notion of derivative and Newton's method

For  $f: \mathbb{R} \rightarrow \mathbb{R}$  smooth: local linear approximation

$$f(x+h) \approx f(x) + f'(x)h, h \in \mathbb{R}$$

General:  $F: V \rightarrow W$ ,  $V, W \cong$  vector space

$$F(\underline{x}+\underline{h}) \approx F(\underline{x}) + DF(\underline{x})(\underline{h}), \forall \underline{h} \in V$$

$\uparrow$   
= linear mapping  $V \rightarrow W$

Note:

- $c \in \mathbb{R} \iff$  linear mapping  $\mathbb{R} \rightarrow \mathbb{R}$ ,  $t \mapsto c \cdot t$
- $V = W = \mathbb{R}^n$ : Jacobi matrix describes linear mapping  $DF(\underline{x}): \mathbb{R}^n \rightarrow \mathbb{R}^n$

Abstract Newton iteration for solving  $F(\underline{x}) = 0$

for  $F: D \subset V \rightarrow W$

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

$\rightarrow$  To state Newton iteration we have to find expressions for  $DF(\underline{x}^{(k)}) \underline{h}, \underline{h} \in V$

## Application: Derivative of matrix inversion

Tool: Product rule for general derivatives

- Vector spaces:  $V, W, U, Z$
- $F: V \rightarrow W, G: V \rightarrow U$  differentiable mappings
- $b: W \times U \rightarrow Z$  bilinear (linear in each argument)
- $T(\underline{x}) := b(F(\underline{x}), G(\underline{x})), T: V \rightarrow Z$

Simplest setting:  $V = W = U = Z = \mathbb{R}$ ,  $b(\underline{z}, \underline{y}) = \underline{z} \cdot \underline{y}$

$$T(\underline{x}) = F(\underline{x}) \cdot G(\underline{x}), \underline{x} \in \mathbb{R}$$

std. prod. rule  $\Rightarrow T'(\underline{x}) = F'(\underline{x}) \cdot G(\underline{x}) + F(\underline{x}) \cdot G'(\underline{x}) \in \mathbb{R}$

Regarded as linear mappings:  $T'(\underline{x}) \underline{h} = F'(\underline{x}) \underline{h} \cdot G(\underline{x}) + F(\underline{x}) \cdot (G'(\underline{x}) \underline{h})$

$\uparrow$  linear mapping  $\mathbb{R} \rightarrow \mathbb{R}$        $\uparrow \forall \underline{h} \in \mathbb{R} \rightarrow$

Product rule:  $F: D \subset V \mapsto W, G: D \subset V \mapsto U$  sufficiently smooth,  $b: W \times U \mapsto Z$  bilinear, i.e., linear in each argument:

linear map  $V \rightarrow W$   
 $\downarrow$   
 $T(\underline{x}) = b(F(\underline{x}), G(\underline{x})) \Rightarrow DT(\underline{x}) \underline{h} = b(DF(\underline{x}) \underline{h}, G(\underline{x})) + b(F(\underline{x}), DG(\underline{x}) \underline{h}),$  (8.4.9)  
 $\forall \underline{h} \in V, \underline{x} \in D. \quad \underline{w} \in W$   
 $\downarrow$   
linear map  $V \rightarrow U$

(17)

Matrix inversion :

$$\text{inv} := \begin{cases} \mathbb{R}_{\text{reg}}^{n,n} & \xrightarrow{\text{invertible } n \times n \text{ matrices}} \mathbb{R}^{n,n} \\ X & \xrightarrow{} X^{-1} \end{cases}$$

 $\Rightarrow \text{D}\text{inv}(X) = \text{linear mapping } \mathbb{R}^{n,n} \rightarrow \mathbb{R}^{n,n}$ 

Here: Basis free approach

Implicit differentiation &amp; product rule

diff.

$$\left[ \begin{array}{l} T(X) := \text{inv}(X) \cdot X = I \quad \left[ \begin{array}{l} \text{Here, } I \stackrel{!}{=} \text{identity matrix} \\ \text{is a const. mapping } \mathbb{R}^{n,n} \rightarrow \mathbb{R}^{n,n} \end{array} \right] \\ \rightarrow DT(X)H = \text{D}\text{inv}(X)(H) \cdot X + \text{inv}(X) \cdot H = 0 \\ \boxed{\text{D}\text{inv}(X)H = -X^{-1}H X^{-1}, H \in \mathbb{R}^{n,n}} \end{array} \right]$$

Product rule:  $F: D \subset V \mapsto W, G: D \subset V \mapsto U$  sufficiently smooth,  $b: W \times U \mapsto Z$  bilinear, ie., linear in each argument:

$$T(x) = b(F(x), G(x)) \Rightarrow DT(x)h = b(DF(x)h, G(x)) + b(F(x), DG(x)h), \quad (8.4.9)$$

$h \in V, x \in D$ .

applied with:  $F \stackrel{!}{=} \text{inv}, G \stackrel{!}{=} \text{id} (\Rightarrow D\text{G}(X)H = H)$

b  $\Leftrightarrow$  matrix multiply

$$V = U = W = Z = \mathbb{R}^{n,n}$$

! We never use that matrices can represent linear maps

Note:  $D\text{id} = \text{id} :$ 

$$\begin{array}{c} \text{id}(x+h) = x + h = \text{id}(x) + \text{id}(h) \\ \Updownarrow \qquad \qquad \qquad \Updownarrow \qquad \qquad \qquad \Updownarrow \\ F(x+h) \approx \quad \quad \quad F(x) + DF(x)(h) \end{array}$$

Newton iteration for matrix inversion

$$\left[ \begin{array}{l} F: \mathbb{R}_{\text{reg}}^{n,n} \rightarrow \mathbb{R}^{n,n}, F(X) = A - X^{-1} \\ [A \in \mathbb{R}_{\text{reg}}^{n,n}] \\ F(X^*) = 0 \iff X^* = A^{-1} \end{array} \right]$$

Derivative of  $F$ :  $F(X) = A - \text{inv}(X)$ 

$$DF(X) \cdot H = 0 + X^{-1}H X^{-1}$$

Newton correction:

$$DF(X^{(k)})S = (X^{(k)})^{-1}S(X^{(k)})^{-1} = -F(X^{(k)}) = -A + (X^{(k)})^{-1}$$

$$\Rightarrow S = -X^{(k)}AX^{(k)} + X^{(k)}$$

$$\Rightarrow \text{Newton it. : } \begin{aligned} X^{(k+1)} &= X^{(k)} - X^{(k)}AX^{(k)} + X^{(k)} \\ &= X^{(k)}(2I - A \cdot X^{(k)}) \end{aligned}$$

### 2.4.3. Termination of Newton Iteration

Asymptotic quadratic cog:

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

$$\Rightarrow \|\underline{x}^{(k)} - \underline{x}^*\| \approx \|\underline{x}^{(k)} - \underline{x}^{(k+1)}\|$$

**correction based termination**

↳ One redundant Newton step ↴

often very few steps only

significant additional cost.

Cheaper: use **simplified Newton correction**  
(for  $n \geq 1$ )

$$\Delta \underline{x}^{(k)} := D F(\underline{x}^{(k-1)})^{-1} F(\underline{x}^{(k)})$$

↑  
LU-dec. available

► Economical correction based termination criterion for Newton's method:

**STOP**, as soon as  $\|\Delta \bar{\underline{x}}^{(k)}\| \leq \tau_{\text{rel}} \|\underline{x}^{(k)}\|$  or  $\|\Delta \bar{\underline{x}}^{(k)}\| \leq \tau_{\text{abs}}$ ,

with **simplified Newton correction**  $\Delta \bar{\underline{x}}^{(k)} := D F(\underline{x}^{(k-1)})^{-1} F(\underline{x}^{(k)})$ .

### 2.4.4. Damped Newton Method

Problem: Annoyingly local convergence

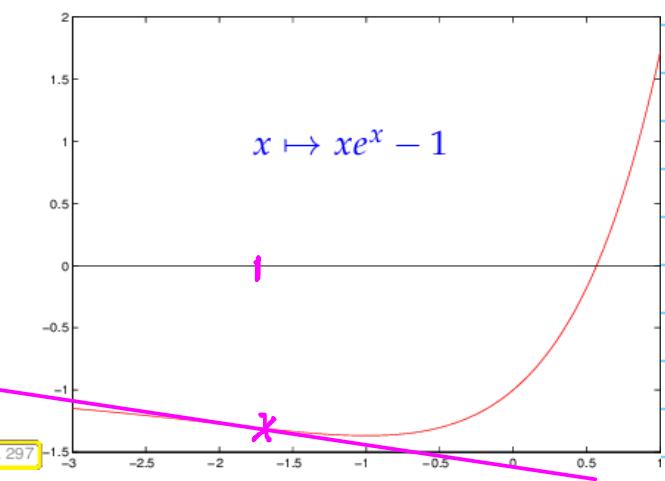
Causes:

①

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

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

because all Newton corrections for  $x^{(k)} < -1$  make the iterates decrease even further.



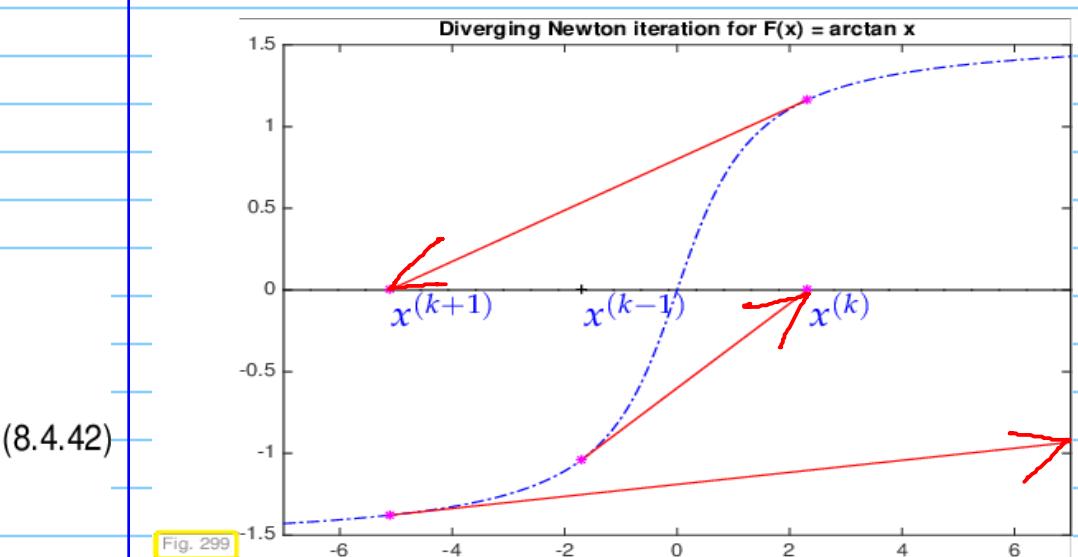
→ Newton correction in the wrong direction

→ Beyond repair

②

Overshooting of  
Newton correction

Remedy:  
**damping**



(19)



we observe "overshooting" of Newton correction

Idea:

damping of Newton correction:

$$\text{With } \lambda^{(k)} > 0: \quad \mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \lambda^{(k)} \mathbf{D}F(\mathbf{x}^{(k)})^{-1} \mathbf{F}(\mathbf{x}^{(k)}). \quad (8.4.47)$$

Terminology:  $\lambda^{(k)}$  = damping factor ,  $\lambda^{(k)} \in [0, 1]$

### Affine invariant damping strategy

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

$$\text{choose "maximal" } 0 < \lambda^{(k)} \leq 1: \quad \|\Delta\bar{\mathbf{x}}(\lambda^{(k)})\| \leq \left(1 - \frac{\lambda^{(k)}}{2}\right) \|\Delta\mathbf{x}^{(k)}\|_2 \quad (8.4.49)$$

where  $\Delta\mathbf{x}^{(k)} := \mathbf{D}F(\mathbf{x}^{(k)})^{-1} \mathbf{F}(\mathbf{x}^{(k)})$  → current Newton correction ,

$\Delta\bar{\mathbf{x}}(\lambda^{(k)}) := \mathbf{D}F(\mathbf{x}^{(k)})^{-1} \mathbf{F}(\mathbf{x}^{(k)} + \lambda^{(k)} \Delta\mathbf{x}^{(k)})$  → tentative simplified Newton correction .

↳ tentative next iterate

\* The same damping factor for  $\mathbf{A}\mathbf{F}(\mathbf{x}) = \mathbf{0} \quad \forall \mathbf{A} \in \mathbb{R}^{n,n}$

"Theory": If quad. cong  $\Rightarrow \|\Delta\bar{\mathbf{x}}(1)\| \leq \frac{1}{2} \|\Delta\mathbf{x}^{(k)}\|$

If NMT passed  $\Rightarrow$  reduce damping :  $\lambda \leftarrow 2\lambda$   
failed  $\Rightarrow$  increase damping :  $\lambda \leftarrow \lambda/2$

### C++11 code 8.4.50: Generic damped Newton method based on natural monotonicity test

```

1 template <typename FuncType, typename JacType, typename VecType>
2 void dampnewton(const FuncType &F, const JacType &DF,
3                  VecType &x, double rtol, double atol)
4 {
5     using index_t = typename VecType::Index;
6     using scalar_t = typename VecType::Scalar;
7     const index_t n = x.size();
8     const scalar_t lmin = 1E-3; // Minimal damping factor
9     scalar_t lambda = 1.0; // Initial and actual damping factor
10    VecType s(n), st(n); // Newton corrections
11    VecType xn(n); // Tentative new iterate
12    scalar_t sn, stn; // Norms of Newton corrections
13
14    do {
15        auto jacfac = DF(x).lu(); // LU-factorize Jacobian
16        s = jacfac.solve(F(x)); // Newton correction
17        sn = s.norm(); // Norm of Newton correction
18        lambda *= 2.0; // tentatively reduce damping
19        do {
20            lambda /= 2; // tentatively increase damping
21            if (lambda < lmin) throw "No convergence: lambda > 0";
22            xn = x - lambda * s; // Tentative next iterate
23            st = jacfac.solve(F(xn)); // Simplified Newton correction
24            stn = st.norm();
25        }
26        while (stn > (1 - lambda / 2) * sn); // Natural monotonicity test
27        x = xn; // Now: xn accepted as new iterate
28        lambda = std::min(2.0 * lambda, 1.0); // Try to mitigate damping
29    }
30    // Termination based on simplified Newton correction
31    while ((stn > rtol * x.norm()) && (stn > atol));
32 }
```