

# 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 2015

(C) Seminar für Angewandte Mathematik, ETH Zürich

URL: <http://www.sam.math.ethz.ch/~hiptmair/tmp/NumCSE/NumCSE15.pdf>

## II. Iterative Methods for Non-linear Systems of Equations

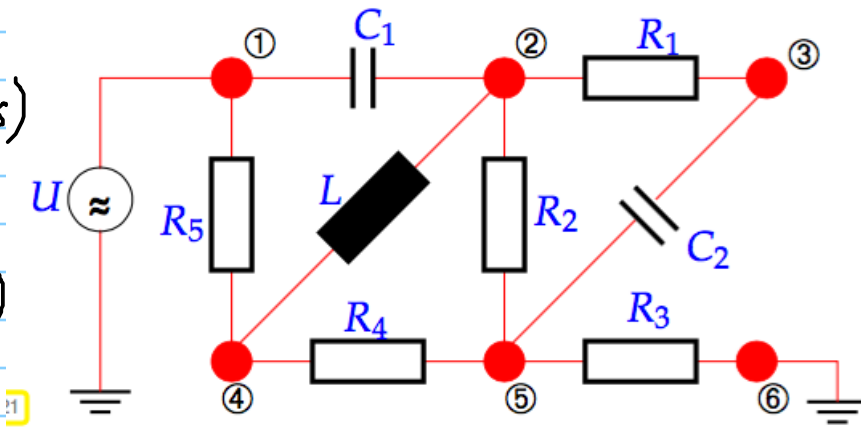
### Case study: Nodal analysis of electrical circuits

- $\hat{=}$  nodes

Circuit elements (connect nodes)  
 $\hookrightarrow$  all linear\*

$\square$   $\hat{=}$  resistor (resistance  $R$ )  
 $\blacksquare$   $\hat{=}$  coil (inductance  $L$ )  
 $\text{--}\text{||}\text{--}$   $\hat{=}$  capacitor (capacity  $C$ )

$\text{---}$   $\hat{=}$  wire



$I_{kj} \hat{=}$  current node  $k \rightarrow$  node  $j$ :  $I_{kj} = -I_{jk}$

Kirchhoffs law:  $\sum_{j \in S(k)} I_{kj} = 0 \quad (1)$

$U_k \hat{=}$  nodal potential

$$* \quad I_{kj} = \alpha (U_k - U_j) \quad (2)$$

$\uparrow \quad \quad \quad \hookrightarrow R^{-1} \text{ or } i\omega C \text{ or } -i/\omega L$

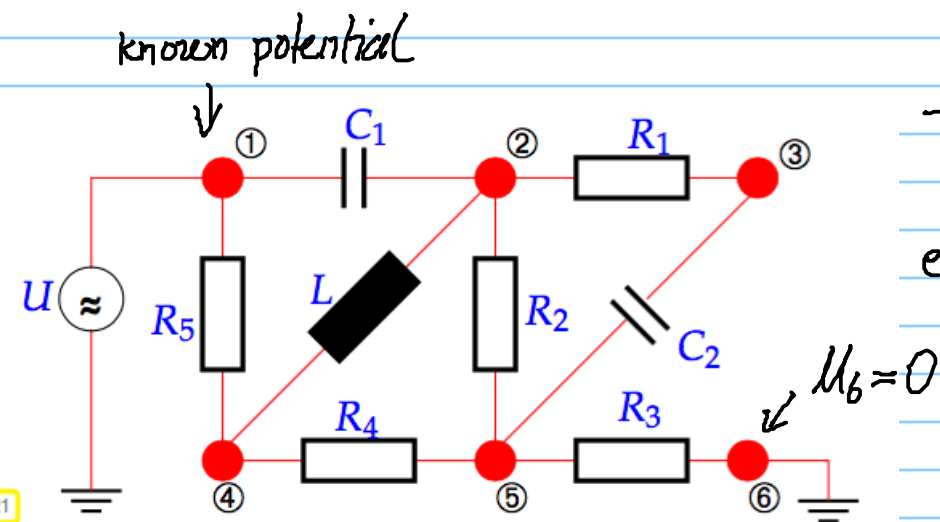
$\in \mathbb{C}$ : complex amplitudes  $\quad \quad \quad [\omega \hat{=}$  angular frequ.]

$$i(t) = \operatorname{Re}(I e^{i\omega t}) = \operatorname{Re}(I) \cos(\omega t) - \operatorname{Im}(I) \sin(\omega t)$$

Nodal analysis:

- (1) for each node (with unknown potential)
- Replace all currents using (2)

2



→ 4 equations from (1)

e.g.  $I_{12} + I_{32} + I_{52} + I_{42} = 0$

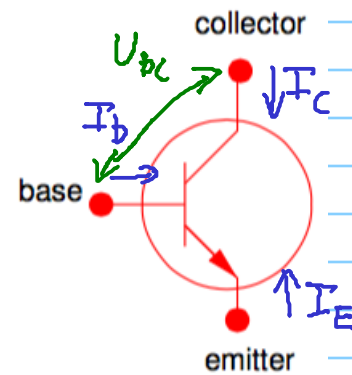
$$\begin{aligned} \textcircled{2}: & \quad i\omega C_1(U_2 - U_1) + R_1^{-1}(U_2 - U_3) - i\omega^{-1}L^{-1}(U_2 - U_4) + R_2^{-1}(U_2 - U_5) = 0, \\ \textcircled{3}: & \quad R_1^{-1}(U_3 - U_2) + i\omega C_2(U_3 - U_5) = 0, \\ \textcircled{4}: & \quad R_5^{-1}(U_4 - U_1) - i\omega^{-1}L^{-1}(U_4 - U_2) + R_4^{-1}(U_4 - U_5) = 0, \\ \textcircled{5}: & \quad R_2^{-1}(U_5 - U_2) + i\omega C_2(U_5 - U_3) + R_4^{-1}(U_5 - U_4) + R_3^{-1}(U_5 - U_6) = 0, \\ & \quad U_1 = U, \quad U_6 = 0. \end{aligned}$$

$$\begin{pmatrix} i\omega C_1 + \frac{1}{R_1} - \frac{i}{\omega L} + \frac{1}{R_2} & -\frac{1}{R_1} & \frac{i}{\omega L} & -\frac{1}{R_2} \\ -\frac{1}{R_1} & \frac{1}{R_1} + i\omega C_2 & 0 & -i\omega C_2 \\ \frac{i}{\omega L} & 0 & \frac{1}{R_5} - \frac{i}{\omega L} + \frac{1}{R_4} & -\frac{1}{R_4} \\ -\frac{1}{R_2} & -i\omega C_2 & -\frac{1}{R_4} & \frac{1}{R_2} + i\omega C_2 + \frac{1}{R_4} + R_3^{-1} \end{pmatrix} \begin{pmatrix} U_2 \\ U_3 \\ U_4 \\ U_5 \end{pmatrix} = \begin{pmatrix} i\omega C_1 U \\ 0 \\ \frac{1}{R_5} U \\ 0 \end{pmatrix}$$

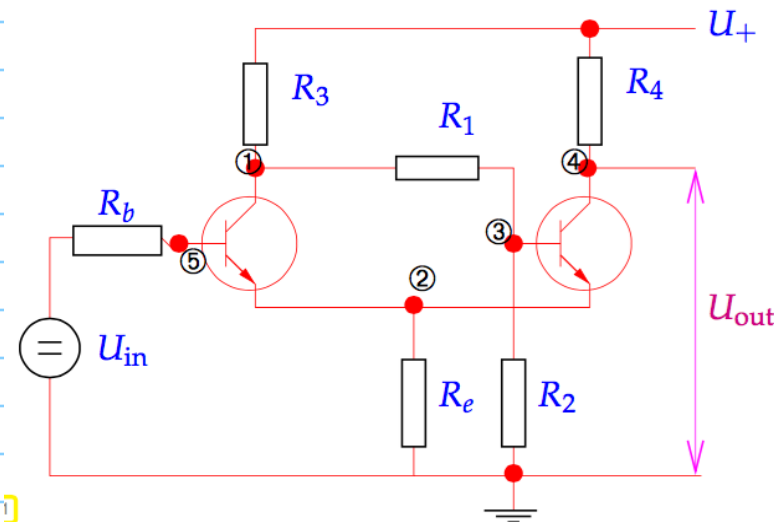
$\hat{=}$  LSE (in  $\mathbb{C}$ )

For large circuits : sparse LSE

Schmitt trigger  $\triangleright$   
NPN transistors:



3-port-circuit element



Ebers-Moll model : *strongly nonlinear*

$$\begin{aligned} I_C &= I_S \left( e^{\frac{U_{BE}}{U_T}} - e^{\frac{U_{BC}}{U_T}} \right) - \frac{I_S}{\beta_R} \left( e^{\frac{U_{BC}}{U_T}} - 1 \right) = I_C(U_{BE}, U_{BC}), \\ I_B &= \frac{I_S}{\beta_F} \left( e^{\frac{U_{BE}}{U_T}} - 1 \right) + \frac{I_S}{\beta_R} \left( e^{\frac{U_{BC}}{U_T}} - 1 \right) = I_B(U_{BE}, U_{BC}), \\ I_E &= I_S \left( e^{\frac{U_{BE}}{U_T}} - e^{\frac{U_{BC}}{U_T}} \right) + \frac{I_S}{\beta_F} \left( e^{\frac{U_{BE}}{U_T}} - 1 \right) = I_E(U_{BE}, U_{BC}). \end{aligned}$$

$$\begin{aligned} \textcircled{1}: & \quad R_3(U_1 - U_+) + R_1(U_1 - U_3) + I_B(U_5 - U_1, U_5 - U_2) = 0, \\ \textcircled{2}: & \quad R_e U_2 + I_E(U_5 - U_1, U_5 - U_2) + I_E(U_3 - U_4, U_3 - U_2) = 0, \\ \textcircled{3}: & \quad R_1(U_3 - U_1) + I_B(U_3 - U_4, U_3 - U_2) = 0, \\ \textcircled{4}: & \quad R_4(U_4 - U_+) + I_C(U_3 - U_4, U_3 - U_2) = 0, \\ \textcircled{5}: & \quad R_b(U_5 - U_{in}) + I_B(U_5 - U_1, U_5 - U_2) = 0. \end{aligned}$$

5 equations  $\leftrightarrow$  5 unknowns  $U_1, U_2, U_3, U_4, U_5$

In short :  $F(x) = 0$  ,  $x = (U_1, \dots, U_5)^T$

③

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

NLSE: Given  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$\text{seek } \underline{x} \in \mathbb{R}^n: F(\underline{x}) = 0$$

- no general theory
- $F \stackrel{\triangle}{=} \text{function } y = F(x) \text{ "black box"}$

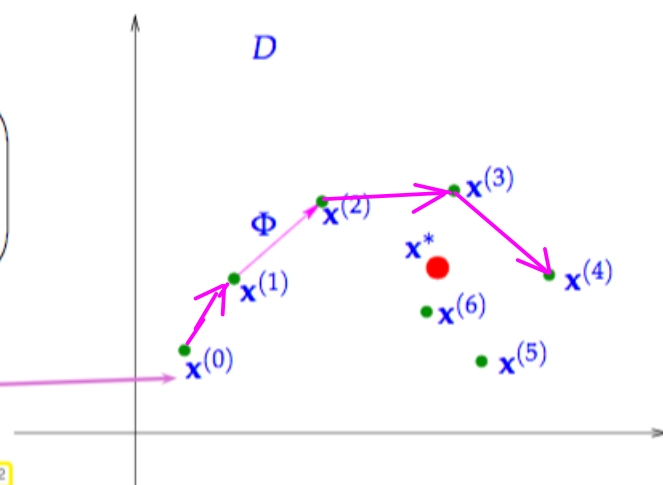
## 2.1. Iterative Methods

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

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

Initial guess

Fig. 62



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

More concrete: **stationary m-point iteration**

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

↳ iteration function

Initial guess:  $\underline{x}^{(0)}, \dots, \underline{x}^{(m-1)} \in \mathbb{R}^n$

Issues: → well defined?  
→ does it converge?  
→ If yes, does it converge to a solution?  
→ How fast (**speed of convergence**)?

[ Initial guess matters much! ]

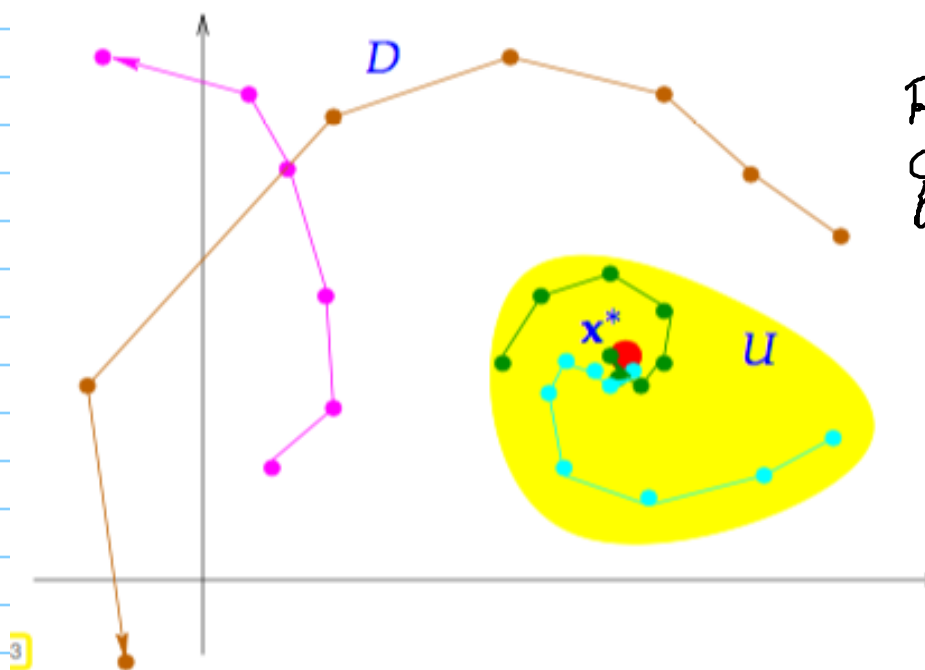
**Definition 2.1.8. Local and global convergence** → [12, Def. 17.1]

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**.



Region  $U$  of local convergence may be very small.

## ④ 2.1.1. Speed of convergence

"Slow methods":

### Definition 2.1.9. Linear convergence

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

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

smallest possible  $L$  in  $(*)$ : rate of (linear) conv.

How to tell linear conv. in numerical test ( $\mathbf{x}^*$  known)

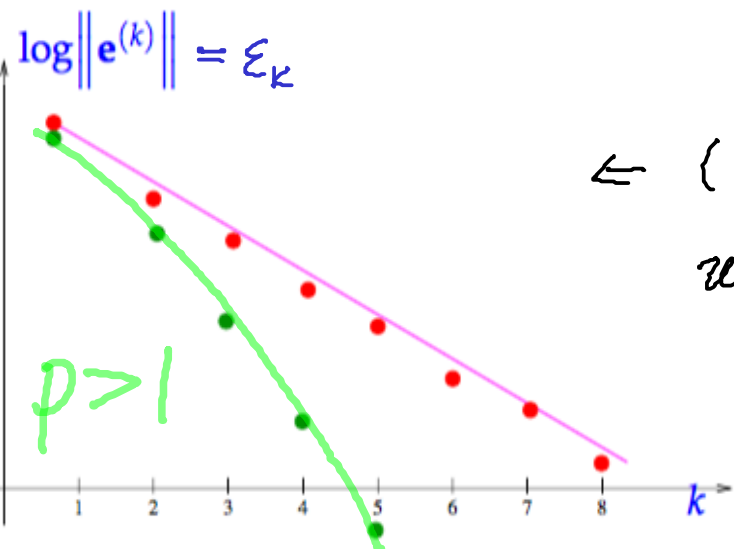
$$\varepsilon_k := \|\mathbf{x}^{(k)} - \mathbf{x}^*\| \text{ known,}$$

$$L.C. \Rightarrow \text{assume } \varepsilon_{k+1} \approx L \varepsilon_k \approx L^2 \varepsilon_{k-1} \dots$$

$$\Rightarrow \varepsilon_k \approx L^k \varepsilon_0$$

$$\log \varepsilon_k \approx k \underbrace{\log L}_{\leq 0} + \log \varepsilon_0$$

$\leftarrow (k, \log \varepsilon_k)$  on a straight line with slope  $\log L < 0$ !



From tabulated values: check, if  $\frac{\varepsilon_k}{\varepsilon_{k-1}} \approx L$   
"Faster convergence"

### Definition 2.1.17. Order of convergence $\rightarrow$ [12, Sect. 17.2], [4, Def. 5.14], [16, Def. 6.1]

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

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

and, in addition,  $C < 1$  in the case  $p = 1$  (linear convergence  $\rightarrow$  Def. 2.1.9).

Identifying conv. of order  $p$  from measured  $\varepsilon_k$

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

$$\xrightarrow{\text{subtract}} \log \varepsilon_{k+1} \approx \log C + p \log \varepsilon_k$$

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

$$p = \frac{\log \varepsilon_{k+1} - \log \varepsilon_k}{\log \varepsilon_k - \log \varepsilon_{k-1}}$$

Famous example:  $\sqrt{\quad}$ -iteration ( $n = 1, m = 1$ )

$$x^{(k+1)} = \frac{1}{2} \left( x^{(k)} + \frac{a}{x^{(k)}} \right), \quad a > 0: x^{(k)} \rightarrow \sqrt{a}$$

$$\underbrace{x^{(k+1)} - \sqrt{a}}_{\varepsilon^{(k+1)}} = \frac{1}{2x^{(k)}} \underbrace{(x^{(k)} - \sqrt{a})^2}_{\varepsilon^{(k)}}, \quad x^{(k)} > \sqrt{a} \text{ for } k \geq 1$$

$$\leq \frac{1}{2\sqrt{a}} \varepsilon^{(k)}$$

⑤  $\Rightarrow$  quadratic convergence  $p=2$

$k$	$x^{(k)}$	$e^{(k)} := x^{(k)} - \sqrt{2}$	$\log \frac{ e^{(k)} }{ e^{(k-1)} } : \log \frac{ e^{(k-1)} }{ e^{(k-2)} }$
0	2.0000000000000000	0.58578643762690485	
1	1.5000000000000000	0.08578643762690485	
2	1.4166666666666665	0.00245310429357137	1.850
3	1.4142156862745096	0.00000212390141452	1.984
4	1.4142135623746898	0.0000000000159472	2.000
5	1.4142135623730949	0.0000000000000022	0.630

[red  $\hat{=}$  correct digits]

roundoff

relative error

$$x^{(k)} = x^* (1 + \delta_k) : \delta_k \leq 10^{-l} \Leftrightarrow x^{(k)} \text{ has } l \text{ correct digits}$$

doubling of no. of correct digits in every step!

Quadratic cvg.

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

$$|\delta_{k+1} x^*| \approx C |\delta_k x^*|^2, x^* \neq 0$$

$$\Rightarrow |\delta_{k+1}| \approx C |x^*| |\delta_k|^2$$

$$C, |x^*| \approx 1, \delta_k = 10^{-l} \Rightarrow \delta_{k+1} \approx 10^{-2l}$$

Compute linear cvg.: fixed number (also a fraction) of extra correct digits in each step

## 2.1.2. Termination

Ideal:  $\hookrightarrow$  STOP, if  $x^{(k)}$  is "just good enough"

$$\|x^{(k)} - x^*\| \leq \tau_{\text{abs}}$$

$\uparrow$

Absolute tolerance

$$\text{or } \|x^{(k)} - x^*\| \leq \tau_{\text{rel}} \cdot \|x^*\|$$

$\uparrow$

relative tolerance

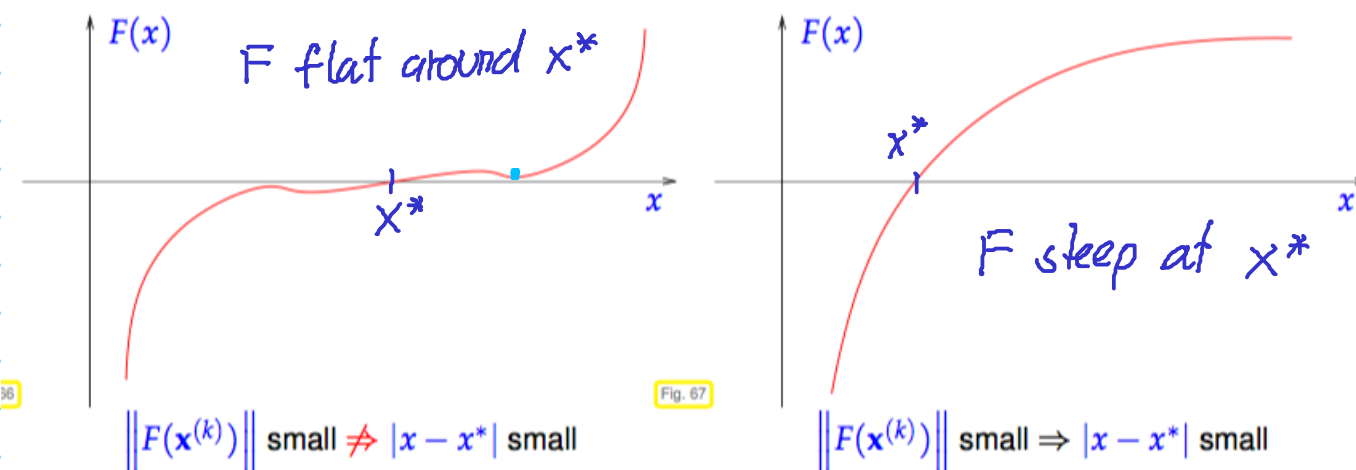
Not practical, because  $x^*$  is not known

Practical: [for solving  $F(x) = 0$ ]

① Residual based termination:

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

tells little about  $\varepsilon_k$





⑥ [8.10.2015]

Rep: Solve  $F(\underline{x}) = 0$ ,  $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ Iterative method:  $\{x^{(k-m+1)}, \dots, x^{(k)}\} \rightarrow x^{(k+1)}$ Cvg of order  $p > 1$ :  $\|e^{(k+1)}\| \leq C \|e^{(k)}\|^p$  (\*)If (\*), cvg. guaranteed, if  $C \|e^{(0)}\|^{p-1} < 1 \rightarrow$  not practicalHint:  $\|e^{(k+1)}\| \leq C \|e^{(k)}\|^{p-1} \|e^{(k)}\|$ 

② Convergence based termination:

STOP, if  $\|x^{(k+1)} - x^{(k)}\| \leq \begin{cases} \tau_{abs} \\ \tau_{rel} \cdot \|x^{(k+1)}\| \end{cases}$  $\rightarrow$  Generically, no guaranteesException: Linearly cvg. iteration with known rate  $L < 1$ 

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

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

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

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

 $\hookrightarrow$  **Reliable** upper bound: can replace  $\|x^{(k)} - x^*\|$  in termination criteriaIf we overestimate  $L \Rightarrow$  still reliable termination

$$\begin{aligned} \text{Remark: } \Rightarrow \|x^{(k)} - x^*\| &\leq L^{k-1} \|x^{(1)} - x^*\| \\ &\leq \frac{L^k}{1-L} \|x^{(1)} - x^{(0)}\| \end{aligned}$$

 $\rightarrow$  Can be used for **a priori** termination

#steps fixed before start of iteration.

```
function x = sqrtit(a)
x_old = -1; x = a;
while (x_old ~= x)
    x_old = x;
    x = 0.5*(x+a/x);
end
```

 $\leftarrow$  "M-based termination" $\Rightarrow$  guarantees relative update  $\leq \text{EPS}$ 

2.2. Fixed point iterations = 1-point methods

 $x^{(k+1)} = \Phi_F(x^{(k)})$  with iteration function  $\Phi_F: M \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ 

$$\left. \begin{aligned} x^* &= \lim_{k \rightarrow \infty} x^{(k)} \\ \Phi &\text{ continuous} \end{aligned} \right\} \Rightarrow \underline{x}^* = \Phi_F(x^*)$$

$\uparrow$  fixed point of  $\Phi_F$

FPI is **consistent**:  $\Phi_F(x) = x \Leftrightarrow F(x) = 0$

⑦ Many  $\Phi$  possible!

Ex:

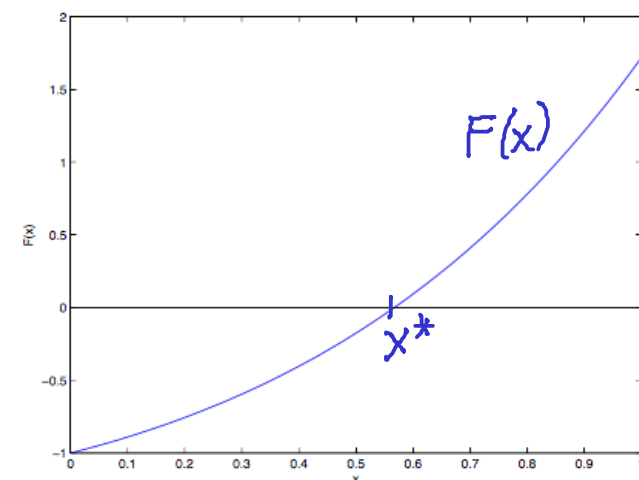
$$F(x) = xe^x - 1, \quad x \in [0, 1].$$

Different fixed point forms:

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

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

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



$$x^{(0)} = 0.5$$

$\Phi_1$

$\Phi_2$

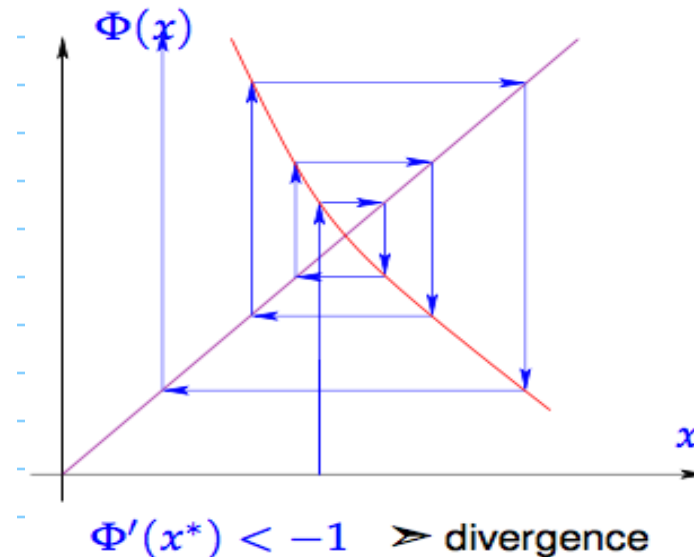
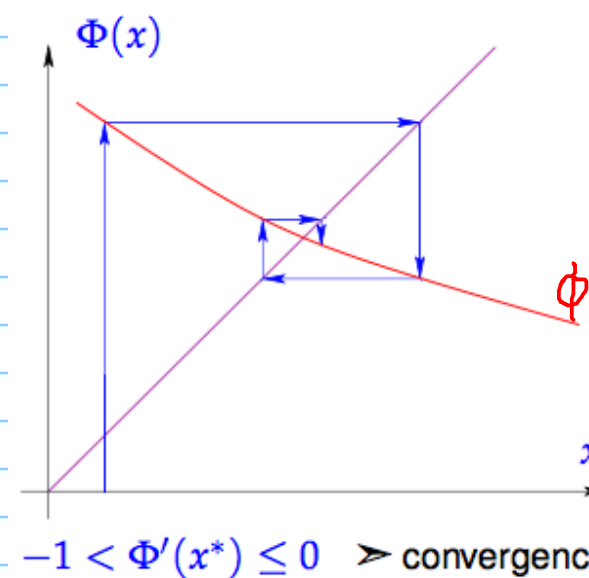
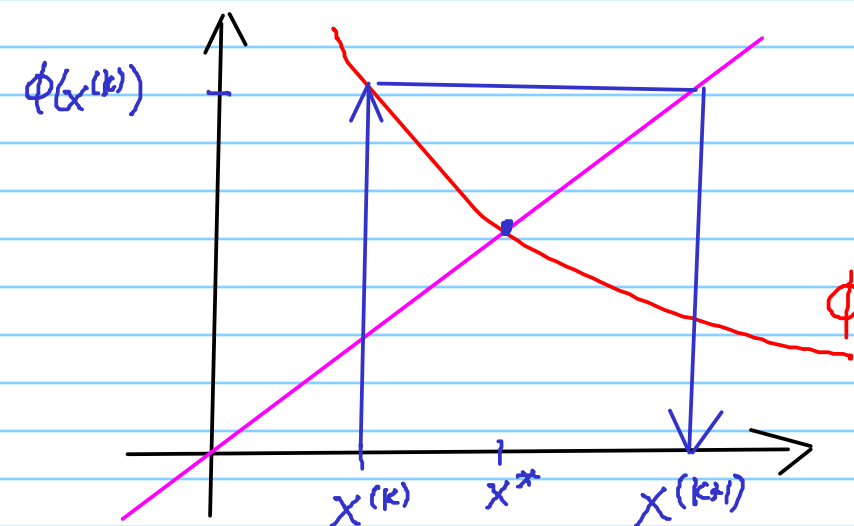
$\Phi_3$

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

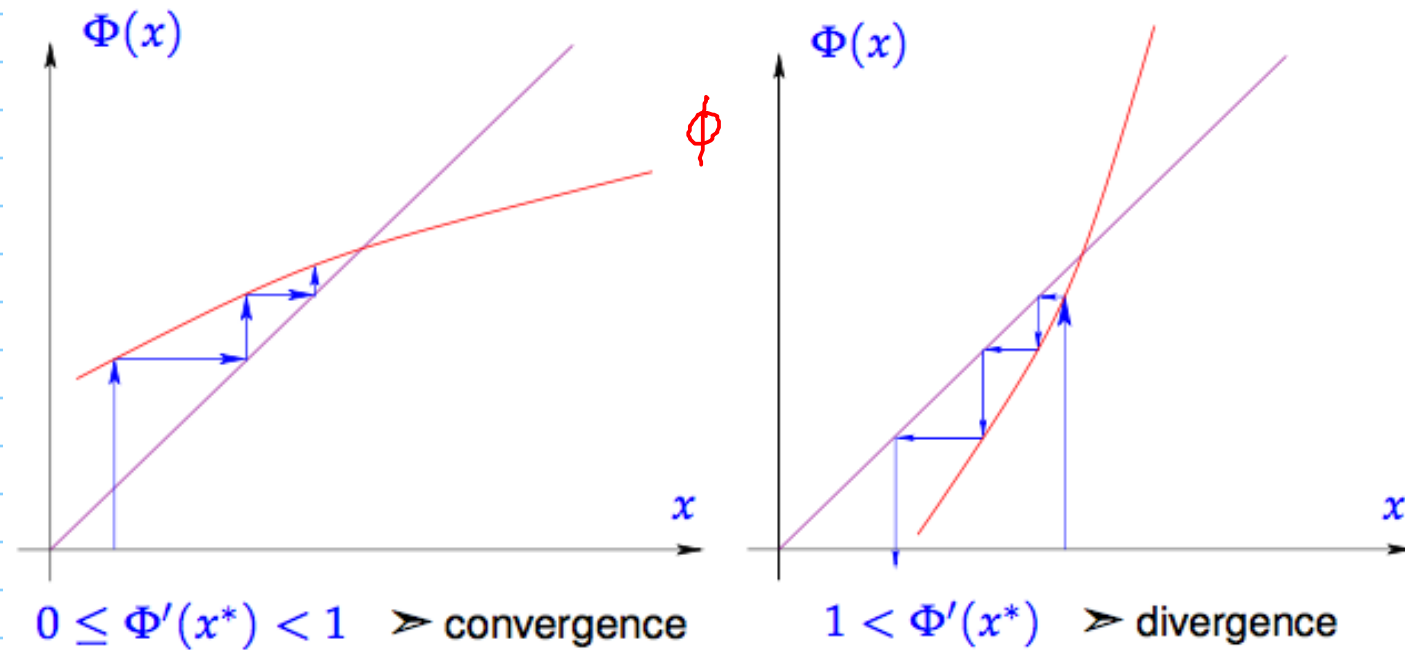
↓  
linear conv.      quadratic conv.      no conv.

How to predict this?

Visualization of FPI in 1D:



8



$\Rightarrow$  Slope of  $x \rightarrow \phi(x)$  at fixed point crucial!

Analysis ( $\phi$  smooth): Taylor expansion around  $x^*$

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

$$= \phi'(x^*)(x^{(k)} - x^*) + \frac{1}{2}\phi''(x^*)(x^{(k)} - x^*)^2 + \frac{1}{6}\phi'''(x^*)(x^{(k)} - x^*)^3 + \dots$$

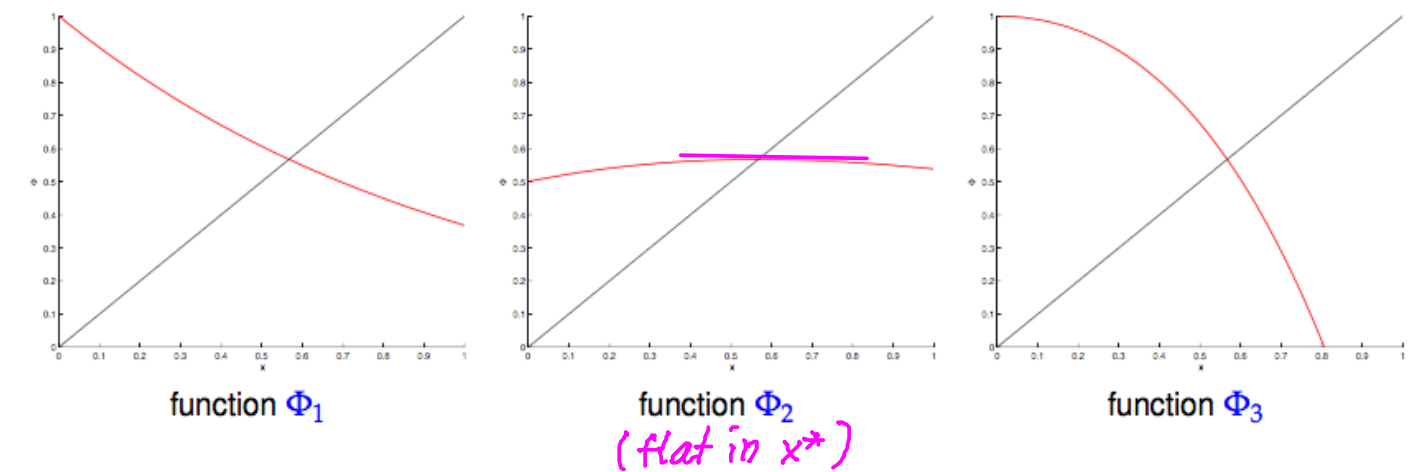
neglect

If  $|x^{(k)} - x^*| \ll 1$  ↑ ↑

If  $|\phi'(x^*)| < 1 \Rightarrow$  local linear conv, rate  $\approx \phi'(x^*)$

If  $\phi'(x^*) = 0 \Rightarrow$  local quadratic conv.

In Example:



Remains true for  $n > 1$ :

**Lemma 2.2.10. Sufficient condition for local linear convergence of fixed point iteration**  $\rightarrow$   
[12, Thm. 17.2], [4, Cor. 5.12]

If  $\Phi : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $\Phi(x^*) = x^*$ ,  $\Phi$  differentiable in  $x^*$ , and  $\|D\Phi(x^*)\| < 1$ , then the fixed point iteration (2.2.2) converges locally and at least linearly.

matrix norm, Def. 1.5.68!

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

If  $D\Phi(x^*) = 0 \Rightarrow$  local quadratic conv.

**Lemma 2.2.12. Sufficient condition for linear convergence of fixed point iteration**

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

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

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



9

## 2.3. $n=1$ : Zero finding

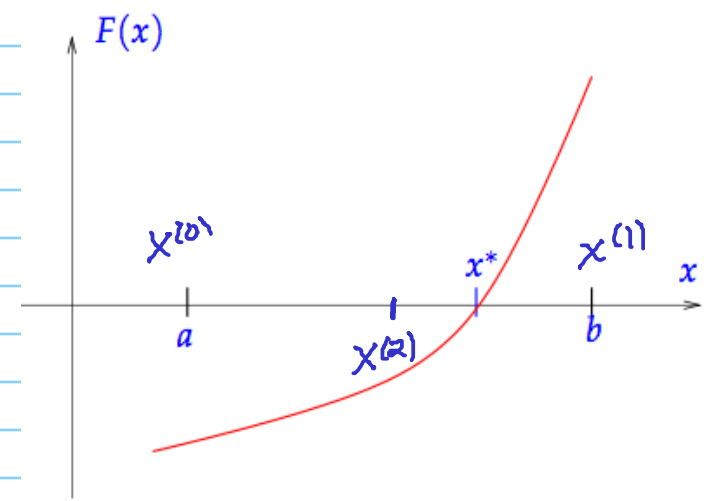
Seek  $x^* \in \mathbb{R} : F(x^*) = 0$ ,  $F: I \subset \mathbb{R} \rightarrow \mathbb{R}$   
continuous

### 2.3.1. Bisection

Interval  $[a, b]$ :

$$\left. \begin{matrix} F(a) < 0 \\ F(b) > 0 \end{matrix} \right\} \Rightarrow \exists x^* \in [a, b] : F(x^*) = 0$$

[intermediate value theorem]



**MATLAB-code 2.3.2: Bisection method for solving  $F(x) = 0$  on  $[a, b]$**

```

1 function x = bisection(F,a,b,tol)
2 % Searching zero of F in [a,b] by bisection
3 if (a>b), t=a; a=b; b=t; end;
4 fa = F(a); fb = F(b);
5 if (fa*fb>0), error('f(a), f(b) same sign'); end;
6 if (fa > 0), v=-1; else v = 1; end
7 x = 0.5*(b+a);
8 while ((b-a > tol) && ((a<x) & (x<b)))
9     if (v*F(x)>0), b=x; else a=x; end;
10    x = 0.5*(a+b);
11 end
    
```

M-based termination

$$\Rightarrow |x^* - x^{(k)}| \leq 2^{-k} |b-a| \quad \text{"kind of lin. conv."}$$

$\rightarrow$  robust method:  $|x^{(k)} - x^*| \leq \text{tol}$  is guaranteed  
 simple method, only F-evaluation required

### 2.3.2. Model Function methods should be simple

- $\hookrightarrow$  In step  $k$ :
- replace  $F$  with  $\tilde{F}_k$
  - $x^{(k+1)} : \tilde{F}_k(x^{(k+1)}) = 0$

#### 2.3.2.1 Newton's method

Assume:  $F$  differentiable

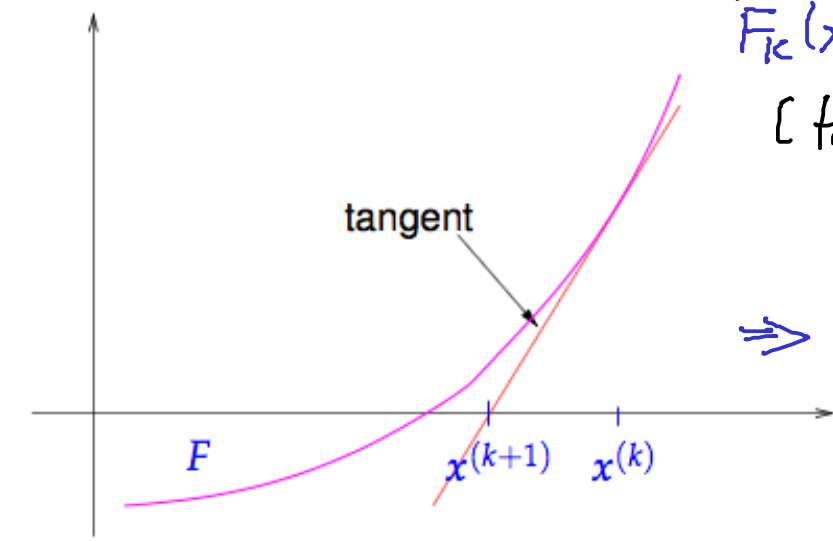
$$\tilde{F}_k(x) = F(x^{(k)}) + F'(x^{(k)})(x - x^{(k)})$$

[tangent in  $(x^{(k)}, F(x^{(k)}))$ ]

$$\tilde{F}_k(x^{(k+1)}) \stackrel{!}{=} 0$$

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

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



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

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

$\{ F(x^*) = 0 \Leftrightarrow \phi(x^*) = x^* \} \Rightarrow \phi'(x^*) = 0 \Rightarrow$  local quad. conv.

10

Example:  $F(x) = x^2 - a \Rightarrow F'(x) = 2x$

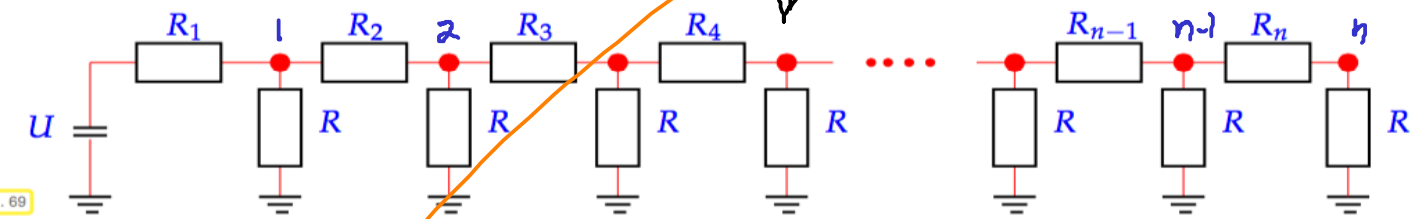
N.I.:  $x^{(k+1)} = x^{(k)} - \frac{(x^{(k)})^2 - a}{2x^{(k)}} = \frac{1}{2} \left( x^{(k)} + \frac{a}{x^{(k)}} \right)$

Drawback:   
 • Local cvg.   
 •  $F'(x)$  not available

Case study: Finding  $F'$

want to achieve target potential here by varying  $R$

Linear circuit:



Nodal analysis  $\Rightarrow$  LSE:  $(A + xI) \underline{u}(x) = \underline{b}$  (\*)   
 $[x = 1/R]$    
 $\uparrow$   $\uparrow$    
 symmetric tridiagonal matrix vector of nodal pot.

$\Rightarrow F(x) = e_k^T \underline{u}(x) - 1 = 0$

$\Rightarrow F'(x) = e_k^T \underline{u}'(x)$

$$A = \begin{bmatrix} \frac{1}{R_1} + \frac{1}{R_2} & -\frac{1}{R_2} & & \\ -\frac{1}{R_2} & \frac{1}{R_2} + \frac{1}{R_3} & -\frac{1}{R_3} & \\ & \ddots & \ddots & \ddots \\ & & & \frac{1}{R_{n-1}} + \frac{1}{R_n} \end{bmatrix} \in \mathbb{R}^{n,n}$$

$\triangleright (A + xI) \underline{u}(x) = \underline{b}$  can be solved with effort  $O(n)$

$\underline{u}'(x)$  by implicit differentiation

dx on (\*) & product rule:

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

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

Newton iteration:

(i) Solve  $(A + x^{(k)}I) \underline{u} = \underline{b}$

(ii) Solve  $(A + x^{(k)}I) \underline{u}' = -\underline{u}$

(iii)  $x^{(k+1)} = x^{(k)} - \frac{e_k^T \underline{u} - 1}{e_k^T \underline{u}'}$

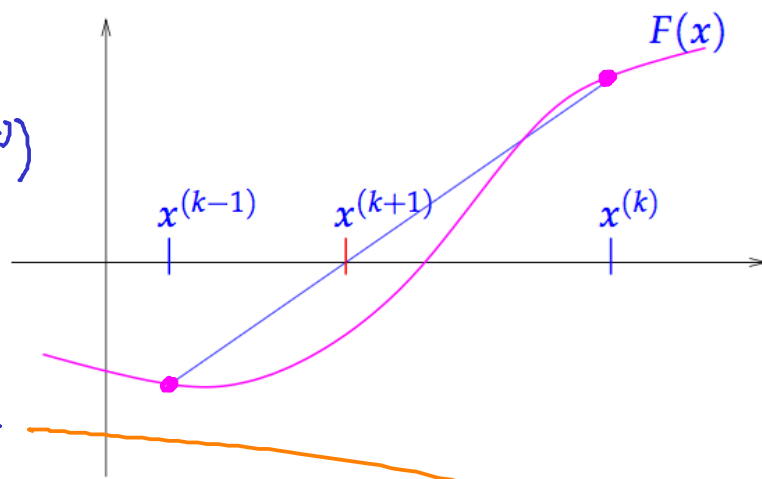
## ⑪ 2.3.2.2. Multi-point method ( $m > 1$ )

Simplest: secant method

$$\tilde{F}_k(x) = F(x^{(k)}) + \Delta s(x - x^{(k)})$$

$$\Delta s = \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$$

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



```
function x = secant(x0,x1,F,rtol,atol)
fo = F(x0);
for i=1:MAXIT
    fn = F(x1);
    s = fn*(x1-x0)/(fn-fo); % correction
    x0 = x1; x1 = x1-s;
    if ((abs(s) <
        max(atol,rtol*min(abs([x0;x1])))))
        x = x1; return; end
    fo = fn;
end
```

Correction based termination

1 F-evaluation  
per step

No derivatives

↓  
black-box suitable

Exp.: Cvg. of 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	0.00000000134427	
11	0.56714329040978	-0.00000000000001	-0.00000000000000	

Fractional (!) order of cvg.  
( $p \approx 1.6$ )

$$S.M.: \quad x^{(k+1)} = \phi(x^{(k)}, x^{(k-1)})$$

$$\phi(x, y) = x - \frac{x-y}{F(x)-F(y)} \cdot F(x)$$

$$F(x^*) = 0: \quad \phi(x^*, x^*) = x^*$$

$$\frac{\partial \phi}{\partial x}(x^*, x^*) = \frac{\partial \phi}{\partial y}(x^*, x^*) = \frac{\partial^2 \phi}{\partial x^2}(x^*, x^*) = \frac{\partial^2 \phi}{\partial y^2}(x^*, x^*) = 0$$

▷ Asymptotic error recursion by Taylor exp.

$$e^{(k+1)} = \phi(x^* + e^{(k)}, x^* + e^{(k-1)}) - x^*$$

$$= K e^{(k)} e^{(k-1)} + \text{"small term"}$$

(12)

$$\left. \begin{aligned} |e^{(k)}| &\approx C|e^{(k-1)}|^p \\ |e^{(k+1)}| &\approx C^{p+1}|e^{(k-1)}|^{p^2} \end{aligned} \right\} \text{ for order } p$$

Plug into error recursion:

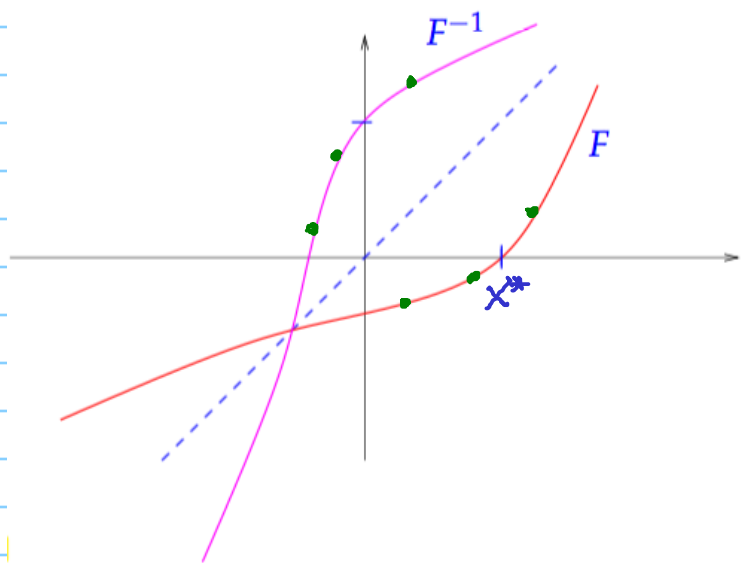
$$|e^{(k+1)}|^{p^2-p-1} \approx K^p C = \text{const.} \quad \forall k$$

$$\Rightarrow p^2 - p - 1 = 0 \Rightarrow p = \frac{1}{2}(1 + \sqrt{5}) \approx 1.6$$

More general: Inverse interpolation method

Assume:  $F: I \subset \mathbb{R} \rightarrow \mathbb{R}$  strictly monotonic

$$\text{Idea: } F(x^*) = 0 \iff F^{-1}(0) = x^*$$



Interpolate  $F^{-1}$  by a polynomial  $p$  of degree  $m-1$  in  $m$  points

$$p(F(x^{(k-j)})) = x^{(k-j)}$$

$$j = 0, \dots, m-1$$

$$x^{(k+1)} := p(0)$$

$m = 2 \Rightarrow$  secant method

$m = 3 \Rightarrow$  quadratic inverse interpolation

MAPLE code:  $p := x \rightarrow a \cdot x^2 + b \cdot x + c;$

$\text{solve}(\{p(f_0)=x_0, p(f_1)=x_1, p(f_2)=x_2\}, \{a, b, c\});$

$\text{assign}(\%); p(0);$

$\rightarrow$  parabola

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

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

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

$\uparrow$   
 $p \approx 1.8$

### ⑬ 2.3.3. Asymptotic efficiency

high accuracy  $\downarrow$   $\rightarrow = \frac{\text{gain}^* (\text{extra correct digits})}{\text{work}}$

work  $\leftrightarrow$  # F-eval. + # F'-eval. =  $W$  (per step)

$$* \|e^{(k)}\| \leq \rho \|e^{(0)}\|, \text{ digits gained } -\log_{10} \rho$$

$$\hookrightarrow \rho \ll 1$$

$$\text{Efficiency} = \frac{|\log \rho|}{W \cdot K(\rho)}$$

$K(\rho)$  = minimal no. of steps for \*

Assume iteration of order  $p > 1$ :

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

$$\|e^{(k)}\| \leq \dots C^{1+p+p^2+\dots+p^{k-1}} \|e^{(0)}\|^{p^k}$$

$$* \text{ is guaranteed, if } C^{\frac{p^k-1}{p-1}} \|e^{(0)}\|^{p^k-1} \leq \rho$$

$$[L_0 := C^{\frac{1}{p-1}} \|e^{(0)}\| < 1] \quad k \geq \frac{\log(\frac{\log \rho}{\log L_0} + 1)}{\log p}$$

$$\log p \cdot k \gtrsim \log |\log \rho| - \log |\log L_0|$$

$$\gtrsim \log |\log \rho|$$

$$\text{Asymp. Efficiency} \approx \frac{\log p}{\log |\log \rho|} \frac{\log \rho}{W}$$

$$\approx \frac{\log p}{W} \cdot \frac{\log \rho}{\log |\log \rho|}$$

$\uparrow$   
independent of  $p$   
and  $W$

( $p=2$ )

( $p \approx 1.6$ )

Newton vs. Secant method.

$$\frac{\log 2}{2W_N} : \frac{\log 1.6}{W_N} \approx 0.7$$

$\triangleright$  S.M. more efficient than Newton



## 14 2.4. Newton's method

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

As before: Model function method based on **local linearization**.

$$F(\underline{x}) \rightarrow \hat{F}_k(\underline{x}) = F(\underline{x}^{(k)}) + \underset{\substack{\uparrow \\ \text{Jacobian} \in \mathbb{R}^{n,n}}}{DF(\underline{x}^{(k)})}(\underline{x} - \underline{x}^{(k)})$$

Newton iteration:  $\underline{x}^{(k+1)} = \underline{x}^{(k)} - DF(\underline{x}^{(k)})^{-1} F(\underline{x}^{(k)})$

```
template <typename FuncType, typename JacType, typename VecType>
void newton(const FuncType &F, const JacType &DFinv,
            VecType &x, double rtol, double atol)
{
    using index_t = typename VecType::Index;
    using scalar_t = typename VecType::Scalar;
    const index_t n = x.size();
    VecType s(n);
    scalar_t sn;
    do {
        s = DFinv(x, F(x)); // compute Newton correction
        x -= s;              // compute next iterate
        sn = s.norm();
    }
    // correction based termination (relative and absolute)
    while ((sn > rtol*x.norm()) && (sn > atol));  $\rightarrow$  correction based t.c.
}
```

Objects of type **JacType** must provide a method

$\underline{z} = \text{VecType operator}(\text{const VecType \&x, const VecType \&f});$   
 $\rightarrow$  Solves  $DF(\underline{x})\underline{z} = \underline{f}$

$n \gg 1$ : Computation of Newton correction can be expensive, asymptotic effort  $O(n^3)$

Remark: **Affine invariance of Newton's method**

Newton it. for  $G_A(\underline{x}) = A \cdot F(\underline{x})$ ,  $A \in \mathbb{R}^{m,n}$  **regular**

$$\rightarrow \underline{x}^{(k+1)} = \underline{x}^{(k)} - \underbrace{[A \cdot DF(\underline{x}^{(k)})]^{-1}}_{\text{Jacobian of } G_A} A F(\underline{x}^{(k)})$$

$$[DF_{G_A}(\underline{x}) = A \cdot DF(\underline{x})] DF(\underline{x}^{(k)})^{-1} F(\underline{x}^{(k)})$$

The same N.I. for all  $A$ !

$\triangleright$  Termination criteria etc. for Newton's method should say STOP at the same index for all  $A$ !

⑮ Case study: Quasi-linear system of equations  
→ Homework Sheet #5

$$A(\underline{x})\underline{x} = \underline{b}$$

"LSE with solution dependent system matrix"

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

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

Sought: Newton's method for  $F(\underline{x}) = A(\underline{x})\underline{x} - \underline{b} \stackrel{!}{=} \underline{0}$

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

Derive:

$$\mathcal{D}\{\underline{x} \rightarrow T\underline{x}\} = T \cdot I = T \quad (\text{Jacobian of a linear mapping})$$

$$\mathcal{D}\{\underline{x} \rightarrow \underline{x}\|\underline{x}\|_2\} = \mathcal{D}\{\underline{x} \rightarrow [x_i \sqrt{x_1^2 + \dots + x_n^2}]_{i=1}^n\}$$

→ Jacobian by partial differentiation (safe & tedious)

Product rule:  $\mathcal{D}\{\underline{x} \rightarrow F(\underline{x}) \cdot G(\underline{x})\}h$

$$= \mathcal{D}F(\underline{x})h \cdot G(\underline{x}) + F(\underline{x}) \mathcal{D}G(\underline{x})h$$

← perturbation vector  $h \in \mathbb{R}^n$

General derivatives → Analysis

$$\phi(\underline{x} + h) = \phi(\underline{x}) + \mathcal{D}\phi(\underline{x})h + o(h)$$

↑ perturbation

↑ "tends → 0 faster than h"

Chain rule

$$\mathcal{D}\{\underline{x} \rightarrow F(G(\underline{x}))\}h = \mathcal{D}F(G(\underline{x}))\mathcal{D}G(\underline{x})h$$

Apply this to  $\underline{x} \rightarrow \underline{x} \cdot \|\underline{x}\|_2$

(i) Product rule:  $\mathcal{D}\{\underline{x} \rightarrow \underline{x}\|\underline{x}\|_2\}h = I \cdot h \cdot \|\underline{x}\|_2 + \underline{x} \mathcal{D}\{\underline{x} \rightarrow \|\underline{x}\|_2\}h$

(ii)  $\|\underline{x}\|_2 = \sqrt{\underline{x}^T \underline{x}}$  [Chain rule:  $F \leftrightarrow \sqrt{\quad}, G \leftrightarrow \underline{x}^T \underline{x}$ ]

$$\mathcal{D}\{\underline{x} \rightarrow \|\underline{x}\|_2\}h = \frac{1}{2\|\underline{x}\|_2} \cdot (\underbrace{h^T \underline{x} + \underline{x}^T h}_{\substack{\uparrow = \uparrow \\ \text{(inner product!)}}}) = \frac{1}{\|\underline{x}\|_2} \underline{x}^T h$$

$$\mathcal{D}\{x \mapsto x \|x\|_2\} h = \|x\|_2 \cdot h + x \cdot \frac{1}{\|x\|_2} x^T h = \underbrace{(\|x\|_2 \cdot I + \frac{x x^T}{\|x\|_2})}_{\text{Jacobian}} h$$

$$\Rightarrow \text{For } F(x) = A(x)x - b$$

$$\mathcal{D}F(x) = \underbrace{I + \|x\|_2 \cdot I}_{A(x)} + \frac{x x^T}{\|x\|_2}$$

$\uparrow$  rank-1 modification of  $A(x)$        $\uparrow$  tensor product: rank-1-matrix

Case study: Matrix inversion a la Newton

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

$$[F(A^{-1}) = A - (A^{-1})^{-1} = A - A = 0] \quad [F: \mathbb{R}^{n,n} \rightarrow \mathbb{R}^{n,n}]$$

What is  $\mathcal{D}F$ ?      product!

$$\text{Inv}(X) := X^{-1} \Leftrightarrow \text{Inv}(X) \cdot X = I$$

Implicit differentiation: "d":  $\mathcal{D}\text{Inv}(X)H \cdot X + \text{Inv}(X) \cdot H = 0$

by product rule

$$\Rightarrow \mathcal{D}\text{Inv}(X)H = -X^{-1} \cdot H \cdot X^{-1} = -\mathcal{D}F(X)H$$

Newton update  $S$  solves:  $[\mathcal{D}F(X^{(k)})S = F(X^{(k)})]$

$$\Leftrightarrow (X^{(k)})^{-1} S (X^{(k)})^{-1} = A - (X^{(k)})^{-1}, S \in \mathbb{R}^{n,n}$$

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

$\Rightarrow$  Newton iteration:

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

$$= X^{(k)} (2I - A X^{(k)})$$

$$X^{(k)} \rightarrow A^{-1} \text{ for } k \rightarrow \infty$$

\* Newton iteration:  $X^{(k+1)} = X^{(k)} - \mathcal{D}F(X^{(k)})^{-1} F(X^{(k)})$

$\triangleq$  Newton update, Newton correction

$$S := \mathcal{D}F(X^{(k)})^{-1} F(X^{(k)}) \Leftrightarrow \mathcal{D}F(X^{(k)})S = F(X^{(k)})$$

## 2.4.2. Convergence of Newton's method

Newton's iteration for  $F(x) = 0$  as fixed point iteration:

$$x^{(k+1)} = x^{(k)} - \mathcal{D}F(x^{(k)})^{-1} F(x^{(k)})$$

$$\Leftrightarrow x^{(k+1)} = \phi(x^{(k)}), \phi(x) = x - \mathcal{D}F(x)^{-1} F(x)$$

$$\mathcal{D}\phi(x)h = \cancel{h} - \mathcal{D}\{x \mapsto \mathcal{D}F(x)^{-1}\}h \cdot F(x) - \mathcal{D}F(x)^{-1} \mathcal{D}F(x)h$$

[product rule!]  $= -\mathcal{D}\{\dots\}h \cdot F(x)$

$$F(x^*) = 0 \Rightarrow \mathcal{D}\phi(x^*) = 0$$

Lemma 2.2.18  $\Rightarrow$  Local quadratic conv.!

### 2.4.3. Termination of Newton iteration

(local) quad. cvg.

$$\|x^{(k+1)} - x^*\| \ll \|x^{(k)} - x^*\| \quad \text{only few steps will suffice}$$

$$\Rightarrow \|x^{(k+1)} - x^{(k)}\| \approx \|x^{(k)} - x^*\| \quad [n \gg 1 : \text{a single step is expensive}]$$

$\Rightarrow$  Correction based termination :

In row form: One redundant step!

$\triangleright$  Idea: use simplified Newton correction

$$\Delta \bar{x}^{(k)} := JF[x^{(k)}]^{-1} F(x^{(k)}) \rightarrow \text{Effort } O(n^2)$$

↑  
LU-decomposition available

$$\text{STOP, if } \|\Delta \bar{x}^{(k)}\| \leq \tau_{\text{rel}} \|x^{(k)}\|$$

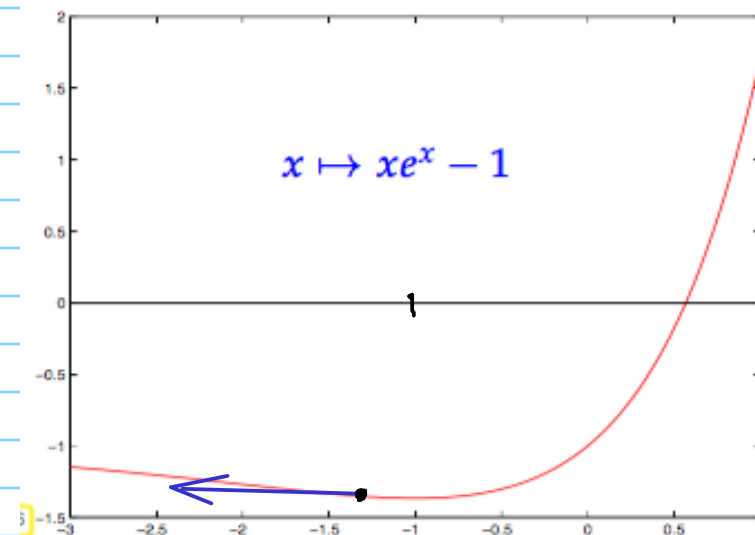
$$\text{or } \|\Delta \bar{x}^{(k)}\| \leq \tau_{\text{abs}}$$

Note: Affine invariant, because this is true for  $\Delta \bar{x}^{(k)}$ !

### 2.4.4. Damped Newton Method

"Usually": convergence really local!

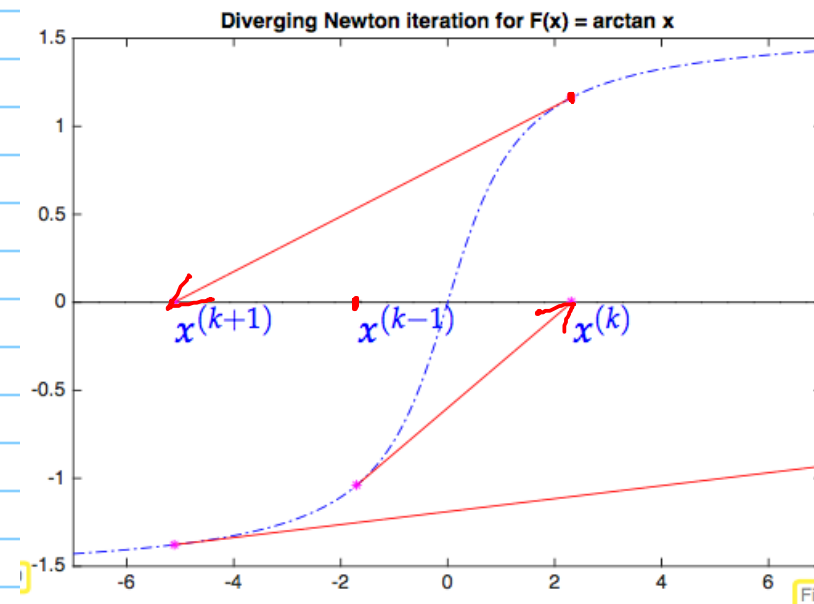
Study in 1D:



$$x^{(0)} < -1$$

$$x^{(10)} \rightarrow -\infty$$

"wrong direction"  
*hopeless*



"Overshooting Newton  
corrections"



we observe "overshooting" of Newton correction

Idea: **damping** of Newton correction:

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

Terminology:  $\lambda^{(k)}$  = damping factor,  $0 < \lambda^{(k)} \leq 1$

A time-tested heuristics: simplified Newton correction

### Affine invariant damping strategy

Choice of damping factor: affine invariant **natural monotonicity test** [7, Ch. 3]:

choose "maximal"  $0 < \lambda^{(k)} \leq 1$ :  $\|\Delta\bar{\mathbf{x}}(\lambda^{(k)})\| \leq (1 - \frac{\lambda^{(k)}}{2}) \|\Delta\mathbf{x}^{(k)}\|_2$  (2.4.49)

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

Practice: If (2.4.49) fails  $\Rightarrow \lambda \leftarrow \frac{1}{2}$

If (2.4.49) satisfied: - 1 step with damping factor  
 -  $2 \cdot \lambda$  for next step

[\*] A convergence monitor: warns user, when method fails.

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

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

We observe that damping is effective and asymptotic quadratic convergence is recovered.

$k$	$\lambda^{(k)}$	$x^{(k)}$	$F(x^{(k)})$
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	-0.00000000000001	-0.00000000000001

### C++11 code 2.4.50: 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;
19         do {
20             lambda /= 2;
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             std::cout << "Inner: |stn| = " << stn << std::endl;
26         }
27         while (stn > (1-lambda/2)*sn); // Natural monotonicity test
28         x = xn; // Now: xn accepted as new iterate
29         lambda = std::min(2.0*lambda, 1.0); // Try to mitigate damping
30     }
31     // Termination based on simplified Newton correction
32     while ((stn > rtol*x.norm()) && (stn > atol));
33 }

```

