

(Full) Newton Method

$\underline{x}^{(e)}$

- Alternatively, we can write:

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

where the vector $\delta \mathbf{x}^{(k)}$ is referred to as the **Newton step**, or **Newton correction** and is the solution of the linear system:

$$\mathbf{J}(\mathbf{x}^{(k)}) \delta \mathbf{x}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)}).$$

*Solve this
system involving
the Jacobian
matrix of \mathbf{F}*

$$\underline{x}^{(r)} = \underline{x}^{(e)} + \delta \underline{x}^{(e)}$$



(Full) Newton Method

- The computational framework presented above can be extended to higher dimensions. Given an initial iterate $\mathbf{x}^{(0)}$, Newton's method for solving (1) takes the form:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{J}(\mathbf{x}^{(k)})^{-1} \mathbf{F}(\mathbf{x}^{(k)}) \quad (6)$$

where $\mathbf{J} \in \mathbb{R}^{N \times N}$ is the Jacobian matrix of \mathbf{F} :

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} \nabla f_1(\mathbf{x})^T \\ \nabla f_2(\mathbf{x})^T \\ \vdots \\ \nabla f_N(\mathbf{x})^T \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_N} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_2(\mathbf{x})}{\partial x_N} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_N(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_N(\mathbf{x})}{\partial x_N} \end{bmatrix} \quad \mathcal{J}_{*1} \quad \mathcal{J}_{*2} \quad \mathcal{J}_{*N}$$

$$\mathbf{J}_{*j} = \mathbf{J}(\underline{\mathbf{x}}) \underline{\mathbf{e}}_j \quad \underline{\mathbf{e}}_j = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \mathbf{J}_{*N}$$

Terminating Newton's Method

- One common criterion used to terminate Newton's method is based on the norm of the nonlinear residual $\|\mathbf{F}(\mathbf{x}^{(k+1)})\|$.
- Typically, a combination of both a **relative error tolerance** τ_r and **absolute error tolerance** τ_a is used to halt the iteration:

$$\underbrace{\|\mathbf{F}(\mathbf{x}^{(k+1)})\|}_{\text{nonlinear residual}} \leq \tau_r \|\mathbf{F}(\mathbf{x}^{(0)})\| + \tau_a. \quad (7)$$

(Full) Newton Method

n^2 derivatives

An algorithm for Newton's method is summarised as follows.

Newton's method

Choose $\mathbf{x}^{(0)}$

While not converged for $k = 0, 1, \dots$

Solve $\mathbf{J}(\mathbf{x}^{(k)})\delta\mathbf{x}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})$ for $\delta\mathbf{x}^{(k)}$

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

End

$$\|\mathbf{F}(\mathbf{x}^{(k+1)})\|_2 < \epsilon$$

Generate \mathcal{T} : n^2 fractions
PLU decomposition.

$$\|\mathbf{x}\|_p = \left[\sum_{i=1}^n |\mathbf{x}_i|^p \right]^{\frac{1}{p}},$$

$$\begin{aligned} p=1 : \quad & \|\mathbf{x}\|_1 = \sum_{i=1}^n |\mathbf{x}_i| \\ p=2 : \quad & \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n |\mathbf{x}_i|^2} \\ p=\infty : \quad & \|\mathbf{x}\|_\infty = \max_{i \in \mathbb{N}} |\mathbf{x}_i|. \end{aligned}$$

Chord Method

An algorithm for the Chord method is summarised as follows.

Chord method

Choose $\mathbf{x}^{(0)}$ and set $\bar{\mathbf{J}} = \mathbf{J}(\mathbf{x}^{(0)})$

While not converged for $K = 0, 1, \dots$

Solve $\bar{\mathbf{J}}\delta\mathbf{x}^{(K)} = -\mathbf{F}(\mathbf{x}^{(K)})$ for $\delta\mathbf{x}^{(K)}$

$\mathbf{x}^{(K+1)} = \mathbf{x}^{(K)} + \delta\mathbf{x}^{(K)}$

End

Shamanskii Method

Alternatively, the Jacobian matrix may be *periodically* updated. For example, alternating a full Newton step with a sequence of Chord steps leads to a method often attributed to Shamanskii.

$$T_{\text{tol}}, \text{tol}, \rho^* = 0.5$$

Shamanskii method

Choose $\mathbf{x}^{(0)}$

While not converged for $k = 0, 1, \dots$
If $k \equiv 0 \pmod m$ or $\rho > \rho^*$

Compute/Update $\bar{\mathbf{J}} = \mathbf{J}(\mathbf{x}^{(k)})$

End

Solve $\bar{\mathbf{J}}\delta\mathbf{x}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})$ for $\delta\mathbf{x}^{(k)}$

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

End

$$\left[\begin{array}{l} \text{err} = \| \mathbf{F}^{(0)} \| \\ \text{tol} = \text{tol} + \text{tol} * \text{err} \end{array} \right]$$

ρ is the ratio
 ρ of nonlinear residues

$$\left[\begin{array}{l} \text{err}_{\text{old}} = \text{err} ; \\ \text{err} = \| \mathbf{F}^{(k)} \| \\ \rho = \frac{\text{err}}{\text{err}_{\text{old}}} ; \end{array} \right]$$

$$\frac{\text{err}}{\text{err}_{\text{old}}} < \rho^* = 0.5$$
$$\text{err} < 0.5 \text{ err}_{\text{old}}$$



Convergence Analysis

Definition (Lipschitz continuity)

Let $D \subset \mathbb{R}^N$ and $\mathbf{J} : D \rightarrow \mathbb{R}^{N \times N}$. Then \mathbf{J} is **Lipschitz continuous** on D , with **Lipschitz constant** γ if

$$\|\mathbf{J}(\mathbf{x}) - \mathbf{J}(\mathbf{y})\| \leq \gamma \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in D,$$

and we write $\mathbf{J} \in \text{Lip}_\gamma(D)$.

The Standard Assumptions

- 1 $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ has a solution \mathbf{x}^* on $D \subset \mathbb{R}^N$.
- 2 $\mathbf{J} \in \text{Lip}_\gamma(D)$.
- 3 $\mathbf{J}(\mathbf{x})$ is nonsingular and $\|\mathbf{J}(\mathbf{x})^{-1}\| \leq \beta \quad \forall \mathbf{x} \in D$.



Convergence Analysis

The following formulation of the fundamental theorem of calculus will be important in the theory that follows.

Theorem

Let $\mathbf{F} : D \rightarrow \mathbb{R}^N$ be differentiable on an open subset $D \subset \mathbb{R}^N$ and let $\mathbf{x} \in D$. Then for any nonzero vector \mathbf{h} , such that $\mathbf{x} + \mathbf{h} \in D$, we have

$$\mathbf{F}(\mathbf{x} + \mathbf{h}) = \mathbf{F}(\mathbf{x}) + \int_0^1 \mathbf{J}(\mathbf{x} + t\mathbf{h}) \mathbf{h} dt$$

where \mathbf{J} is the Jacobian matrix of \mathbf{F} .



Finite Difference Jacobian

- If \mathbf{F} has a complicated form, or if it is only available as a computer code rather than an analytic formula, it may not be possible to compute the Jacobian matrix analytically. In this case finite difference approximations to the derivatives can be used instead.
- Typically a first order forward difference approximation is used, whereby the j^{th} column of the Jacobian matrix is approximated as follows:

$$\mathbf{J}_{*,j}(\mathbf{x}) = \mathbf{J}(\mathbf{x})\mathbf{e}_j \approx \tilde{\mathbf{J}}_{*,j}(\mathbf{x}) = \frac{\mathbf{F}(\mathbf{x} + h\mathbf{e}_j) - \mathbf{F}(\mathbf{x})}{h} \quad (9)$$

for a suitable shift value $h > 0$, where \mathbf{e}_j is the j^{th} column of the $N \times N$ identity matrix.

- This strategy can be applied to either Newton's method, or the Chord and Shamanskii methods.



Idea: choose $\underline{x} = \varepsilon \underline{e}_j$, $\underline{e}_j = (0, \dots, 0, 1, 0, \dots, 0)^T$.

\underline{x} is a small number.

$$F(\underline{x} + \varepsilon \underline{e}_j) = F(\underline{x}) + \int_0^1 \varepsilon J(\underline{x} + t\varepsilon \underline{e}_j) \underline{e}_j \cdot dt$$

$$J_{*j}(\underline{x}) = J(\underline{x}) \underline{e}_j \approx \frac{F(\underline{x} + \varepsilon \underline{e}_j) - F(\underline{x})}{\varepsilon} = \tilde{J}_{*j}(\underline{x})$$

what is the error in j th column:
 First error in \underline{x} $\tilde{J}_{*j}(\underline{x}) < \varepsilon$?

$$\| J_{*j}(\underline{x}) - \tilde{J}_{*j}(\underline{x}) \| = J(\underline{x}) \underline{e}_j - \left\{ \frac{F(\underline{x} + \varepsilon \underline{e}_j) - F(\underline{x})}{\varepsilon} \right\}$$

$$= J(\underline{x}) \underline{e}_j - \int_0^1 J(\underline{x} + t\varepsilon \underline{e}_j) \underline{e}_j \cdot dt$$

$$= \int_0^1 J(\underline{x}) \underline{e}_j \cdot dt - \int_0^1 J(\underline{x} + t\varepsilon \underline{e}_j) \underline{e}_j \cdot dt$$

$$\begin{aligned}
\|\mathcal{T}^*j - \tilde{\mathcal{T}}^*j\| &= \left\| \int_0^1 [\mathcal{T}(\underline{x}) - \mathcal{T}(\underline{x} + t\varepsilon e_j \cdot)] \varepsilon e_j dt \right\| \quad \text{(using matrix-vector multiplication)} \\
&\leq \int_0^1 \|(\mathcal{T}(\underline{x}) - \mathcal{T}(\underline{x} + t\varepsilon e_j \cdot)) \varepsilon e_j\| dt \\
&\quad \text{matrix-vector multiplication} \\
&\leq \int_0^1 \|\mathcal{T}(\underline{x}) - \mathcal{T}(\underline{x} + t\varepsilon e_j \cdot)\| dt \\
&\quad \text{continuous function} \\
&\leq \int_0^1 \gamma \|\underline{x} - \underline{x} - t\varepsilon e_j\| dt \quad \left\{ \begin{array}{l} \text{Assume that} \\ \mathcal{T} \in \text{Lip}_r(D) \end{array} \right\} \\
&\leq \varepsilon \int_0^1 \gamma |t| dt = \frac{\varepsilon \gamma}{2} \\
&\leq \|\mathcal{T} - \tilde{\mathcal{T}}\|_2^2 = \sum_{j=1}^n \|(\mathcal{T} - \tilde{\mathcal{T}}) e_j\|_2^2 \\
&\leq n \frac{\sum \gamma^2}{2} = \boxed{\left(\|\mathcal{T} - \tilde{\mathcal{T}}\|_2 \leq \sqrt{n} \frac{\varepsilon \gamma}{2} \right)}
\end{aligned}$$

$$F(\bar{x}) = \sum_{i=1}^n f_i(\bar{x})$$

$$\bar{F}(\bar{x}) = \frac{F(\bar{x} + \varepsilon e_i) - F(\bar{x})}{\varepsilon}$$

$$\bar{f}_i(\bar{x}) = \frac{f_i(\bar{x} + \varepsilon e_i) - f_i(\bar{x})}{\varepsilon}$$

$$\begin{pmatrix} \bar{f}_1(\bar{x}) & \bar{f}_2(\bar{x}) & \dots & \bar{f}_n(\bar{x}) \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(\bar{x}) & \frac{\partial f_1}{\partial x_2}(\bar{x}) & \dots & \frac{\partial f_1}{\partial x_n}(\bar{x}) \\ \frac{\partial f_2}{\partial x_1}(\bar{x}) & \frac{\partial f_2}{\partial x_2}(\bar{x}) & \dots & \frac{\partial f_2}{\partial x_n}(\bar{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1}(\bar{x}) & \frac{\partial f_n}{\partial x_2}(\bar{x}) & \dots & \frac{\partial f_n}{\partial x_n}(\bar{x}) \end{pmatrix}$$

Finite Difference Jacobian

Theorem

Let $\mathbf{F} : D \rightarrow \mathbb{R}^N$ for some open subset $D \subset \mathbb{R}^N$ and let $\mathbf{J} \in Lip_\gamma(D)$, then

$$\|\mathbf{J}_{*,j}(\mathbf{x}) - \tilde{\mathbf{J}}_{*,j}(\mathbf{x})\| \leq \frac{\gamma}{2} h$$

where $\tilde{\mathbf{J}}_{*,j}$ is computed using finite differences as described in equation (9).

One reasonable choice for the shift parameter is

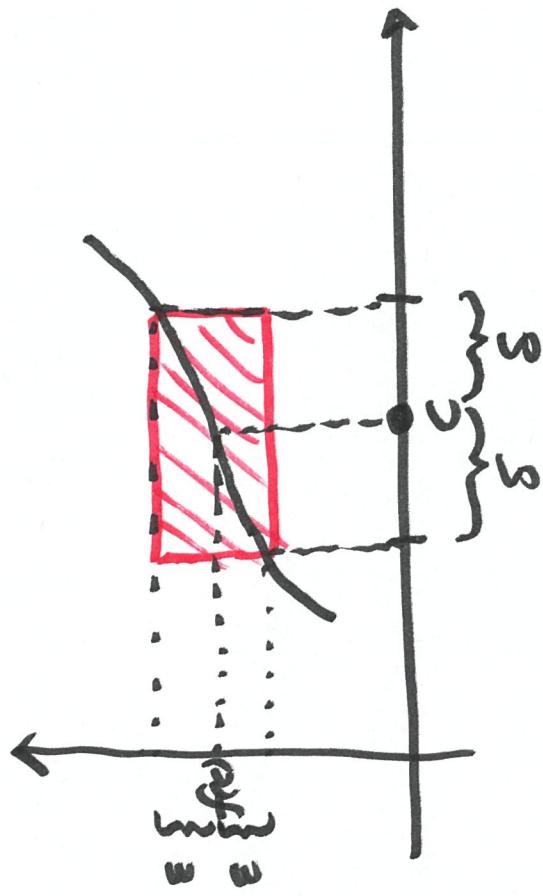
$$h = \begin{cases} \sqrt{\epsilon} \|\mathbf{x}\|_2 & \mathbf{x} \neq 0 \\ \sqrt{\epsilon} & \mathbf{x} = 0 \end{cases},$$

where ϵ is the unit roundoff (denoted by `eps` and equal to $2.2204e-16$ in MATLAB).



Definition: Let $D \subset \mathbb{R}$ and $f: D \rightarrow \mathbb{R}$

then f is continuous at c if for every $\epsilon > 0$ there is a $\delta > 0$ such that whenever $x \in D$ and $|x - c| < \delta$ then $|f(x) - f(c)| < \epsilon$



for $|x - c| < \delta$
for $f(x)$ must be in
the shaded region

Note: δ depends on ϵ
and c .

Lipschitz continuity:

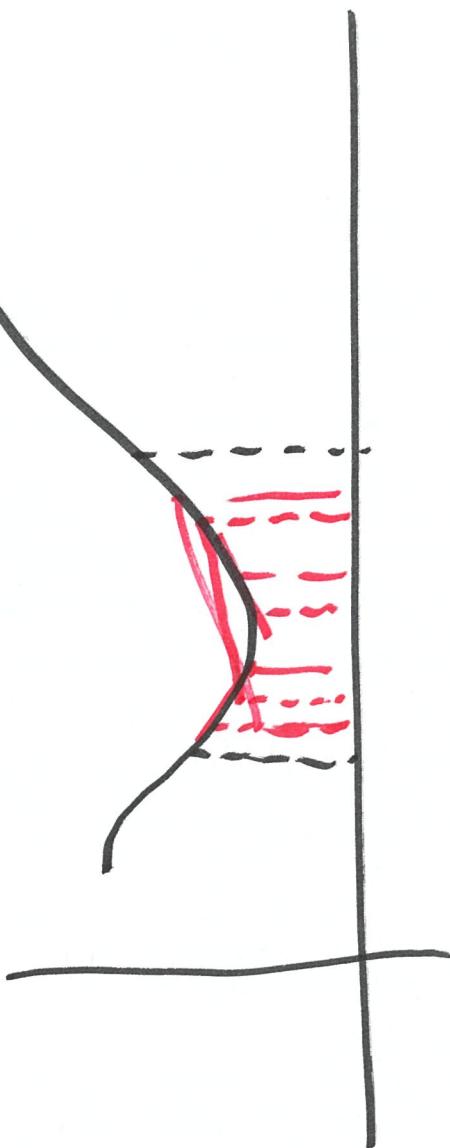
$f: D \rightarrow \mathbb{R}$ such that there exists a number γ such that for all $x, y \in D$

$$|f(x) - f(y)| \leq \gamma |x - y|$$

Interpretation:

$$\left| \frac{f(x) - f(y)}{x - y} \right| \leq \gamma$$

slope of secant lines.



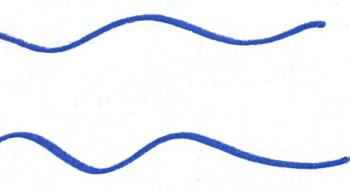
Tridiagonal Jacobian

Consider a tridiagonal Jacobian:

$$\mathbf{J} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & & & & & & & \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & & & & & & \\ & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} & & & & & \\ & & \frac{\partial f_4}{\partial x_3} & \frac{\partial f_4}{\partial x_4} & & & & \\ & & & \frac{\partial f_5}{\partial x_4} & \frac{\partial f_5}{\partial x_5} & & & \\ & & & & \frac{\partial f_6}{\partial x_5} & \frac{\partial f_6}{\partial x_6} & & \\ & & & & & \frac{\partial f_7}{\partial x_6} & \frac{\partial f_7}{\partial x_7} & \\ & & & & & & \frac{\partial f_7}{\partial x_7} & \frac{\partial f_7}{\partial x_8} \\ & & & & & & & \ddots \end{bmatrix}$$



$\mathbf{J} =$



We now show that a Jacobian matrix having this structure can be approximated using only 3 additional function evaluations.

Tridiagonal Jacobian

Define the vectors in \mathbb{R}^N :

$$\mathbf{s}_1 = (1, 0, 0, 1, 0, 0, 1, \dots)^T, \quad \mathbf{s}_2 = (0, 1, 0, 0, 1, 0, 0, \dots)^T,$$

$$\mathbf{s}_3 = (0, 0, 1, 0, 0, 1, 0, \dots)^T,$$

and consider the products $\mathbf{J}\mathbf{s}_1$, $\mathbf{J}\mathbf{s}_2$ and $\mathbf{J}\mathbf{s}_3$

- $\mathbf{J}\mathbf{s}_1 = \mathbf{J}_{*,1} + \mathbf{J}_{*,4} + \mathbf{J}_{*,7} + \dots = \sum_{j=1:3:N} \mathbf{J}_{*,j}$ is the sum of columns 1, 4, 7, ..., and so on.
- $\mathbf{J}\mathbf{s}_2 = \mathbf{J}_{*,2} + \mathbf{J}_{*,5} + \mathbf{J}_{*,8} + \dots = \sum_{j=2:3:N} \mathbf{J}_{*,j}$ is the sum of columns 2, 5, 8, ..., and so on.
- $\mathbf{J}\mathbf{s}_3 = \mathbf{J}_{*,3} + \mathbf{J}_{*,6} + \mathbf{J}_{*,9} + \dots = \sum_{j=3:3:N} \mathbf{J}_{*,j}$ is the sum of columns 3, 6, 9, ..., and so on.

where for example $1 : 3 : N$ refers to MATLAB's colon notation.



Tridiagonal Jacobian

The key observation is that each of the above summations involve vectors that do not possess entries in the same row. Therefore:

$$\mathbf{J}\mathbf{s}_1 = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} \\ \frac{\partial f_2}{\partial x_1} \\ \frac{\partial f_3}{\partial x_4} \\ \frac{\partial f_4}{\partial x_4} \\ \frac{\partial f_5}{\partial x_4} \\ \vdots \end{bmatrix}, \quad \mathbf{J}\mathbf{s}_2 = \begin{bmatrix} \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_2} \\ \frac{\partial f_3}{\partial x_2} \\ \frac{\partial f_4}{\partial x_5} \\ \frac{\partial f_5}{\partial x_5} \\ \frac{\partial f_6}{\partial x_5} \\ \frac{\partial f_7}{\partial x_7} \\ \vdots \end{bmatrix}, \quad \mathbf{J}\mathbf{s}_3 = \begin{bmatrix} \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_3} \\ \frac{\partial f_3}{\partial x_3} \\ \frac{\partial f_4}{\partial x_6} \\ \frac{\partial f_5}{\partial x_6} \\ \frac{\partial f_6}{\partial x_6} \\ \frac{\partial f_7}{\partial x_9} \\ \vdots \end{bmatrix}.$$

Tridiagonal Jacobian

For a given vector $\mathbf{x} \in \mathbb{R}^N$, each of the above vectors can be approximated as follows:

$$\mathbf{J}(\mathbf{x})\mathbf{s}_j \approx \frac{\mathbf{F}(\mathbf{x} + h_j \mathbf{s}_j) - \mathbf{F}(\mathbf{x})}{h_j} \quad (10)$$

where the shift parameter is chosen as

$$h_j = \begin{cases} \frac{\sqrt{\epsilon} \|\mathbf{x}\|_2}{\|\mathbf{s}_j\|_2} & \mathbf{x} \neq 0 \\ \frac{\sqrt{\epsilon}}{\|\mathbf{s}_j\|_2} & \mathbf{x} = 0. \end{cases}$$

The above strategy can also be extended to a general banded matrix with bandwidth greater than 3.

