Non-linear problems in n variable

Lectures for PHD course on Numerical Optimization

Enrico Bertolazzi

DII – Università di Trento





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes			

Outline

1	The Newton Raphson		
2	The Frobenius matrix norm		
3	The Broyden method		
4	The dumped Broyden method		
5	Stopping criteria and q -order estimation	n	
			m m
Enrico Bert	rolazzi — Non-linear problems in $\it n$ variable		2/78
Notes			

The problem to solve

Problem

Given $\mathbf{F}:D\subseteq\mathbb{R}^n\mapsto\mathbb{R}^n$

Find $x_{\star} \in D$ for which $\mathbf{F}(x_{\star}) = 0$.

Example

Let

$$\mathbf{F}(\mathbf{x}) = \begin{pmatrix} x_1^2 + x_2^3 + 7 \\ x_1 + x_2 + 1 \end{pmatrix}$$

which has $\mathbf{F}(x_{\star}) = \mathbf{0}$ for $x_{\star} = (1, -2)^T$.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

INOTES		

1 The Newton Raphson

- 2 The Frobenius matrix norm
- 3 The Broyden method
- 4 The dumped Broyden method
- 5 Stopping criteria and q-order estimation





4.70

Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

IN	0	Ť	е	5

The Newton procedure

Consider the following map

$$\mathbf{F}(\mathbf{x}) = \begin{pmatrix} x_1^2 + x_2^3 + 7 \\ x_1 + x_2 + 1 \end{pmatrix}$$

we known an approximation of a root $x_0 \approx (1.1, -1.9)^T$.

lacksquare Setting $oldsymbol{x}_1 = oldsymbol{x}_0 + oldsymbol{p}$ we obtain 1

$$\mathbf{F}(\boldsymbol{x}_0 + \boldsymbol{p}) = \begin{pmatrix} 1.351 \\ 0.2 \end{pmatrix} + \begin{pmatrix} 2.2 & 10.83 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \end{pmatrix} + \vec{\boldsymbol{\mathcal{O}}}(\|\boldsymbol{p}\|^2)$$

if x_0 is a good approximation of a root of $\mathbf{F}(x)$ then $\vec{\mathcal{O}}(\|\mathbf{p}\|^2)$ is a small vector.



 1 Here $\vec{\mathcal{O}}(x)$ means $(\mathcal{O}(x),\ldots,\mathcal{O}(x))^{T}$



Enrico Bertolazzi — Non-linear problems in n variable

Notes	

The Newton procedure

■ Neglecting $\vec{\mathcal{O}}(\|\boldsymbol{p}\|^2)$ and solving

$$\begin{pmatrix} 1.351 \\ 0.2 \end{pmatrix} + \begin{pmatrix} 2.2 & 10.83 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \end{pmatrix} = \mathbf{0}$$

we obtain $p = (-0.094438, -0.105562)^T$.

Now we set

$$m{x}_1 = m{x}_0 + m{p} = egin{pmatrix} 1.005562 \\ -2.0055612 \end{pmatrix}$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

6/78

Notes

The Newton procedure

Considering

$$\mathbf{F}(x_1 + q) = \begin{pmatrix} -0.05576 \\ 810^{-7} \end{pmatrix} + \begin{pmatrix} 2.0111 & 12.0668 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} q_1 \\ q_2 \end{pmatrix} + \mathbf{\mathcal{O}}(\|\mathbf{q}\|^2)$$

■ Neglecting $\vec{\mathcal{O}}(\|q\|^2)$ and solving

$$\begin{pmatrix} -0.05576 \\ 810^{-7} \end{pmatrix} + \begin{pmatrix} 2.0111 & 12.0668 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} q_1 \\ q_2 \end{pmatrix} = \mathbf{0}$$

we obtain $q = (-0.0055466, 0.0055458)^T$.

Now we set $x_2 = x_1 + q = (1.000015, -2.000015)^T$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

///8

Notes

Modern point of view of Newton scheme (1/2)

The previous procedure can be resumed as follows:

- 1 Consider the following function $\mathbf{F}(x)$. We known an approximation of a root x_0 .
- 2 Expand by Taylor series

$$\mathbf{F}(oldsymbol{x}) = \mathbf{F}(oldsymbol{x}_0) +
abla \mathbf{F}(oldsymbol{x}_0) (oldsymbol{x} - oldsymbol{x}_0) + oldsymbol{\mathcal{O}}(\|oldsymbol{x} - oldsymbol{x}_0\|^2)$$

3 Drop the term $\vec{\mathcal{O}}(\|x-x_0\|^2)$ and solve

$$\mathbf{0} = \mathbf{F}(\boldsymbol{x}_0) + \nabla \mathbf{F}(\boldsymbol{x}_0)(\boldsymbol{x} - \boldsymbol{x}_0)$$

Call x_1 this solution.

Repeat 1-3 with x_1 , x_2 , x_3 , ...





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Modern point of view of Newton scheme (2/2)

Algorithm (Newton iterative scheme)

Let x_0 assigned, then for k = 0, 1, 2, ...

1 Solve for p_k :

$$abla \mathbf{F}(oldsymbol{x}_k) oldsymbol{p}_k + \mathbf{F}(oldsymbol{x}_k) = \mathbf{0}$$

2 Update

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{p}_k$$





9/78

Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Standard Assumptions

In the study of convergence of numerical scheme, some standard regularity assumption are assumed for the function ${\bf F}(x)$.

Assumption (Standard Assumptions)

The function $\mathbf{F}: D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ is continuous, differentiable with Lipschitz derivative $\nabla \mathbf{F}(x)$. i.e.

$$\|\nabla \mathbf{F}(\mathbf{x}) - \nabla \mathbf{F}(\mathbf{y})\| \le \gamma \|\mathbf{x} - \mathbf{y}\| \qquad \forall \mathbf{x}, \mathbf{y} \in D \subset \mathbb{R}^n$$

Lemma (Taylor like expansion)

Let $\mathbf{F}(x)$ satisfy the standard assumptions, then

$$\|\mathbf{F}(y) - \mathbf{F}(x) - \nabla \mathbf{F}(x)(y - x)\| \le \frac{\gamma}{2} \|x - y\|^2 \quad \forall x, y \in D \subset \mathbb{R}^n$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Proof: From basic Calculus:

$$\mathbf{F}(\boldsymbol{y}) - \mathbf{F}(\boldsymbol{x}) = \int_0^1 \nabla \mathbf{F}(\boldsymbol{x} + t(\boldsymbol{y} - \boldsymbol{x}))(\boldsymbol{y} - \boldsymbol{x}) dt$$

subtracting on both side $abla \mathbf{F}(x)(y-x)$ we have

$$\mathbf{F}(y) - \mathbf{F}(x) - \nabla \mathbf{F}(x)(y - x) =$$

$$\int_{0}^{1} \left[\nabla \mathbf{F}(x + t(y - x)) - \nabla \mathbf{F}(x) \right] (y - x) dt$$

and taking the norm

$$\|\mathbf{F}(\boldsymbol{y}) - \mathbf{F}(\boldsymbol{x}) - \nabla \mathbf{F}(\boldsymbol{x})(\boldsymbol{y} - \boldsymbol{x})\| \le \int_0^1 \gamma t \|\boldsymbol{y} - \boldsymbol{x}\|^2 dt$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Lemma (Jacobian norm control)

Let $\mathbf{F}(x)$ satisfying standard assumptions, and $\nabla \mathbf{F}(x_\star)$ non singular. Then there exists $\delta>0$ such that for all $\|x-x_\star\|\leq \delta$ we have

$$2^{-1} \|\nabla \mathbf{F}(\boldsymbol{x})\| \le \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\| \le 2 \|\nabla \mathbf{F}(\boldsymbol{x})\|$$

and

$$2^{-1} \|\nabla \mathbf{F}(\mathbf{x})^{-1}\| \le \|\nabla \mathbf{F}(\mathbf{x}_{\star})^{-1}\| \le 2 \|\nabla \mathbf{F}(\mathbf{x})^{-1}\|$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

$$\|\nabla \mathbf{F}(\boldsymbol{x})\| \le \|\nabla \mathbf{F}(\boldsymbol{x}) - \nabla \mathbf{F}(\boldsymbol{x}_{\star})\| + \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\|$$

$$\le \gamma \|\boldsymbol{x} - \boldsymbol{x}_{\star}\| + \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\|$$

$$\le (3/2) \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\| \le 2 \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\|$$

again choosing $\gamma \delta \leq 2^{-1} \| \nabla \mathbf{F}(\boldsymbol{x}_\star) \|$

$$\|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\| \leq \|\nabla \mathbf{F}(\boldsymbol{x}_{\star}) - \nabla \mathbf{F}(\boldsymbol{x})\| + \|\nabla \mathbf{F}(\boldsymbol{x})\|$$
$$\leq \gamma \|\boldsymbol{x} - \boldsymbol{x}_{\star}\| + \|\nabla \mathbf{F}(\boldsymbol{x})\|$$
$$\leq 2^{-1} \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\| + \|\nabla \mathbf{F}(\boldsymbol{x})\|$$

so that $2^{-1} \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})\| \leq \|\nabla \mathbf{F}(\boldsymbol{x})\|$.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Proof: (Proof. (2/3)) From the continuity of the determinant there exists a neighbor with $\nabla \mathbf{F}(x)$ non singular for all $\|x - x_{\star}\| \leq \delta$.

$$\begin{aligned} \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} - \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\| \\ & \leq \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star}) - \nabla \mathbf{F}(\boldsymbol{x}) \right\| \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\| \\ & \leq \gamma \left\| \boldsymbol{x} - \boldsymbol{x}_{\star} \right\| \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\| \end{aligned}$$

and choosing δ such that $\gamma\delta\left\|\nabla\mathbf{F}(\boldsymbol{x}_{\star})^{-1}\right\|\leq2^{-1}$ we have

$$\|\nabla \mathbf{F}(\boldsymbol{x})^{-1} - \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1}\| \le 2^{-1} \|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\|$$

and using this last inequality

$$\|\nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1}\| \leq \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} - \nabla \mathbf{F}(\boldsymbol{x})^{-1}\| + \|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\|$$
$$\leq (3/2) \|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\| \leq 2 \|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\|$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes			

Proof: (Proof.

(3/3)) Using last inequality again

$$\|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\| \leq \|\nabla \mathbf{F}(\boldsymbol{x})^{-1} - \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1}\| + \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1}\|$$
$$\leq 2^{-1} \|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\| + \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1}\|$$

so that

$$2^{-1} \|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\| \le \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1}\|$$

choosing δ such that for all $\|x-x_\star\| \leq \delta$ we have $\nabla \mathbf{F}(x)$ non singular and $\gamma \delta \leq 2^{-1} \|\nabla \mathbf{F}(x_\star)\|$ and $\gamma \delta \|\nabla \mathbf{F}(x_\star)^{-1}\| \leq 2^{-1}$ then the inequality of the lemma are true.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Theorem (Local Convergence of Newton method)

Let $\mathbf{F}(x)$ satisfying standard assumptions, and x_{\star} a simple root (i.e. $\nabla \mathbf{F}(x_{\star})$ non singular). Then, if $\|x_0 - x_{\star}\| \leq \delta$ with $C\delta \leq 1$ where

$$C = \gamma \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\|$$

then, the sequence generated by Newton method satisfies:

- $\| m{x}_k m{x}_\star \| \leq \delta \; ext{for} \; k = 0, 1, 2, 3, \dots$
- $\| oldsymbol{x}_{k+1} oldsymbol{x}_{\star} \| \leq C \| oldsymbol{x}_k oldsymbol{x}_{\star} \|^2 ext{ for } k = 0, 1, 2, 3, \dots$
- $\lim_{k\to\infty} oldsymbol{x}_k = oldsymbol{x}_\star.$
- The point 2 of the theorem is the second q-order of convergence of Newton method.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Proof: Consider a Newton step with $\|x_k - x_\star\| \leq \delta$ and

$$egin{aligned} oldsymbol{x}_{k+1} - oldsymbol{x}_{\star} &= oldsymbol{x}_k - oldsymbol{x}_{\star} -
abla \mathbf{F}(oldsymbol{x}_k)^{-1} ig[\mathbf{F}(oldsymbol{x}_k) - \mathbf{F}(oldsymbol{x}_k) - \mathbf{F}(oldsymbol{x}_k) + \mathbf{F}(oldsymbol{x}_{\star}) ig] \ &=
abla \mathbf{F}(oldsymbol{x}_k)^{-1} ig[
abla \mathbf{F}(oldsymbol{x}_k) (oldsymbol{x}_k - oldsymbol{x}_{\star}) - \mathbf{F}(oldsymbol{x}_k) + \mathbf{F}(oldsymbol{x}_{\star}) ig] \end{aligned}$$

taking the norm and using Taylor like lemma

$$\|x_{k+1} - x_{\star}\| \le 2^{-1} \gamma \|x_k - x_{\star}\|^2 \|\nabla \mathbf{F}(x_k)^{-1}\|$$

from Jacobian norm control lemma (slide 12) there exist a δ such that $2 \left\| \nabla \mathbf{F}(x_k)^{-1} \right\| \ge \left\| \nabla \mathbf{F}(x_\star)^{-1} \right\|$ for all $\|x_k - x_\star\| \le \delta$. Reducing eventually δ such that $\gamma \delta \left\| \nabla \mathbf{F}(x_\star)^{-1} \right\| \le 1$ we have

$$\|\boldsymbol{x}_{k+1} - \boldsymbol{x}_{\star}\| \le \gamma \|\nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1}\| \delta \|\boldsymbol{x}_{k} - \boldsymbol{x}_{\star}\|^{2} \le \|\boldsymbol{x}_{k} - \boldsymbol{x}_{\star}\|,$$

So that by induction we prove point 1. Point 2 and 3 follows trivially.





Enrico Bertolazzi — Non-linear problems in n variable

Notes		

Theorem (Newton-Kantorovich)

Let $\mathbf{F}: D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ be a differentiable mapping and let $x_0 \in D$ be such that $\nabla \mathbf{F}(x_0)$ is nonsingular. Let be

$$B(\mathbf{x}_0, \rho) = \{ \mathbf{y} \mid ||\mathbf{x}_0 - \mathbf{y}|| < \rho \},$$

$$\alpha = ||\nabla \mathbf{F}(\mathbf{x}_0)^{-1} \mathbf{F}(\mathbf{x}_0)||,$$

Moreover

- $\overline{B(x_0,\rho)}\subset D;$
- $\qquad \qquad \left\| \nabla \mathbf{F}(\boldsymbol{x}_0)^{-1} (\mathbf{F}(\boldsymbol{x}) \mathbf{F}(\boldsymbol{x}_0)) \right\| \leq \omega \, \|\boldsymbol{x} \boldsymbol{x}_0\| \quad \text{for all} \quad \boldsymbol{x} \in D;$
- $\kappa := \alpha \omega \leq 2^{-1}$;

If the radius ρ is large enough, i.e.

$$\hat{\rho} := \frac{1 - \sqrt{1 - 2\kappa}}{\omega} \le \rho$$

Then:



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Theorem (cont.)

- $lackbox{f F}(x)$ has a zero $x_\star \in \overline{B(x_0,\hat
 ho)}$;
- The open ball $B(x_0, \hat{\rho})$ does not contain any zero of $\mathbf{F}(x)$ different from x_{\star} ;
- The Newton iterative procedure produce sequences belonging to $B(x_0, \hat{\rho})$ that converge to x_* ;
- If $\kappa < 2^{-1}$ then for Newton's method, we have

$$\|\boldsymbol{x}_k - \boldsymbol{x}_\star\| \leq \frac{2\beta\lambda^{2^k}}{1 - \lambda^{2^k}}$$

where

$$\beta = \frac{\sqrt{1 - 2\kappa}}{\omega}, \qquad \lambda = \frac{1 - \kappa - \sqrt{1 - 2\kappa}}{\kappa}$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Proof:

P. Deuflhard and G. Heindl

Affine Invariant Convergence Theorems for Newton's Method and Extensions to Related Methods SIAM Journal on Numerical Analysis, 16, 1979.

Florian A. Potra

The Kantorovich Theorem and interior point methods Math. Program., Ser. A 102, 2005.

J.M. Ortega

The Newton-Kantorovich theorem Amer. Math. Monthly 75, 1968.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

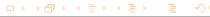
- Newton method converge normally only when x_0 is near x_\star a root of the nonlinear system.
- A way to make a more robust non linear solver is to use the techniques developed for minimization to make a globally convergent nonlinear solver.
- In particular if we consider the merit function

$$f(\boldsymbol{x}) = \frac{1}{2} \left\| \mathbf{F}(\boldsymbol{x}) \right\|^2$$

we have that $\mathsf{f}(x) \geq 0$ and if x_\star is such that $\mathsf{f}(x_\star) = 0$ than we have that

- 1 x_{\star} is a global minimum of f(x);
- $\mathbf{F}(x_{\star}) = \mathbf{0}$, i.e. is a solution of the nonlinear system $\mathbf{F}(x)$.
- \blacksquare So that finding a global minimum of the merit function $\mathbf{f}(x)$ is the same of finding a solution of the nonlinear system $\mathbf{F}(x).$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes			

- We can apply for example the gradient method to the merit function f(x). This produce a slow method.
- Instead, we can use the Newton method to produce a search direction. The resulting method is the following
 - 1 Compute the search direction by solving $abla \mathbf{F}(oldsymbol{x}_k) oldsymbol{d}_k + \mathbf{F}(oldsymbol{x}_k) = \mathbf{0};$
 - 2 Find an approximate solution of the problem $\alpha_k = \operatorname{arg\,min}_{\alpha > 0} \|\mathbf{F}(\boldsymbol{x}_k + \alpha \boldsymbol{d}_k)\|^2$;
 - 3 Update the solution $x_{k+1} = x_k + \alpha_k d_k$.
- lacktriangle The previous algorithm work if the direction d_k is a descent direction.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

22/78	

Notes		

Is d_k a descent direction?

Lemma

The direction d computed as a solution of the problem

$$\nabla \mathbf{F}(x)d + \mathbf{F}(x) = \mathbf{0}$$

is a descent direction.

Proof:

Consider the gradient of $f(x) = (1/2) \|\mathbf{F}(x)\|^2$:

$$\frac{\partial f(\boldsymbol{x})}{\partial x_k} = \frac{1}{2} \frac{\partial \|\mathbf{F}(\boldsymbol{x})\|^2}{\partial x_k} = \frac{1}{2} \frac{\partial}{\partial x_k} \sum_{i=1}^n F_i(\boldsymbol{x})^2 = \sum_{i=1}^n \frac{\partial F_i(\boldsymbol{x})}{\partial x_k} F_i(\boldsymbol{x})$$

this can be written as $\nabla f(x) = \mathbf{F}(x)^T \nabla \mathbf{F}(x)$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes	

Is d_k a descent direction?

Proof: Now we check $\nabla f(x)d$:

$$\nabla f(\boldsymbol{x})\boldsymbol{d} = \mathbf{F}(\boldsymbol{x})^T \nabla \mathbf{F}(\boldsymbol{x}) \boldsymbol{d}$$

$$= -\mathbf{F}(\boldsymbol{x})^T \nabla \mathbf{F}(\boldsymbol{x}) \nabla \mathbf{F}(\boldsymbol{x})^{-1} \mathbf{F}(\boldsymbol{x})$$

$$= -\mathbf{F}(\boldsymbol{x})^T \mathbf{F}(\boldsymbol{x})$$

$$= -\|\mathbf{F}(\boldsymbol{x})\|^2 < 0$$

 $\hfill\Box$ This lemma prove that Newton direction is a descent direction.





- - - -

Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes

Angle between d_k and $\nabla f(\boldsymbol{x}_k)$

Is the angle bounded from $\pi/2$?

Let θ_k the angle between $\nabla f(x_k)$ and d_k , then we have

$$\cos \theta_k = -\frac{\nabla f(\boldsymbol{x}_k) \boldsymbol{d}_k}{\|\mathbf{F}(\boldsymbol{x}_k)\| \|\nabla \mathbf{F}(\boldsymbol{x}_k)^{-1} \mathbf{F}(\boldsymbol{x}_k)\|}$$

$$= \frac{\|\mathbf{F}(\boldsymbol{x}_k)\|}{\|\nabla \mathbf{F}(\boldsymbol{x}_k)^{-1} \mathbf{F}(\boldsymbol{x}_k)\|}$$

$$\geq \frac{\|\mathbf{F}(\boldsymbol{x}_k)\|}{\|\nabla \mathbf{F}(\boldsymbol{x}_k)^{-1}\| \|\mathbf{F}(\boldsymbol{x}_k)\|}$$

$$\geq \|\nabla \mathbf{F}(\boldsymbol{x}_k)^{-1}\|^{-1}$$

so that, if for example $\|\nabla \mathbf{F}(x)^{-1}\|$ is bounded from below then the angle θ_k is strictly less then $\pi/2$ radiants. By the Zoutendijk theorem then the globalized Newton scheme is globally convergent.



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Algorithm (The globalized Newton method)

```
k \leftarrow 0; x \text{ assigned}; \\ f \leftarrow \mathbf{F}(x); \\ \textbf{while} \ \|f\| > \epsilon \ \textbf{do} \\ - \text{Evaluate search direction} \\ \text{Solve} \quad \nabla \mathbf{F}(x)d + \mathbf{F}(x) = \mathbf{0}; \\ - \text{Evaluate dumping factor } \lambda \\ \lambda \approx \arg\min_{\alpha > 0} \|\mathbf{F}(x + \alpha d_k)\|^2 \qquad \text{by line-search}; \\ - \text{perform step} \\ x \leftarrow x + \lambda d; \\ f \leftarrow \mathbf{F}(x); \\ k \leftarrow k + 1; \\ \textbf{end while}
```





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

- 1 The Newton Raphson
- 2 The Frobenius matrix norm
- 3 The Broyden method
- 4 The dumped Broyden method
- 5 Stopping criteria and q-order estimation





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Definition

The Frobenius norm $\left\|\cdot\right\|_F$ of a matrix $A\in\mathbb{R}^{n\times m}$ is defined as follows:

$$\|A\|_F = \left(\sum_{i=1}^n \sum_{j=1}^m A_{ij}^2\right)^{1/2}$$

is a matrix norm, i.e. it satisfy:

- $\blacksquare \ \| \boldsymbol{A} \|_F \geq 0 \text{ and } \| \boldsymbol{A} \|_F = 0 \Longleftrightarrow \boldsymbol{A} = \mathbf{0};$
- $2 \|\lambda A\|_F = |\lambda| \|A\|_F;$
- $\|A + B\|_F \le \|A\|_F + \|B\|_F;$
- $\|AB\|_F \leq \|A\|_F \|B\|_F;$

The Frobenius norm is the length of the vector A if we consider A as a vector in \mathbb{R}^{n^2} .





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

The first two point of the Frobenius norm $\|\cdot\|_F$ are trivial, to prove point 3 and 4 we need two classical inequality:

Cauchy-Schwartz inequality

$$\sum_{i=1}^{n} a_i b_i \le \left(\sum_{i=1}^{n} a_i^2\right)^{1/2} \left(\sum_{i=1}^{n} b_i^2\right)^{1/2}$$

The inequality is strict unless $a_i = \lambda b_i$ for $i = 1, 2, \dots, n$.

Triangular inequality

$$\left(\sum_{i=1}^{n} (a_i + b_i)^2\right)^{1/2} \le \left(\sum_{i=1}^{n} a_i^2\right)^{1/2} + \left(\sum_{i=1}^{n} b_i^2\right)^{1/2}$$

The inequality is strict unless $a_i = \lambda b_i$ for $i = 1, 2, \dots, n$.



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

NOTES		

Proof of $\|A+B\|_F \leq \|A\|_F + \|B\|_F$. By using triangular inequality

$$\|\mathbf{A} + \mathbf{B}\|_{F} = \left(\sum_{i,j=1}^{n} (A_{ij} + B_{ij})^{2}\right)^{1/2}$$

$$\leq \left(\sum_{i,j=1}^{n} A_{ij}^{2}\right)^{1/2} + \left(\sum_{i,j=1}^{n} B_{ij}^{2}\right)^{1/2}$$

$$= \|\mathbf{A}\|_{F} + \|\mathbf{B}\|_{F}.$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes

Proof of $\|AB\|_F \leq \|A\|_F \, \|B\|_F$. By using Cauchy–Schwartz inequality with

$$\|\mathbf{A}\mathbf{B}\|_{F} = \left(\sum_{i,j=1}^{n} \left(\sum_{k=1}^{n} A_{ik} B_{kj}\right)^{2}\right)^{1/2}$$

$$\leq \left(\sum_{i,j=1}^{n} \left(\sum_{k=1}^{n} A_{ik}^{2}\right) \left(\sum_{k'=1}^{n} B_{k'j}^{2}\right)\right)^{1/2}$$

$$= \left(\left(\sum_{i=1}^{n} \sum_{k=1}^{n} A_{ik}^{2}\right) \left(\sum_{j=1}^{n} \sum_{k'=1}^{n} B_{k'j}^{2}\right)\right)^{1/2}$$

$$= \|\mathbf{A}\|_{F} \|\mathbf{B}\|_{F}.$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes

Let $u, w \in \mathbb{R}^m$ column vector then the following equality is true:

$$\left\| \boldsymbol{u} \boldsymbol{w}^T \right\|_F \leq \left\| \boldsymbol{u} \right\|_2 \left\| \boldsymbol{w} \right\|_2$$

Proof:

$$\|uw^{T}\|_{F}^{2} = \sum_{i,j=1}^{n} u_{i}^{2} w_{j}^{2}$$

$$= \left(\sum_{i=1}^{n} u_{i}^{2}\right) \left(\sum_{j=1}^{n} w_{j}^{2}\right)$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Let $A \in \mathbb{R}^{n \times m}$ and $x \in \mathbb{R}^m$ column vector then the following inequality is true:

$$\left\| \boldsymbol{A} \boldsymbol{x} \right\|_2 \leq \left\| \boldsymbol{A} \right\|_F \left\| \boldsymbol{x} \right\|_2$$

Proof: By using Cauchy-Schwarz inequality

$$\|Ax\|_{2}^{2} = \sum_{i=1}^{n} \left(\sum_{j=1}^{m} A_{ij}x_{j}\right)^{2} \leq \sum_{i=1}^{n} \left(\sum_{j=1}^{m} A_{ij}^{2}\right) \left(\sum_{k} x_{k}^{2}\right)$$

$$= \|A\|_{F}^{2} \|x\|_{2}^{2}$$



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Let $a, b \in \mathbb{R}^n$ and $x, y \in \mathbb{R}^m$ orthonormal vector. i.e. $x^Ty = 0$ and $\|x\|_2 = \|y\|_2 = 1$, then the following equality is true

$$\left\| m{a}m{x}^T + m{b}m{y}^T
ight\|_F^2 = \left\| m{a}
ight\|_2^2 + \left\| m{b}
ight\|_2^2$$

Proof:

$$\|\boldsymbol{a}\boldsymbol{x}^{T} + \boldsymbol{b}\boldsymbol{y}^{T}\|_{F}^{2} = \sum_{i=1}^{n} \sum_{j=1}^{m} (a_{i}x_{j} + b_{i}y_{j})^{2}$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{m} (a_{i}^{2}x_{j}^{2} + b_{i}^{2}y_{j}^{2} + 2a_{i}x_{j}b_{i}y_{j})$$

$$= \|\boldsymbol{a}\|_{2}^{2} \|\boldsymbol{x}\|_{2}^{2} + \|\boldsymbol{b}\|_{2}^{2} \|\boldsymbol{y}\|_{2}^{2} + 2(\boldsymbol{a}^{T}\boldsymbol{b}) \underbrace{(\boldsymbol{x}^{T}\boldsymbol{y})}_{=0}$$





4) Q (4

Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Let $A \in \mathbb{R}^{n \times m}$ and v_1 , v_2 , . . . , $v_n \in \mathbb{R}^m$ a base of orthonormal vector for \mathbb{R}^m , then

$$\|m{A}\|_F^2 = \sum_{k=1}^n \|m{A}m{v}_k\|_2^2$$

Proof: consider a generic vector $u = \alpha_1 v_1 + \cdots + \alpha_m v_m$ and notice that

$$\left(\sum_{k=1}^{m} \boldsymbol{v}_{k} \boldsymbol{v}_{k}^{T}\right) \boldsymbol{u} = \left(\sum_{k=1}^{m} \boldsymbol{v}_{k} \boldsymbol{v}_{k}^{T}\right) \left(\sum_{j=1}^{m} \alpha_{j} \boldsymbol{v}_{j}\right) = \sum_{k=1}^{m} \sum_{j=1}^{m} \boldsymbol{v}_{k} \boldsymbol{v}_{k}^{T} \boldsymbol{v}_{j} \alpha_{j}$$

$$= \sum_{k=1}^{m} \alpha_{k} \boldsymbol{v}_{k} = \boldsymbol{u}$$





35/78

Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes			

Proof: Thus

$$oldsymbol{I} = \sum_{k=1}^m oldsymbol{v}_k oldsymbol{v}_k^T$$

Using this relation we can write

$$\left\|oldsymbol{A}
ight\|_F^2 = \left\|oldsymbol{A}oldsymbol{I}
ight\|_F^2 = \left\|oldsymbol{A}oldsymbol{I}
ight\|_F^2 = \left\|\sum_{k=1}^m oldsymbol{w}_koldsymbol{v}_k^T
ight\|_F^2 = \left\|\sum_{k=1}^m oldsymbol{w}_koldsymbol{v}_k^T
ight\|_F^2 = \left\|\sum_{k=1}^m oldsymbol{w}_koldsymbol{v}_k^T
ight\|_F^2$$

where $oldsymbol{w}_k = oldsymbol{A} oldsymbol{v}_k$. Using the previous lemma we have

$$\|m{A}\|_F^2 = \sum_{k=1}^m \|m{w}_k\|_2^2 = \sum_{k=1}^m \|m{A}m{v}_k\|_2^2$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

- 2 The Frobenius matrix norm
- 3 The Broyden method
- 4 The dumped Broyden method
- 5 Stopping criteria and q-order estimation





27/70

Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

- Newton method is a fast (q-order 2) numerical scheme to approximate the root of a function $\mathbf{F}(x)$ but needs the knowledge of the Jacobian $\nabla \mathbf{F}(x)$.
- Sometimes Jacobian is not available or too expensive to compute, in this case a numerical procedure to approximate the root which does not use derivative is mandatory.
- The Newton scheme find successively the root of the affine approximation

$$L_k(\mathbf{x}) \doteq \nabla \mathbf{F}(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \mathbf{F}(\mathbf{x}_k) = \mathbf{0}$$

lacksquare Substituting the Jacobian in the affine approximation by $oldsymbol{A}_k$

$$M_k(\mathbf{x}) \doteq \mathbf{A}_k(\mathbf{x} - \mathbf{x}_k) + \mathbf{F}(\mathbf{x}_k) = \mathbf{0}$$

and solving successively this affine model produces the family of different methods:



Enrico Bertolazzi — Non-linear problems in n variable

Notes			

Algorithm (Generic Secant iterative scheme)

Let x_0 and A_0 assigned, then for k = 0, 1, 2, ...

1 Solve for p_k :

$$M_k(\boldsymbol{p}_k + \boldsymbol{x}_k) = \boldsymbol{A}_k \boldsymbol{p}_k + \mathbf{F}(\boldsymbol{x}_k) = \mathbf{0}$$

2 Update the root approximation

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{p}_k$$

3 Update the affine model and produce A_{k+1} .





*) \ (*

Notes		

- 1 The update of $M_k \to M_{k+1}$ determine the algorithm.
- 2 A simple update is the forcing of a number of the secant relation:

$$M_{k+1}(x_{k+1-\ell}) = \mathbf{F}(x_{k+1-\ell}), \qquad \ell = 1, 2, \dots, m$$

notice that $M_{k+1}(\boldsymbol{x}_{k+1}) = \mathbf{F}(\boldsymbol{x}_{k+1})$ for all \boldsymbol{A}_{k+1} .

- 3 If $A_{k+1} \in \mathbb{R}^{n \times n}$ and m=n and $d_\ell=x_{k+1-\ell}-x_{k+1}$ are linearly independent then we have enough linear relation to determine A_{k+1} .
- Unfortunately vectors d_ℓ tends to become linearly dependent so that this approach is very ill conditioned.
- A more feasible approach uses less secant relation and other conditions to determine M_{k+1} .





Enrico Bertolazzi — Non-linear problems in n variable

Notes		

- 1 The update of $M_k \to M_{k+1}$ in Broyden scheme is the following:

 - 2 $M_{k+1}(x) M_k(x)$ is small in some sense;
- 2 The first condition imply

$$A_{k+1}(x_k - x_{k+1}) + \mathbf{F}(x_{k+1}) = \mathbf{F}(x_k)$$

which set n linear equation that do not determine the n^2 coefficients of \mathbf{A}_{k+1} .

3 The second condition become

$$M_{k+1}(x) - M_k(x) = (A_{k+1} - A_k)(x - x_k)$$

$$|||M_{k+1}(x) - M_k(x)|| \le |||A_{k+1} - A_k|| |||x - x_k||$$

where $\|\cdot\|$ is some norm. The term $\|x-x_k\|$ is not controllable, so a condition should be $\|A_{k+1}-A_k\|$ is minimum.



Enrico Bertolazzi — Non-linear problems in n variable

Notes			

Defining

$$oldsymbol{y}_k = \mathbf{F}(oldsymbol{x}_{k+1}) - \mathbf{F}(oldsymbol{x}_k), \qquad oldsymbol{s}_k = oldsymbol{x}_{k+1} - oldsymbol{x}_k$$

the Broyden scheme find the update A_{k+1} which satisfy:

- $\|m{A}_{k+1} m{A}_k\| \leq \|m{B} m{A}_k\|$ for all $m{B}$ such that $m{B}m{s}_k = m{y}_k$.
- 2 If we choose for the norm $\|\cdot\|$ the Frobenius norm $\|\cdot\|_F$

$$\left\|\boldsymbol{A}\right\|_{F} = \left(\sum_{i,j=1}^{n} A_{ij}^{2}\right)^{1/2}$$

then the problem admits a unique solution.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

INOTES		

With the Frobenius matrix norm it is possible to solve the following problem

Lemma

Let $A \in \mathbb{R}^{n \times n}$ and $s,y \in \mathbb{R}^n$ with $s \neq 0$ and $As \neq y$. Consider the set

$$\mathcal{B} = ig\{ B \in \mathbb{R}^{n imes n} \, | \, Bs = y ig\}$$

then there exists a unique matrix $B \in \mathcal{B}$ such that

$$\|A-B\|_F \leq \|A-C\|_F$$
 for all $C \in \mathcal{B}$

moreover B has the following form

$$oldsymbol{B} = oldsymbol{A} + rac{(y - As)s^T}{s^Ts}$$

i.e. B is a rank one perturbation of the matrix A.



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes

Proof: (Proof.

(1/4)) First of all notice that

$$rac{1}{s^Ts}ys^T\in\mathcal{B} \qquad iggl[rac{1}{s^Ts}ys^Tiggr]s=y$$

so that set ${\cal B}$ is not empty. Next we reformulate the problem as a constrained minimum problem:

$$\underset{\boldsymbol{B} \in \mathbb{R}^{n \times n}}{\operatorname{arg \, min}} \quad \frac{1}{2} \sum_{i,j=1}^{n} (A_{ij} - B_{ij})^{2} \quad \text{subject to } \boldsymbol{Bs} = \boldsymbol{y}.$$

The solution is a stationary point of the Lagrangian:

$$g(\mathbf{B}, \lambda) = \frac{1}{2} \sum_{i,j=1}^{n} (A_{ij} - B_{ij})^2 + \sum_{i=1}^{n} \lambda_i \left(\sum_{j=1}^{n} B_{ij} s_j - y_i \right)$$



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

) 4 (

Notes		

Proof: (Proof.

(2/4)) taking the gradient we have

$$\frac{\partial}{\partial B_{ij}}g(\boldsymbol{B},\boldsymbol{\lambda}) = A_{ij} - B_{ij} + \lambda_i s_j = 0$$

$$\frac{\partial}{\partial \lambda_i} g(\boldsymbol{B}, \boldsymbol{\lambda}) = \sum_{j=1}^n B_{ij} s_j - y_j = 0$$

The previous equality can be written in matrix form

$$B = A + \lambda s^T$$
 $Bs = y$

so that we can solve for λ

$$Bs = As + \lambda s^T s = y \qquad \lambda = rac{y - As}{s^T s}$$

next we prove that B is the unique minimum.



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Proof: (Proof. (3/4)) The matrix B is at minimum distance, in fact

$$\left\|oldsymbol{B}-oldsymbol{A}
ight\|_F = \left\|oldsymbol{A} + rac{(oldsymbol{y} - oldsymbol{A} oldsymbol{s}^T}{oldsymbol{s}^Toldsymbol{s}} - oldsymbol{A}
ight\|_F = \left\|rac{(oldsymbol{y} - oldsymbol{A} oldsymbol{s}) oldsymbol{s}^T}{oldsymbol{s}^Toldsymbol{s}}
ight\|_F$$

for all $C \in \mathcal{B}$ we have Cs = y so that

$$egin{align} \|B-A\|_F &= \left\|rac{(Cs-As)s^T}{s^Ts}
ight\|_F = \left\|(C-A)rac{ss^T}{s^Ts}
ight\|_F \ &\leq \|C-A\|_F \left\|rac{ss^T}{s^Ts}
ight\|_F = \|C-A\|_F \end{aligned}$$

because in general

$$\left\|\boldsymbol{u}\boldsymbol{v}^T\right\|_F = \left(\sum_{i,j=1}^n u_i^2 v_j^2\right)^{\frac{1}{2}} = \left(\sum_{i=1}^n u_i^2 \sum_{j=1}^n v_j^2\right)^{\frac{1}{2}} = \left\|\boldsymbol{u}\right\| \left\|\boldsymbol{v}\right\|$$



Enrico Bertolazzi — Non-linear problems in n variable

46/78



Proof: (Proof. (4/4)) Let B' and B'' two different minimum. Then $\frac{1}{2}(B'+B'')\in\mathcal{B}$ moreover

$$\left\|oldsymbol{A} - rac{1}{2}(oldsymbol{B}' + oldsymbol{B}'')
ight\|_F \leq rac{1}{2}\left\|oldsymbol{A} - oldsymbol{B}'
ight\|_F + rac{1}{2}\left\|oldsymbol{A} - oldsymbol{B}''
ight\|_F$$

If the inequality is strict we have a contradiction. From the Cauchy–Schwartz inequality we have an equality only when $A-B'=\lambda(A-B'')$ so that

$$B' - \lambda B'' = (1 - \lambda)A$$

and

$$B's - \lambda B''s = (1 - \lambda)As \Rightarrow (1 - \lambda)y = (1 - \lambda)As$$

due to $As \neq y$ this is true only when $\lambda = 1$, i.e. B' = B''.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Corollary

The update

$$oldsymbol{A}_{k+1} = oldsymbol{A}_k + rac{(oldsymbol{y}_k - oldsymbol{A}_k oldsymbol{s}_k) oldsymbol{s}_k^T}{oldsymbol{s}_k^T oldsymbol{s}_k}$$

satisfy the secant condition:

$$A_{k+1}s_k = y_k$$

moreover, A_{k+1} is the nearest matrix in the Frobenius norm that satisfy the secant condition.

Remark

Different the norm produce different results and in general you can loose uniqueness of the update.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Algorithm (The Broyden method)

```
k\leftarrow 0; x_0 and A_0 assigned (for example A_0=
abla \mathbf{F}(x_0)); f_0\leftarrow \mathbf{F}(x_0); while \|f_k\|>\epsilon do Solve for s_k the linear system A_ks_k+f_k=0; x_{k+1}=x_k+s_k; f_{k+1}=\mathbf{F}(x_{k+1}); y_k=f_{k+1}-f_k; Update: A_{k+1}=A_k+rac{(y_k-A_ks_k)s_k^T}{s_k^Ts_k}; k\leftarrow k+1; end while
```





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

49//8

Notice that $y_k - A_k s_k = f_{k+1} - f_k + f_k$ so that the update can be written as $A_{k+1} \leftarrow A_k + f_{k+1} s_k^T / s_k^T s_k$ and y_k can be eliminated.

Algorithm (The Broyden method (alternative version))

```
k \leftarrow 0; x and A assigned (for example A = \nabla \mathbf{F}(x)); f \leftarrow \mathbf{F}(x); while \|f\| > \epsilon do Solve for s the linear system As + f = 0; x \leftarrow x + s; f \leftarrow \mathbf{F}(x); Update: A \leftarrow A + \frac{fs^T}{s^Ts}; k \leftarrow k + 1; end while
```





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

50/78

Broyden algorithm properties

Theorem

Let $\mathbf{F}(x)$ satisfy the standard regularity conditions with $\nabla \mathbf{F}(x_{\star})$ nonsingular. Then there exists positive constants ϵ , δ such that if $\|x_0 - x_\star\| \le \epsilon$ and $\|A_0 - \nabla F(x_\star)\| \le \delta$, then the sequence $\{x_k\}$ generated by the Broyden method is well defined and converge q-superlinearly to x_{\star} , i.e.

$$\lim_{k \to \infty} \frac{\|x_{k+1} - x_k\|}{\|x_k - x_\star\|} = 0$$



C.G.Broyden, J.E.Dennis, J.J.Moré On the local and super-linear convergence of quasi-Newton methods.

J. Inst. Math. Appl, 6 222–236, 1973.





Notes		

Broyden algorithm properties

Theorem

Let $\mathbf{F}(x) = Ax - b$ where $A \in \mathbb{R}^{n \times n}$. Then the Broyden method converge in at most 2n steps.

Theorem

Let $\mathbf{F}: \mathbb{R}^n \mapsto \mathbb{R}^n$ satisfy the standard regularity conditions with $\nabla \mathbf{F}(x_{\star})$ nonsingular. Then there exists positive constants ϵ , δ such that if $\|x_0 - x_\star\| \le \epsilon$ and $\|A_0 - \nabla \mathbf{F}(x_\star)\| \le \delta$, then the sequence $\{x_k\}$ generated by the Broyden method satisfy

$$\|\boldsymbol{x}_{k+2n} - \boldsymbol{x}_{\star}\| \le C \|\boldsymbol{x}_k - \boldsymbol{x}_{\star}\|^2$$



D.M. Gay

Some convergence properties of Broyden's method. SIAM Journal of Numerical Analysis, 16 623-630, 1979.





Notes		

Reorganizing Broyden update

- lacksquare Broyden method needs to solve a linear system for $oldsymbol{A}_k$ at each step
- This can be onerous in terms of CPU cost
- lacksquare it is possible to update directly the inverse of $m{A}_k$ i.e. it is possible to update $m{H}_k = m{A}_k^{-1}$.
- lacksquare The update of $m{A}_k$ solve the problem of efficiency but do not alleviate the memory occupation
- The matrix A_k can be written as a product of simple matrix, this can save memory if the update are lesser respect to the system dimension.





E2 /70

Notes		

Sherman-Morrison formula

Sherman-Morrison formula permit to explicity write the inverse of a matrix perturbed with a rank 1 matrix

Proposition (Sherman–Morrison formula)

$$(A + uv^T)^{-1} = A^{-1} - \frac{1}{\alpha}A^{-1}uv^TA^{-1}$$

where

$$\alpha = 1 + \boldsymbol{v}^T \boldsymbol{A}^{-1} \boldsymbol{u}$$

The Sherman–Morrison formula can be checked by a direct calculation.





54/70

Notes		

Application of Sherman-Morrison formula

(1/2)

■ From the Broyden update formula

$$oldsymbol{A}_{k+1} = oldsymbol{A}_k + rac{oldsymbol{f}_{k+1} oldsymbol{s}_k^T}{oldsymbol{s}_k^T oldsymbol{s}_k}$$

■ By using Sherman–Morrison formula

$$egin{aligned} oldsymbol{A}_{k+1}^{-1} &=& oldsymbol{A}_k^{-1} - rac{1}{eta_k} oldsymbol{A}_k^{-1} oldsymbol{f}_{k+1} oldsymbol{s}_k^T oldsymbol{A}_k^{-1} \ eta_k &=& oldsymbol{s}_k^T oldsymbol{s}_k + oldsymbol{s}_k^T oldsymbol{A}_k^{-1} oldsymbol{f}_{k+1} \end{aligned}$$

lacksquare By setting $oldsymbol{H}_k = oldsymbol{A}_k^{-1}$ we have the update formula for $oldsymbol{H}_k$:

$$egin{aligned} oldsymbol{H}_{k+1} &= oldsymbol{H}_k - rac{1}{eta_k} oldsymbol{H}_k oldsymbol{f}_{k+1} oldsymbol{s}_k^T oldsymbol{H}_k \end{aligned} egin{aligned} eta_k &= oldsymbol{s}_k^T oldsymbol{s}_k + oldsymbol{s}_k^T oldsymbol{H}_k oldsymbol{f}_{k+1} \end{aligned}$$



Enrico Bertolazzi — Non-linear problems in n variable

Notes		

Application of Sherman-Morrison formula

(2/2)

■ The update formula for H_k :

$$egin{aligned} m{H}_{k+1} &= m{H}_k - rac{1}{eta_k} m{H}_k m{f}_{k+1} m{s}_k^T m{H}_k \ eta_k &= m{s}_k^T m{s}_k + m{s}_k^T m{H}_k m{f}_{k+1} \end{aligned}$$

- Can be reorganized as follows
 - $oxed{1}$ Compute $oldsymbol{z}_{k+1} = oldsymbol{H}_k oldsymbol{f}_{k+1}$;

 - 2 Compute $\beta_k = s_k^T s_k + s_k^T z_{k+1}$; 3 Compute $\boldsymbol{H}_{k+1} = \left(\boldsymbol{I} \beta_k^{-1} z_{k+1} s_k^T\right) \boldsymbol{H}_k$;





Notes			

The Broyden method with inverse updated

Algorithm (The Broyden method (updating inverse))

```
k \leftarrow 0; x_0 assigned; f_0 \leftarrow \mathbf{F}(x_0); H_0 \leftarrow I or better H_0 \leftarrow \nabla \mathbf{F}(x_0)^{-1}; while \|f_k\| > \epsilon do - perform step s_k = -H_k f_k; x_{k+1} = x_k + s_k; f_{k+1} = \mathbf{F}(x_{k+1}); - update H z_{k+1} = H_k f_{k+1}; \beta_k = s_k^T s_k + s_k^T z_{k+1}; H_{k+1} = (I - \beta_k^{-1} z_{k+1} s_k^T) H_k; k \leftarrow k+1; end while
```





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Note:	S
-------	---

- If n is very large then the storing of \mathbf{H}_k can be very expensive.
- Moreover when n is very large we hope to find a good solution with a number m of iteration with $m \ll n$
- So that instead of storing H_k we can decide to store the vectors z_k and s_k plus the scalars β_k . With this vectors and scalars we can write

$$oldsymbol{H}_k = ig(oldsymbol{I} - eta_{k-1} oldsymbol{z}_k oldsymbol{s}_{k-1}^Tig) \cdots ig(oldsymbol{I} - eta_1 oldsymbol{z}_2 oldsymbol{s}_1^Tig) ig(oldsymbol{I} - eta_0 oldsymbol{z}_1 oldsymbol{s}_0^Tig) oldsymbol{H}_0$$

- Assuming $H_0 = I$ or can be computed on the fly we must store only 2nm + m real number instead of n^2 saving a lot of memory.
- However we can do better. It is possible to eliminate z_k ad store only n m + m real numbers.





Enrico Bertolazzi — Non-linear problems in n variable

Notes		

A step of the broyden iterative scheme can be rewritten as

$$egin{aligned} oldsymbol{d}_k &= -oldsymbol{H}_k oldsymbol{f}_k \ oldsymbol{x}_{k+1} &= oldsymbol{x}_k + oldsymbol{d}_k \ oldsymbol{f}_{k+1} &= oldsymbol{F}(oldsymbol{x}_{k+1}) \ oldsymbol{z}_{k+1} &= oldsymbol{H}_k oldsymbol{f}_{k+1} \ oldsymbol{H}_{k+1} &= igg(oldsymbol{I} - rac{oldsymbol{z}_{k+1} oldsymbol{d}_k^T}{oldsymbol{d}_k^T oldsymbol{d}_k + oldsymbol{d}_k^T oldsymbol{z}_{k+1}} igg) oldsymbol{H}_k \end{aligned}$$

- 2 you can notice that z_k and d_k are similar and contains a lot of common information.
- It is possible exploring the iteration to eliminate z_k from the update formula of H_k so that we can store the whole sequence without the vectors z_k .





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes

$$egin{aligned} -m{d}_{k+1} &= m{H}_{k+1}m{f}_{k+1} = m{igg(I - rac{m{z}_{k+1}m{d}_k^T}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}}m{igg)}m{H}_km{f}_{k+1} \ &= m{igg(I - rac{m{z}_{k+1}m{d}_k^T}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}}m{igg)}m{z}_{k+1} \ &= m{z}_{k+1} - rac{m{z}_{k+1}m{d}_k^Tm{z}_{k+1}}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}} \ &= rac{m{d}_k^Tm{d}_k}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}}m{z}_{k+1} \end{aligned}$$

substituting in the update formula for $oldsymbol{H}_{k+1}$ we obtain

$$oldsymbol{H}_{k+1} \leftarrow igg(oldsymbol{I} + rac{oldsymbol{d}_{k+1} oldsymbol{d}_k^T}{oldsymbol{d}_k^T oldsymbol{d}_k}igg)oldsymbol{H}_k$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Note	∋s
------	----

Substituting into the step of the broyden iterative scheme and assuming $oldsymbol{d}_k$ known

$$egin{aligned} m{x}_{k+1} &= m{x}_k + m{d}_k \ m{f}_{k+1} &= m{F}(m{x}_{k+1}) \ m{z}_{k+1} &= m{H}_k m{f}_{k+1} \ m{d}_{k+1} &= -rac{m{d}_k^T m{d}_k}{m{d}_k^T m{d}_k + m{d}_k^T m{z}_{k+1}} m{z}_{k+1} \ m{H}_{k+1} &= igg(m{I} + rac{m{d}_{k+1} m{d}_k^T}{m{d}_k^T m{d}_k}igg) m{H}_k \end{aligned}$$

notice that x_{k+1} , f_{k+1} and z_{k+1} are not used in H_{k+1} so that only d_k and its length need to be stored.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes

Algorithm (The Broyden method with low memory usage)

```
k \leftarrow 0; x assigned; f \leftarrow \mathbf{F}(x); H_0 \leftarrow \nabla \mathbf{F}(x)^{-1}; d_0 \leftarrow -H_0 f; \ell_0 \leftarrow d_0^T d_0; while \|f\| > \epsilon do - perform step x \leftarrow x + d_k; f \leftarrow \mathbf{F}(x); - evaluate H_k f z \leftarrow H_0 f; for j = 0, 1, \ldots, k-1 do z \leftarrow z + \left[ (d_j^T z)/\ell_j \right] d_{j+1}; end for - update H_{k+1} d_{k+1} = -\left[ \ell_k/(\ell_k + d_k^T z) \right] z; \ell_{k+1} = d_{k+1}^T d_{k+1}; \ell_k \leftarrow \ell_k + 1; end while
```



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

62/78

Outline

- 1 The Newton Raphson
- 2 The Frobenius matrix norm
- 3 The Broyden method
- 4 The dumped Broyden method
- 5 Stopping criteria and q-order estimation





63/78

Notes		

Algorithm (The dumped Broyden method)

```
k \leftarrow 0; x_0 assigned; f_0 \leftarrow \mathbf{F}(x_0); H_0 \leftarrow \nabla \mathbf{F}(x_0)^{-1}; while \|f_k\| > \epsilon do - compute search direction d_k = -H_k f_k; Approximate \arg\min_{\lambda>0} \|\mathbf{F}(x_k + \lambda d_k)\|^2 by line-search; - perform step s_k = \lambda_k d_k; x_{k+1} = x_k + s_k; f_{k+1} = \mathbf{F}(x_{k+1}); y_k = f_{k+1} - f_k; - update H_{k+1} H_{k+1} = H_k + \frac{(s_k - H_k y_k) s_k^T}{s_k^T H_k y_k} H_k; k \leftarrow k+1; end while
```



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

64/78

Notice that

$$H_k y_k = H_k f_{k+1} - H_k f_k = z_{k+1} + d_k$$
, and $s_k = \lambda_k d_k$

and

$$egin{aligned} oldsymbol{H}_{k+1} &= oldsymbol{H}_k + rac{(oldsymbol{s}_k - oldsymbol{H}_k oldsymbol{y}_k) oldsymbol{s}_k^T oldsymbol{H}_k oldsymbol{y}_k}{oldsymbol{s}_k^T oldsymbol{H}_k oldsymbol{y}_k} oldsymbol{H}_k + rac{(oldsymbol{\lambda}_k oldsymbol{d}_k - oldsymbol{z}_{k+1} - oldsymbol{d}_k) oldsymbol{\lambda}_k^T oldsymbol{g}_k oldsymbol{H}_k oldsymbol{H}_k \\ &= igg(oldsymbol{I} + rac{(oldsymbol{\lambda}_k oldsymbol{d}_k - oldsymbol{z}_{k+1} - oldsymbol{d}_k) oldsymbol{d}_k^T oldsymbol{g}_k oldsymbol{H}_k oldsy$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

65/78

A step of the broyden iterative scheme can be rewritten as

$$egin{aligned} oldsymbol{d}_k &= -oldsymbol{H}_k oldsymbol{f}_k \ oldsymbol{x}_{k+1} &= oldsymbol{x}_k + \lambda_k oldsymbol{d}_k \ oldsymbol{f}_{k+1} &= oldsymbol{F}(oldsymbol{x}_{k+1}) \ oldsymbol{z}_{k+1} &= oldsymbol{H}_k oldsymbol{f}_{k+1} \ oldsymbol{H}_{k+1} &= igg(oldsymbol{I} - rac{(oldsymbol{z}_{k+1} + (1 - \lambda_k) oldsymbol{d}_k) oldsymbol{d}_k^T}{oldsymbol{d}_k^T oldsymbol{d}_k + oldsymbol{d}_k^T oldsymbol{z}_{k+1}} igg) oldsymbol{H}_k \end{aligned}$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes

$$egin{aligned} -m{d}_{k+1} &= m{H}_{k+1}m{f}_{k+1} \ &= \left(m{I} - rac{(m{z}_{k+1} + (1-\lambda_k)m{d}_k)m{d}_k^T}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}}
ight)m{H}_km{f}_{k+1} \ &= \left(m{I} - rac{(m{z}_{k+1} + (1-\lambda_k)m{d}_k)m{d}_k^T}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}}
ight)m{z}_{k+1} \ &= m{z}_{k+1} - rac{(m{z}_{k+1} + (1-\lambda_k)m{d}_k)m{d}_k^Tm{z}_{k+1}}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}} \ &= rac{(m{d}_k^Tm{d}_k)m{z}_{k+1} + (\lambda_k - 1)(m{d}_k^Tm{z}_{k+1})m{d}_k}{m{d}_k^Tm{d}_k + m{d}_k^Tm{z}_{k+1}} \end{aligned}$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

67/78

Solving for z_{k+1}

$$oldsymbol{z}_{k+1} = -oldsymbol{d}_{k+1} - rac{(oldsymbol{d}_k^Toldsymbol{z}_{k+1})}{oldsymbol{d}_k^Toldsymbol{d}_k}(oldsymbol{d}_{k+1} + (\lambda_k - 1)oldsymbol{d}_k)$$

and adding on both side $(1 - \lambda_k)d_k$

$$z_{k+1} + (1 - \lambda_k)d_k = -(d_{k+1} + (\lambda_k - 1)d_k)\left(1 + \frac{(d_k^T z_{k+1})}{d_k^T d_k}\right)$$
$$= -(d_{k+1} + (\lambda_k - 1)d_k)\frac{d_k^T d_k + d_k^T z_{k+1}}{d_k^T d_k}$$

and substituting in \boldsymbol{H}_{k+1} we have

$$oldsymbol{H}_{k+1} = igg(oldsymbol{I} + rac{(oldsymbol{d}_{k+1} + (\lambda_k - 1)oldsymbol{d}_k)oldsymbol{d}_k^T}{oldsymbol{d}_k^Toldsymbol{d}_k}igg)oldsymbol{H}_k$$





Enrico Bertolazzi — Non-linear problems in n variable

Notes

Substituting into the step of the broyden iterative scheme and assuming d_k known

$$egin{aligned} m{x}_{k+1} &= m{x}_k + \lambda_k m{d}_k \ m{f}_{k+1} &= m{F}(m{x}_{k+1}) \ m{z}_{k+1} &= m{H}_k m{f}_{k+1} \ m{d}_{k+1} &= -rac{(m{d}_k^T m{d}_k) m{z}_{k+1} + (\lambda_k - 1) (m{d}_k^T m{z}_{k+1}) m{d}_k}{m{d}_k^T m{d}_k + m{d}_k^T m{z}_{k+1}} \ m{H}_{k+1} &= m{igg(m{I} + rac{(m{d}_{k+1} + (\lambda_k - 1) m{d}_k) m{d}_k^T}{m{d}_k^T m{d}_k} m{m{H}}_k \ \end{pmatrix}} m{H}_k \end{aligned}$$

notice that x_{k+1} , f_{k+1} and z_{k+1} are not used in H_{k+1} so that only d_k and its length need to be stored.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

Algorithm (The dumped Broyden method)



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

70.70



Some additional reference

- C. G. Broyden
 - A Class of Methods for Solving Nonlinear Simultaneous Equations

Mathematics of Computation, 19, No. 92, pp. 577–593

- C.G. Broyden
 On the discovery of the "good Broyden" method
 Mathematical Programming, 87, Number 2, 2000
- E. Bertolazzi, F. Biral and M. Da Lio Symbolic-numeric efficient solution of optimal control problems for multibody systems Journal of Computational and Applied Mathematics, 185, 2006





71/78

Notes		

Outline

- The Newton Raphson
- 2 The Frobenius matrix norm
- 3 The Broyden method
- 4 The dumped Broyden method
- 5 Stopping criteria and q-order estimation





72/78

Notes			

q-convergent sequence stopping criteria (1/2)

- 1 Consider an iterative scheme that produce a sequence $\{x_k\}$ which converge to α with q-order p.
- 2 This means that there exists a constant C such that

$$|x_{k+1} - \alpha| \le C |x_k - \alpha|^p$$
 for $k \ge m$

If
$$\lim_{k \to \infty} \frac{|x_{k+1} - \alpha|}{|x_k - \alpha|^p}$$
 exists and is say C we have

$$|x_{k+1} - \alpha| \approx C |x_k - \alpha|^p$$
 for large k

We can use this last expression to obtain an error estimate for the error and the values of p if unknown using the only known values.



Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes		

q-convergent sequence stopping criteria (2/2)

If $|x_{k+1} - \alpha| \leq C |x_k - \alpha|^p$ we can write:

$$|x_k - \alpha| \le |x_k - x_{k+1}| + |x_{k+1} - \alpha|$$

$$\le |x_k - x_{k+1}| + C|x_k - \alpha|^p$$

$$\downarrow$$

$$|x_k - \alpha| \le \frac{|x_k - x_{k+1}|}{1 - C|x_k - \alpha|^{p-1}}$$

If x_k is so near the solution such that $C|x_k-\alpha|^{p-1}\leq \frac{1}{2}$ then

$$|x_k - \alpha| \le 2|x_k - x_{k+1}|$$

This justify the stopping criteria

$$|x_{k+1} - x_k| \le \tau$$

Absolute tolerance

$$|x_{k+1} - x_k| \le \tau \max\{|x_k|, |x_{k+1}|\}$$
 Relative tolerance



Enrico Bertolazzi — Non-linear problems in n variable

Estimation of the q-order

- Consider an iterative scheme that produce a sequence $\{x_k\}$ which converge to α with q-order p.
- 2 If $|x_{k+1} \alpha| \approx C |x_k \alpha|^p$ then the ratio:

$$\log \frac{|x_{k+1} - \alpha|}{|x_k - \alpha|} \approx \log \frac{C |x_k - \alpha|^p}{|x_k - \alpha|} = (p - 1) \log C^{\frac{1}{p-1}} |x_k - \alpha|$$

and analogously

$$\log \frac{|x_{k+2} - \alpha|}{|x_{k+1} - \alpha|} \approx \log \frac{C^{1+p} |x_k - \alpha|^{p^2}}{C |x_k - \alpha|^p} = p(p-1) \log C^{\frac{1}{p-1}} |x_k - \alpha|$$

3 From this two ratio we can deduce p as

$$\log \frac{|x_{k+2} - \alpha|}{|x_{k+1} - \alpha|} / \log \frac{|x_{k+1} - \alpha|}{|x_k - \alpha|} \approx p$$





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable

Notes			

Estimation of the q-order

1 The ratio

$$\log \frac{|x_{k+2} - \alpha|}{|x_{k+1} - \alpha|} / \log \frac{|x_{k+1} - \alpha|}{|x_k - \alpha|} \approx p$$

uses the error which is not known.

If we are near the solution we can use the estimation $|x_k - \alpha| \approx |x_{k+1} - x_k|$ so that

$$\log \frac{|x_{k+2} - x_{k+3}|}{|x_{k+1} - x_{k+2}|} / \log \frac{|x_{k+1} - x_{k+2}|}{|x_k - x_{k+1}|} \approx p$$

so that 3 iteration are enough to estimate the q-order of a sequence.





Enrico Bertolazzi — Non-linear problems in $\,n\,$ variable



Estimation of the q-order

If the the step length is proportional to the value of f(x) as in Newton-Raphson scheme, i.e. $|x_k - \alpha| \approx M |f(x_k)|$ we can simplify the previous formula as:

$$\log \frac{|f(x_{k+2})|}{|f(x_{k+1})|} / \log \frac{|f(x_{k+1})|}{|f(x_k)|} \approx p$$

2 Such estimation are useful to check code implementation. In fact if we expect order p and we see order $r \neq p$ there is something wrong in the implementation or in the theory!





77/78

No)t	е	S
----	----	---	---

References

- J. Stoer and R. Bulirsch Introduction to numerical analysis Springer-Verlag, Texts in Applied Mathematics, 12, 2002.
- J. E. Dennis, Jr. and Robert B. Schnabel
 Numerical Methods for Unconstrained Optimization and
 Nonlinear Equations
 SIAM, Classics in Applied Mathematics, 16, 1996.
- Jorge Nocedal, and Stephen J. Wright Numerical optimization Springer, 2006





78/78

Notes		