

Non-linear problems in n variable

Lectures for PHD course on
Unconstrained Numerical Optimization

Enrico Bertolazzi

DIMS - Università di Trento

May 2008



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



The problem to solve

Problem

Given $\mathbf{F} : D \subseteq \mathbb{R}^n \mapsto \mathbb{R}^n$
Find $\mathbf{x}_* \in D$ for which $\mathbf{F}(\mathbf{x}_*) = \mathbf{0}$.

Example

Let

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

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



The Newton Raphson

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



The Newton procedure

(1/3)

- 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 know an approximation of a root $\mathbf{x}_0 \approx (1.1, -1.9)^T$.

- Setting $\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{p}$ we obtain ¹

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

if \mathbf{x}_0 is a good approximation of a root of $\mathbf{F}(\mathbf{x})$ then $\tilde{\mathcal{O}}(\|\mathbf{p}\|^2)$ is a small vector.

¹Here $\tilde{\mathcal{O}}(\mathbf{x})$ means $(\mathcal{O}(x_1), \dots, \mathcal{O}(x_n))^T$

The Newton procedure

(3/3)

- Considering

$$\mathbf{F}(\mathbf{x}_1 + \mathbf{q}) = \begin{pmatrix} -0.05576 \\ 8 \cdot 10^{-7} \end{pmatrix} + \begin{pmatrix} 2.0111 & 12.0668 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} q_1 \\ q_2 \end{pmatrix} + \tilde{\mathcal{O}}(\|\mathbf{q}\|^2)$$

- Neglecting $\tilde{\mathcal{O}}(\|\mathbf{q}\|^2)$ and solving

$$\begin{pmatrix} -0.05576 \\ 8 \cdot 10^{-7} \end{pmatrix} + \begin{pmatrix} 2.0111 & 12.0668 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} q_1 \\ q_2 \end{pmatrix} = \mathbf{0}$$

we obtain $\mathbf{q} = (-0.0055466, 0.0055458)^T$.

- Now we set $\mathbf{x}_2 = \mathbf{x}_1 + \mathbf{q} = (1.000015, -2.000015)^T$

The Newton procedure

(2/3)

- Neglecting $\tilde{\mathcal{O}}(\|\mathbf{p}\|^2)$ 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 $\mathbf{p} = (-0.094438, -0.105562)^T$.

- Now we set

$$\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{p} = \begin{pmatrix} 1.005562 \\ -2.0055612 \end{pmatrix}$$

The Newton procedure: a modern point of view

(1/2)

The previous procedure can be resumed as follows:

- Consider the following function $\mathbf{F}(\mathbf{x})$. We know an approximation of a root \mathbf{x}_0 .
- Expand by Taylor series

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

- Drop the term $\tilde{\mathcal{O}}(\|\mathbf{x} - \mathbf{x}_0\|^2)$ and solve

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

Call \mathbf{x}_1 this solution.

- Repeat 1 - 3 with $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$

The Newton procedure: a modern point of view

(2/2)

Algorithm (Newton iterative scheme)

Let \mathbf{x}_0 assigned, then for $k = 0, 1, 2, \dots$

- Solve for \mathbf{p}_k :

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

- Update

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



Proof.

From basic Calculus:

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

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

$$\begin{aligned} \mathbf{F}(\mathbf{y}) - \mathbf{F}(\mathbf{x}) - \nabla \mathbf{F}(\mathbf{x})(\mathbf{y} - \mathbf{x}) &= \\ \int_0^1 [\nabla \mathbf{F}(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla \mathbf{F}(\mathbf{x})](\mathbf{y} - \mathbf{x}) dt & \end{aligned}$$

and taking the norm

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



Standard Assumptions

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

Assumption (Standard Assumptions)

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

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

Lemma (Taylor like expansion)

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

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



Lemma (Jacobian norm control)

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

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

and

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



Proof.

(1/3).

From standard assumptions choosing $\gamma\delta \leq 2^{-1} \|\nabla\mathbf{F}(\mathbf{x}_*)\|$

$$\begin{aligned}\|\nabla\mathbf{F}(\mathbf{x})\| &\leq \|\nabla\mathbf{F}(\mathbf{x}) - \nabla\mathbf{F}(\mathbf{x}_*)\| + \|\nabla\mathbf{F}(\mathbf{x}_*)\| \\ &\leq \gamma \|\mathbf{x} - \mathbf{x}_*\| + \|\nabla\mathbf{F}(\mathbf{x}_*)\| \\ &\leq (3/2) \|\nabla\mathbf{F}(\mathbf{x}_*)\| \leq 2 \|\nabla\mathbf{F}(\mathbf{x}_*)\|\end{aligned}$$

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

$$\begin{aligned}\|\nabla\mathbf{F}(\mathbf{x}_*)\| &\leq \|\nabla\mathbf{F}(\mathbf{x}_*) - \nabla\mathbf{F}(\mathbf{x})\| + \|\nabla\mathbf{F}(\mathbf{x})\| \\ &\leq \gamma \|\mathbf{x} - \mathbf{x}_*\| + \|\nabla\mathbf{F}(\mathbf{x})\| \\ &\leq 2^{-1} \|\nabla\mathbf{F}(\mathbf{x}_*)\| + \|\nabla\mathbf{F}(\mathbf{x})\|\end{aligned}$$

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



Proof.

(2/3).

From the continuity of the determinant there exists a neighbor with $\nabla\mathbf{F}(\mathbf{x})$ non singular for all $\|\mathbf{x} - \mathbf{x}_*\| \leq \delta$.

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

and choosing δ such that $\gamma\delta \|\nabla\mathbf{F}(\mathbf{x}_*)^{-1}\| \leq 2^{-1}$ we have

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

and using this last inequality

$$\begin{aligned}\|\nabla\mathbf{F}(\mathbf{x}_*)^{-1}\| &\leq \|\nabla\mathbf{F}(\mathbf{x}_*)^{-1} - \nabla\mathbf{F}(\mathbf{x})^{-1}\| + \|\nabla\mathbf{F}(\mathbf{x})^{-1}\| \\ &\leq (3/2) \|\nabla\mathbf{F}(\mathbf{x})^{-1}\| \leq 2 \|\nabla\mathbf{F}(\mathbf{x})^{-1}\|\end{aligned}$$



Proof.

(3/3).

Using last inequality again

$$\begin{aligned}\|\nabla\mathbf{F}(\mathbf{x})^{-1}\| &\leq \|\nabla\mathbf{F}(\mathbf{x})^{-1} - \nabla\mathbf{F}(\mathbf{x}_*)^{-1}\| + \|\nabla\mathbf{F}(\mathbf{x}_*)^{-1}\| \\ &\leq 2^{-1} \|\nabla\mathbf{F}(\mathbf{x})^{-1}\| + \|\nabla\mathbf{F}(\mathbf{x}_*)^{-1}\|\end{aligned}$$

so that

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

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



Theorem (Local Convergence of Newton method)

Let $\mathbf{F}(\mathbf{x})$ satisfying standard assumptions, and \mathbf{x}_* a simple root (i.e. $\nabla\mathbf{F}(\mathbf{x}_*)$ non singular). Then, if $\|\mathbf{x}_0 - \mathbf{x}_*\| \leq \delta$ with $C\delta \leq 1$ where

$$C = \gamma \|\nabla\mathbf{F}(\mathbf{x}_*)^{-1}\|$$

then, the sequence generated by Newton method satisfies:

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



Proof.

Consider a Newton step with $\|x_k - x_*\| \leq \delta$ and

$$\begin{aligned} x_{k+1} - x_* &= x_k - x_* - \nabla F(x_k)^{-1} [F(x_k) - F(x_*)] \\ &= \nabla F(x_k)^{-1} [\nabla F(x_k)(x_k - x_*) - F(x_k) + F(x_*)] \end{aligned}$$

taking the norm and using Taylor like lemma

$$\|x_{k+1} - x_*\| \leq 2^{-1}\gamma \|x_k - x_*\|^2 \|\nabla F(x_k)^{-1}\|$$

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

$$\|x_{k+1} - x_*\| \leq \gamma \|\nabla F(x_*)^{-1}\| \delta \|x_k - x_*\|^2 \leq \|x_k - x_*\|,$$

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

Theorem (Newton-Kantorovich)

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

$$\begin{aligned} B(x_0, \rho) &= \{y \mid \|x_0 - y\| < \rho\}, \\ \alpha &= \|\nabla F(x_0)^{-1} F(x_0)\|, \end{aligned}$$

Moreover

- $\overline{B(x_0, \rho)} \subset D$;
- $\|\nabla F(x_0)^{-1} (F(x) - F(x_0))\| \leq \omega \|x - x_0\|$ for all $x \in D$;
- $\kappa := \alpha\omega \leq 2^{-1}$;

If the radius ρ is large enough, i.e.

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

Then:

Theorem (cont.)

- $F(x)$ has a zero $x_* \in \overline{B(x_0, \hat{\rho})}$;
- The open ball $B(x_0, \hat{\rho})$ does not contain any zero of $F(x)$ different from x_* ;
- 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

$$\|x_k - x_*\| \leq \frac{2\beta\lambda^{2^k}}{1 - \lambda^{2^k}}$$

where

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

Proof.

- P. Deuffhard 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.

- Newton method converge normally only when x_0 is near x_* 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**

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

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

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



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



Is \mathbf{d}_k a descent direction?

(1/2)

Lemma

The direction \mathbf{d} computed as a solution of the problem

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

is a descent direction.

Proof.

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

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

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

(cont.)



Is \mathbf{d}_k a descent direction?

(2/2)

Proof.

Now we check $\nabla f(\mathbf{x}) \mathbf{d}$:

$$\begin{aligned} \nabla f(\mathbf{x}) \mathbf{d} &= \mathbf{F}(\mathbf{x})^T \nabla \mathbf{F}(\mathbf{x}) \mathbf{d} \\ &= -\mathbf{F}(\mathbf{x})^T \nabla \mathbf{F}(\mathbf{x}) \nabla \mathbf{F}(\mathbf{x})^{-1} \mathbf{F}(\mathbf{x}) \\ &= -\mathbf{F}(\mathbf{x})^T \mathbf{F}(\mathbf{x}) \\ &= -\|\mathbf{F}(\mathbf{x})\|^2 < 0 \end{aligned}$$

□

This lemma prove that **Newton direction** is a descent direction.



Is the angle between \mathbf{d}_k and $\nabla f(\mathbf{x}_k)$ bounded from $\pi/2$?

Let θ_k the angle between $\nabla f(\mathbf{x}_k)$ and \mathbf{d}_k , then we have

$$\begin{aligned}\cos \theta_k &= -\frac{\nabla f(\mathbf{x}_k)\mathbf{d}_k}{\|\mathbf{F}(\mathbf{x}_k)\| \|\nabla \mathbf{F}(\mathbf{x}_k)^{-1}\mathbf{F}(\mathbf{x}_k)\|} \\ &= \frac{\|\mathbf{F}(\mathbf{x}_k)\|}{\|\nabla \mathbf{F}(\mathbf{x}_k)^{-1}\mathbf{F}(\mathbf{x}_k)\|} \\ &\geq \frac{\|\mathbf{F}(\mathbf{x}_k)\|}{\|\nabla \mathbf{F}(\mathbf{x}_k)^{-1}\| \|\mathbf{F}(\mathbf{x}_k)\|} \\ &\geq \|\nabla \mathbf{F}(\mathbf{x}_k)^{-1}\|^{-1}\end{aligned}$$

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



Algorithm (The globalized Newton method)

```

k ← 0; x assigned;
f ← F(x);
while ‖f‖ > ε do
  — Evaluate search direction
  Solve ∇F(x)d + F(x) = 0;
  — Evaluate dumping factor λ
  λ ≈ arg min_{α>0} ‖F(x + αd_k)‖2   by line-search;
  — perform step
  x ← x + λd;
  f ← F(x);
  k ← k + 1;
end while

```



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



The Frobenius matrix norm

Definition

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

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

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

- 1 $\|\mathbf{A}\|_F \geq 0$ and $\|\mathbf{A}\|_F = 0 \iff \mathbf{A} = \mathbf{0}$;
- 2 $\|\lambda \mathbf{A}\|_F = |\lambda| \|\mathbf{A}\|_F$;
- 3 $\|\mathbf{A} + \mathbf{B}\|_F \leq \|\mathbf{A}\|_F + \|\mathbf{B}\|_F$;
- 4 $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$;

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



The Frobenius matrix norm

(2/4)

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

Cauchy-Schwartz inequality

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

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

Triangular inequality

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

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



The Frobenius matrix norm

(3/4)

Proof of $\|\mathbf{A} + \mathbf{B}\|_F \leq \|\mathbf{A}\|_F + \|\mathbf{B}\|_F$.

By using triangular inequality

$$\begin{aligned} \|\mathbf{A} + \mathbf{B}\|_F &= \left(\sum_{i,j=1}^n (A_{ij} + B_{ij})^2 \right)^{1/2} \\ &\leq \left(\sum_{i,j=1}^n A_{ij}^2 \right)^{1/2} + \left(\sum_{i,j=1}^n B_{ij}^2 \right)^{1/2} \\ &= \|\mathbf{A}\|_F + \|\mathbf{B}\|_F. \end{aligned}$$



The Frobenius matrix norm

(4/4)

Proof of $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$.

By using Cauchy-Schwartz inequality with

$$\begin{aligned} \|\mathbf{A}\mathbf{B}\|_F &= \left(\sum_{i,j=1}^n \left(\sum_{k=1}^n A_{ik} B_{kj} \right)^2 \right)^{1/2} \\ &\leq \left(\sum_{i,j=1}^n \left(\sum_{k=1}^n A_{ik}^2 \right) \left(\sum_{k'=1}^n B_{k'j}^2 \right) \right)^{1/2} \\ &= \left(\left(\sum_{i=1}^n \sum_{k=1}^n A_{ik}^2 \right) \left(\sum_{j=1}^n \sum_{k'=1}^n B_{k'j}^2 \right) \right)^{1/2} \\ &= \|\mathbf{A}\|_F \|\mathbf{B}\|_F. \end{aligned}$$



Lemma

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

$$\|\mathbf{u}\mathbf{w}^T\|_F \leq \|\mathbf{u}\|_2 \|\mathbf{w}\|_2$$

Proof.

$$\begin{aligned} \|\mathbf{u}\mathbf{w}^T\|_F^2 &= \sum_{i,j=1}^n u_i^2 w_j^2 \\ &= \left(\sum_{i=1}^n u_i^2 \right) \left(\sum_{j=1}^n w_j^2 \right) \end{aligned}$$



Lemma

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

$$\|A\mathbf{x}\|_2 \leq \|A\|_F \|\mathbf{x}\|_2$$

Proof.

By using Cauchy-Schwarz inequality

$$\begin{aligned} \|A\mathbf{x}\|_2^2 &= \sum_{i=1}^n \left(\sum_{j=1}^m A_{ij} x_j \right)^2 \leq \sum_{i=1}^n \left(\sum_{j=1}^m A_{ij}^2 \right) \left(\sum_k x_k^2 \right) \\ &= \|A\|_F^2 \|\mathbf{x}\|_2^2 \end{aligned}$$

□

Lemma

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

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

Proof.

$$\begin{aligned} \|\mathbf{a}\mathbf{x}^T + \mathbf{b}\mathbf{y}^T\|_F^2 &= \sum_{i=1}^n \sum_{j=1}^m (a_i x_j + b_i y_j)^2 \\ &= \sum_{i=1}^n \sum_{j=1}^m (a_i^2 x_j^2 + b_i^2 y_j^2 + 2a_i x_j b_i y_j) \\ &= \|\mathbf{a}\|_2^2 \|\mathbf{x}\|_2^2 + \|\mathbf{b}\|_2^2 \|\mathbf{y}\|_2^2 + 2(\mathbf{a}^T \mathbf{b}) \underbrace{(\mathbf{x}^T \mathbf{y})}_{=0} \end{aligned}$$

□

Lemma

Let $A \in \mathbb{R}^{n \times m}$ and $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n \in \mathbb{R}^m$ a base of orthonormal vector for \mathbb{R}^m , then

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

Proof.

consider a generic vector $\mathbf{u} = \alpha_1 \mathbf{v}_1 + \dots + \alpha_m \mathbf{v}_m$ and notice that

$$\begin{aligned} \left(\sum_{k=1}^m \mathbf{v}_k \mathbf{v}_k^T \right) \mathbf{u} &= \left(\sum_{k=1}^m \mathbf{v}_k \mathbf{v}_k^T \right) \left(\sum_{j=1}^m \alpha_j \mathbf{v}_j \right) = \sum_{k=1}^m \sum_{j=1}^m \mathbf{v}_k \mathbf{v}_k^T \mathbf{v}_j \alpha_j \\ &= \sum_{k=1}^m \alpha_k \mathbf{v}_k = \mathbf{u} \end{aligned}$$

(cont.)

Proof.

Thus

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

Using this relation we can write

$$\|A\|_F^2 = \|AI\|_F^2 = \left\| A \left(\sum_{k=1}^m \mathbf{v}_k \mathbf{v}_k^T \right) \right\|_F^2 = \left\| \sum_{k=1}^m \mathbf{w}_k \mathbf{v}_k^T \right\|_F^2 =$$

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

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

□

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



The Broyden method

(1/5)

- Newton method is a **fast** (q -order 2) numerical scheme to approximate the root of a function $\mathbf{F}(\mathbf{x})$ but needs the knowledge of the Jacobian $\nabla \mathbf{F}(\mathbf{x})$.
- 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

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

- Substituting the Jacobian in the affine approximation by \mathbf{A}_k

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

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



The Broyden method

(2/5)

Algorithm (Generic Secant iterative scheme)

Let \mathbf{x}_0 and \mathbf{A}_0 assigned, then for $k = 0, 1, 2, \dots$

- 1 Solve for \mathbf{p}_k :

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

- 2 Update the root approximation

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

- 3 Update the affine model and produce \mathbf{A}_{k+1} .



The Broyden method

(3/5)

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

$$M_{k+1}(\mathbf{x}_{k+1-\ell}) = \mathbf{F}(\mathbf{x}_{k+1-\ell}), \quad \ell = 1, 2, \dots, m$$

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

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



The Broyden method

(4/5)

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

- $M_{k+1}(\mathbf{x}_k) = \mathbf{F}(\mathbf{x}_k)$;
- $M_{k+1}(\mathbf{x}) - M_k(\mathbf{x})$ is small in some sense;

- The first condition imply

$$\mathbf{A}_{k+1}(\mathbf{x}_k - \mathbf{x}_{k+1}) + \mathbf{F}(\mathbf{x}_{k+1}) = \mathbf{F}(\mathbf{x}_k)$$

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

- The second condition become

$$M_{k+1}(\mathbf{x}) - M_k(\mathbf{x}) = (\mathbf{A}_{k+1} - \mathbf{A}_k)(\mathbf{x} - \mathbf{x}_k)$$

$$\|M_{k+1}(\mathbf{x}) - M_k(\mathbf{x})\| \leq \|\mathbf{A}_{k+1} - \mathbf{A}_k\| \|\mathbf{x} - \mathbf{x}_k\|$$

where $\|\cdot\|$ is some norm. The term $\|\mathbf{x} - \mathbf{x}_k\|$ is not controllable, so a condition should be $\|\mathbf{A}_{k+1} - \mathbf{A}_k\|$ is minimum.



The Broyden method

(5/5)

- Defining

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

the Broyden scheme find the update \mathbf{A}_{k+1} which satisfy:

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

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

then the problem admits a unique solution.



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

Lemma

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

$$\mathcal{B} = \{ \mathbf{B} \in \mathbb{R}^{n \times n} \mid \mathbf{B}\mathbf{s} = \mathbf{y} \}$$

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

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

moreover \mathbf{B} has the following form

$$\mathbf{B} = \mathbf{A} + \frac{(\mathbf{y} - \mathbf{A}\mathbf{s})\mathbf{s}^T}{\mathbf{s}^T\mathbf{s}}$$

i.e. \mathbf{B} is a rank one perturbation of the matrix \mathbf{A} .



Proof.

(1/4)

First of all notice that

$$\frac{1}{\mathbf{s}^T\mathbf{s}}\mathbf{y}\mathbf{s}^T \in \mathcal{B} \quad \left[\frac{1}{\mathbf{s}^T\mathbf{s}}\mathbf{y}\mathbf{s}^T \right] \mathbf{s} = \mathbf{y}$$

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

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

The solution is a stationary point of the Lagrangian:

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



Proof.

(2/4).

taking the gradient we have

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

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

The previous equality can be written in matrix form

$$\mathbf{B} = \mathbf{A} + \boldsymbol{\lambda} \mathbf{s}^T \quad \mathbf{B} \mathbf{s} = \mathbf{y}$$

so that we can solve for $\boldsymbol{\lambda}$

$$\mathbf{B} \mathbf{s} = \mathbf{A} \mathbf{s} + \boldsymbol{\lambda} \mathbf{s}^T \mathbf{s} = \mathbf{y} \quad \boldsymbol{\lambda} = \frac{\mathbf{y} - \mathbf{A} \mathbf{s}}{\mathbf{s}^T \mathbf{s}}$$

next we prove that \mathbf{B} is the **unique minimum**.

Proof.

(3/4).

The matrix \mathbf{B} is at minimum distance, in fact

$$\|\mathbf{B} - \mathbf{A}\|_F = \left\| \mathbf{A} + \frac{(\mathbf{y} - \mathbf{A} \mathbf{s}) \mathbf{s}^T}{\mathbf{s}^T \mathbf{s}} - \mathbf{A} \right\|_F = \left\| \frac{(\mathbf{y} - \mathbf{A} \mathbf{s}) \mathbf{s}^T}{\mathbf{s}^T \mathbf{s}} \right\|_F$$

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

$$\begin{aligned} \|\mathbf{B} - \mathbf{A}\|_F &= \left\| \frac{(\mathbf{C} \mathbf{s} - \mathbf{A} \mathbf{s}) \mathbf{s}^T}{\mathbf{s}^T \mathbf{s}} \right\|_F = \left\| (\mathbf{C} - \mathbf{A}) \frac{\mathbf{s} \mathbf{s}^T}{\mathbf{s}^T \mathbf{s}} \right\|_F \\ &\leq \|\mathbf{C} - \mathbf{A}\|_F \left\| \frac{\mathbf{s} \mathbf{s}^T}{\mathbf{s}^T \mathbf{s}} \right\|_F = \|\mathbf{C} - \mathbf{A}\|_F \end{aligned}$$

because in general

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



Proof.

(4/4).

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

$$\left\| \mathbf{A} - \frac{1}{2}(\mathbf{B}' + \mathbf{B}'') \right\|_F \leq \frac{1}{2} \|\mathbf{A} - \mathbf{B}'\|_F + \frac{1}{2} \|\mathbf{A} - \mathbf{B}''\|_F$$

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

$$\mathbf{B}' - \lambda \mathbf{B}'' = (1 - \lambda) \mathbf{A}$$

and

$$\mathbf{B}' \mathbf{s} - \lambda \mathbf{B}'' \mathbf{s} = (1 - \lambda) \mathbf{A} \mathbf{s} \Rightarrow (1 - \lambda) \mathbf{y} = (1 - \lambda) \mathbf{A} \mathbf{s}$$

due to $\mathbf{A} \mathbf{s} \neq \mathbf{y}$ this is true only when $\lambda = 1$, i.e. $\mathbf{B}' = \mathbf{B}''$. \square 

Corollary

The update

$$\mathbf{A}_{k+1} = \mathbf{A}_k + \frac{(\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k) \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{s}_k}$$

satisfy the secant condition:

$$\mathbf{A}_{k+1} \mathbf{s}_k = \mathbf{y}_k$$

moreover, \mathbf{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.



The Brodyen method

(1/2)

Algorithm (The Brodyen method)

$k \leftarrow 0$; \mathbf{x}_0 and \mathbf{A}_0 assigned (for example $\mathbf{A}_0 = \nabla \mathbf{F}(\mathbf{x}_0)$);
 $\mathbf{f}_0 \leftarrow \mathbf{F}(\mathbf{x}_0)$;
while $\|\mathbf{f}_k\| > \epsilon$ **do**
 Solve for \mathbf{s}_k the linear system $\mathbf{A}_k \mathbf{s}_k + \mathbf{f}_k = \mathbf{0}$;
 $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$;
 $\mathbf{f}_{k+1} = \mathbf{F}(\mathbf{x}_{k+1})$;
 $\mathbf{y}_k = \mathbf{f}_{k+1} - \mathbf{f}_k$;
 Update: $\mathbf{A}_{k+1} = \mathbf{A}_k + \frac{(\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k) \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{s}_k}$;
 $k \leftarrow k + 1$;
end while



Brodyen algorithm properties

(1/2)

Theorem

Let $\mathbf{F}(\mathbf{x})$ satisfy the standard regularity conditions with $\nabla \mathbf{F}(\mathbf{x}_*)$ nonsingular. Then there exists positive constants ϵ , δ such that if $\|\mathbf{x}_0 - \mathbf{x}_*\| \leq \epsilon$ and $\|\mathbf{A}_0 - \nabla \mathbf{F}(\mathbf{x}_*)\| \leq \delta$, then the sequence $\{\mathbf{x}_k\}$ generated by the Brodyen method is well defined and converge q -superlinearly to \mathbf{x}_* , i.e.

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|}{\|\mathbf{x}_k - \mathbf{x}_*\|} = 0$$



C.G. Brodyen, 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.



The Brodyen method

(2/2)

Notice that $\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k = \mathbf{f}_{k+1} - \mathbf{f}_k + \mathbf{f}_k$ so that the update can be written as $\mathbf{A}_{k+1} \leftarrow \mathbf{A}_k + \mathbf{f}_{k+1} \mathbf{s}_k^T / \mathbf{s}_k^T \mathbf{s}_k$ and \mathbf{y}_k can be eliminated.

Algorithm (The Brodyen method (alternative version))

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



Brodyen algorithm properties

(2/2)

Theorem

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

Theorem

Let $\mathbf{F} : \mathbb{R}^n \mapsto \mathbb{R}^n$ satisfy the standard regularity conditions with $\nabla \mathbf{F}(\mathbf{x}_*)$ nonsingular. Then there exists positive constants ϵ , δ such that if $\|\mathbf{x}_0 - \mathbf{x}_*\| \leq \epsilon$ and $\|\mathbf{A}_0 - \nabla \mathbf{F}(\mathbf{x}_*)\| \leq \delta$, then the sequence $\{\mathbf{x}_k\}$ generated by the Brodyen method satisfy

$$\|\mathbf{x}_{k+2n} - \mathbf{x}_*\| \leq C \|\mathbf{x}_k - \mathbf{x}_*\|^2$$



D.M. Gay

Some convergence properties of Brodyen's method.

SIAM Journal of Numerical Analysis, 16 623–630, 1979.



Reorganizing Broyden update

- Broyden method needs to solve a linear system for A_k at each step
- This can be onerous in terms of CPU cost
- it is possible to update directly the inverse of A_k i.e. it is possible to update $H_k = A_k^{-1}$.
- The update of 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.



Application of Sherman-Morrison formula

(1/2)

- From the Broyden update formula

$$A_{k+1} = A_k + \frac{f_{k+1} s_k^T}{s_k^T s_k}$$

- By using Sherman-Morrison formula

$$A_{k+1}^{-1} = A_k^{-1} - \frac{1}{\beta_k} A_k^{-1} f_{k+1} s_k^T A_k^{-1}$$

$$\beta_k = s_k^T s_k + s_k^T A_k^{-1} f_{k+1}$$

- By setting $H_k = A_k^{-1}$ we have the update formula for H_k :

$$H_{k+1} = H_k - \frac{1}{\beta_k} H_k f_{k+1} s_k^T H_k$$

$$\beta_k = s_k^T s_k + s_k^T H_k f_{k+1}$$



Sherman-Morrison formula

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

Proposition (Sherman-Morrison formula)

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

where

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

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



Application of Sherman-Morrison formula

(2/2)

- The update formula for H_k :

$$H_{k+1} = H_k - \frac{1}{\beta_k} H_k f_{k+1} s_k^T H_k$$

$$\beta_k = s_k^T s_k + s_k^T H_k f_{k+1}$$

- Can be reorganized as follows

- Compute $z_{k+1} = H_k f_{k+1}$;
- Compute $\beta_k = s_k^T s_k + s_k^T z_{k+1}$;
- Compute $H_{k+1} = (I - \beta_k^{-1} z_{k+1} s_k^T) H_k$;



The Broyden method with inverse updated

Algorithm (The Broyden method (updating inverse))

```

k ← 0; x0 assigned;
f0 ← F(x0);
H0 ← I or better H0 ← ∇F(x0)-1;
while ||fk|| > ε do
  — perform step
  sk = -Hkfk;
  xk+1 = xk + sk;
  fk+1 = F(xk+1);
  — update H
  zk+1 = Hkfk+1;
  βk = skTsk + skTzk+1;
  Hk+1 = (I - βk-1zk+1skT)Hk;
  k ← k + 1;
end while

```



- If n is very large then the storing of 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 $m \lll n$
- 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

$$H_k = (I - \beta_{k-1}z_{k-1}s_{k-1}^T) \cdots (I - \beta_1z_1s_1^T)(I - \beta_0z_0s_0^T)H_0$$

- 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 $nm + m$ real numbers.

Elimination of z_k

(1/3)

- A step of the broyden iterative scheme can be rewritten as

$$\begin{aligned}
 d_k &= -H_k f_k \\
 x_{k+1} &= x_k + d_k \\
 f_{k+1} &= F(x_{k+1}) \\
 z_{k+1} &= H_k f_{k+1} \\
 H_{k+1} &= \left(I - \frac{z_{k+1}d_k^T}{d_k^T d_k + d_k^T z_{k+1}} \right) H_k
 \end{aligned}$$

- you can notice that z_k and d_k are similar and contains a lot of common information.
- It is easier 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 .

Elimination of z_k

(2/3)

$$\begin{aligned}
 -d_{k+1} &= H_{k+1}f_{k+1} = \left(I - \frac{z_{k+1}d_k^T}{d_k^T d_k + d_k^T z_{k+1}} \right) H_k f_{k+1} \\
 &= \left(I - \frac{z_{k+1}d_k^T}{d_k^T d_k + d_k^T z_{k+1}} \right) z_{k+1} \\
 &= z_{k+1} - \frac{z_{k+1}d_k^T z_{k+1}}{d_k^T d_k + d_k^T z_{k+1}} \\
 &= \frac{d_k^T d_k}{d_k^T d_k + d_k^T z_{k+1}} z_{k+1}
 \end{aligned}$$

substituting in the update formula for H_{k+1} we obtain

$$H_{k+1} \leftarrow \left(I + \frac{d_{k+1}d_k^T}{d_k^T d_k} \right) H_k$$



Elimination of z_k

(3/3)

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

$$\mathbf{x}_{k+1} = \mathbf{x}_k + d_k$$

$$\mathbf{f}_{k+1} = \mathbf{F}(\mathbf{x}_{k+1})$$

$$\mathbf{z}_{k+1} = \mathbf{H}_k \mathbf{f}_{k+1}$$

$$d_{k+1} = -\frac{d_k^T d_k}{d_k^T d_k + d_k^T z_{k+1}} z_{k+1}$$

$$\mathbf{H}_{k+1} = \left(\mathbf{I} + \frac{d_{k+1} d_k^T}{d_k^T d_k} \right) \mathbf{H}_k$$

notice that \mathbf{x}_{k+1} , \mathbf{f}_{k+1} and \mathbf{z}_{k+1} are not used in \mathbf{H}_{k+1} so that only d_k and its length need to be stored.



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



Algorithm (The Broyden method with low memory usage)

```

k ← 0; x assigned;
f ← F(x); H0 ← ∇F(x)-1; d0 ← -H0f; ℓ0 ← d0Td0;
while ‖f‖ > ε do
  — perform step
  x ← x + dk;
  f ← F(x);
  — evaluate Hkf
  z ← Hkf;
  for j = 0, 1, ..., k - 1 do
    z ← z + [(djTz)/ℓj]dj+1;
  end for
  — update Hk+1
  dk+1 = -[ℓk/(ℓk + dkTz)]z;
  ℓk+1 = dk+1Tdk+1;
  k ← k + 1;
end while

```



Algorithm (The dumped Broyden method)

```

k ← 0; x0 assigned;
f0 ← F(x0); H0 ← ∇F(x0)-1;
while ‖fk‖ > ε do
  — compute search direction
  dk = -Hkfk;
  Approximate arg minλ>0 ‖F(xk + λdk)‖2 by line-search;
  — perform step
  sk = λkdk;
  xk+1 = xk + sk;
  fk+1 = F(xk+1);
  yk = fk+1 - fk;
  — update Hk+1
  Hk+1 = Hk +  $\frac{(s_k - H_k y_k) s_k^T}{s_k^T H_k y_k} H_k$ ;
  k ← k + 1;
end while

```



Elimination of z_k

(1/5)

Notice that

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

and

$$\begin{aligned} H_{k+1} &= H_k + \frac{(s_k - H_k y_k) s_k^T}{s_k^T H_k y_k} H_k \\ &= H_k + \frac{(\lambda_k d_k - z_{k+1} - d_k) \lambda_k d_k^T}{\lambda_k d_k^T (z_{k+1} + d_k)} H_k \\ &= \left(I + \frac{(\lambda_k d_k - z_{k+1} - d_k) d_k^T}{d_k^T (z_{k+1} + d_k)} \right) H_k \\ &= \left(I - \frac{(z_{k+1} + (1 - \lambda_k) d_k) d_k^T}{d_k^T d_k + d_k^T z_{k+1}} \right) H_k \end{aligned}$$

Elimination of z_k

(2/5)

A step of the broyden iterative scheme can be rewritten as

$$\begin{aligned} d_k &= -H_k f_k \\ x_{k+1} &= x_k + \lambda_k d_k \\ f_{k+1} &= F(x_{k+1}) \\ z_{k+1} &= H_k f_{k+1} \\ H_{k+1} &= \left(I - \frac{(z_{k+1} + (1 - \lambda_k) d_k) d_k^T}{d_k^T d_k + d_k^T z_{k+1}} \right) H_k \end{aligned}$$

Elimination of z_k

(3/5)

$$\begin{aligned} -d_{k+1} &= H_{k+1} f_{k+1} \\ &= \left(I - \frac{(z_{k+1} + (1 - \lambda_k) d_k) d_k^T}{d_k^T d_k + d_k^T z_{k+1}} \right) H_k f_{k+1} \\ &= \left(I - \frac{(z_{k+1} + (1 - \lambda_k) d_k) d_k^T}{d_k^T d_k + d_k^T z_{k+1}} \right) z_{k+1} \\ &= z_{k+1} - \frac{(z_{k+1} + (1 - \lambda_k) d_k) d_k^T z_{k+1}}{d_k^T d_k + d_k^T z_{k+1}} \\ &= \frac{(d_k^T d_k) z_{k+1} + (\lambda_k - 1) (d_k^T z_{k+1}) d_k}{d_k^T d_k + d_k^T z_{k+1}} \end{aligned}$$

Elimination of z_k

(4/5)

Solving for z_{k+1}

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

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

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

and substituting in H_{k+1} we have

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



Elimination of z_k

(5/5)

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

$$x_{k+1} = x_k + \lambda_k d_k$$

$$f_{k+1} = F(x_{k+1})$$

$$z_{k+1} = H_k f_{k+1}$$

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

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

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.



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



Algorithm (The damped Broyden method)

```

k ← 0; x assigned;
f ← F(x); H0 ← ∇F(x)-1; d0 ← -H0f; ℓ0 ← d0Td0;
while ‖fk‖ > ε do
  Approximate arg minλ>0 ‖F(x + λdk)‖2 by line-search;
  — perform step
  x ← x + λkdk;
  f ← F(x);
  — evaluate Hkf
  z ← Hkf;
  for j = 0, 1, ..., k - 1 do
    z ← z + [(djTz)/ℓj](dj+1 + (λj - 1)dj);
  — update Hk+1
  dk+1 = -[ℓkz + (λk - 1)(dkTz)dk]/(ℓk + dkTz);
  ℓk+1 = dk+1Tdk+1;
  k ← k + 1;
end while

```



Outline

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



Stopping criteria for q -convergent sequences

(1/2)

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

$$|x_{k+1} - \alpha| \leq C |x_k - \alpha|^p \quad \text{for } k \geq m$$

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

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

- 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.

Stopping criteria q -convergent sequences

(2/2)

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

$$\begin{aligned} |x_k - \alpha| &\leq |x_k - x_{k+1}| + |x_{k+1} - \alpha| \\ &\leq |x_k - x_{k+1}| + C |x_k - \alpha|^p \\ &\Downarrow \end{aligned}$$

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

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

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

- This justify the stopping criteria

$$|x_{k+1} - x_k| \leq \tau \quad \text{Absolute tolerance}$$

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

Estimation of the q -order

(1/3)

- Consider an iterative scheme that produce a sequence $\{x_k\}$ which converge to α with q -order p .
- 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

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

- From this two ratio we can deduce p as

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

Estimation of the q -order

(2/3)

- The ratio

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

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

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

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

- if 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:

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

- 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!



- 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

