Non-linear problems in n variable Lectures for PHD course on Numerical optimization

Enrico Bertolazzi

DIMS - Universitá di Trento

November 21 - December 14, 2011



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





Problem

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

Example

Let

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

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

・ロト ・聞ト ・ ほト ・ ほト





- 3 The dumped Broyden method
- 4 Stopping criteria and q-order estimation



The Newton procedure

The Newton Raphson

The Newton procedure

• Consider the following map

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

we known an approximation of a root $\boldsymbol{x}_0 \approx (1.1, -1.9)^T$. • Setting $\boldsymbol{x}_1 = \boldsymbol{x}_0 + \boldsymbol{p}$ we obtain ¹

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

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

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

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

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



B ▶ < B ▶

(3/3)

The Newton procedure

• Considering

$$\mathbf{F}(\boldsymbol{x}_{1} + \boldsymbol{q}) = \begin{pmatrix} -0.05576\\8\,10^{-7} \end{pmatrix} + \begin{pmatrix} 2.0111 & 12.0668\\1 & 1 \end{pmatrix} \begin{pmatrix} q_{1}\\q_{2} \end{pmatrix} + \vec{\boldsymbol{\mathcal{O}}}(\|\boldsymbol{q}\|^{2})$$

• Neglecting $ec{\mathcal{O}}(\| \pmb{q} \|^2)$ and solving

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

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

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



(1/2)

The Newton procedure: a modern point of view

The previous procedure can be resumed as follows:

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

$${f F}({m x}) = {f F}({m x}_0) +
abla {f F}({m x}_0) ({m x} - {m x}_0) + {m ec {m O}}(\|{m x} - {m x}_0\|^2)$$

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.

9 Repeat
$$1-3$$
 with $oldsymbol{x}_1$, $oldsymbol{x}_2$, $oldsymbol{x}_3$, \ldots

(2/2)

The Newton procedure: a modern point of view

Algorithm (Newton iterative scheme)

Let x_0 assigned, then for $k = 0, 1, 2, \dots$

• Solve for p_k :

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

Opdate

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



Standard Assumptions

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

Assumption (Standard Assumptions)

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

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

Lemma (Taylor like expansion)

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

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



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 $\nabla \mathbf{F}(\boldsymbol{x})(\boldsymbol{y}-\boldsymbol{x})$ we have

$$\begin{split} \mathbf{F}(\boldsymbol{y}) - \mathbf{F}(\boldsymbol{x}) - \nabla \mathbf{F}(\boldsymbol{x})(\boldsymbol{y} - \boldsymbol{x}) = \\ \int_0^1 \big[\nabla \mathbf{F}(\boldsymbol{x} + t(\boldsymbol{y} - \boldsymbol{x})) - \nabla \mathbf{F}(\boldsymbol{x}) \big] (\boldsymbol{y} - \boldsymbol{x}) \, dt \end{split}$$

and taking the norm

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

E

<ロ> (日) (日) (日) (日) (日)

Lemma (Jacobian norm control)

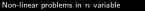
Let $\mathbf{F}(\mathbf{x})$ satisfying standard assumptions, and $\nabla \mathbf{F}(\mathbf{x}_{\star})$ non singular. Then there exists $\delta > 0$ such that for all $\|\mathbf{x} - \mathbf{x}_{\star}\| \leq \delta$ we have

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

and

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





Proof.

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

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

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



(2/3).

Proof.

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

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

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

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

and using this last inequality

$$\begin{split} \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\| &\leq \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} - \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| + \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| \\ &\leq (3/2) \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| \leq 2 \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| \end{split}$$

- 一司

(3/3).

Proof.

Using last inequality again

$$\begin{split} \left| \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| &\leq \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} - \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\| + \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\| \\ &\leq 2^{-1} \left\| \nabla \mathbf{F}(\boldsymbol{x})^{-1} \right\| + \left\| \nabla \mathbf{F}(\boldsymbol{x}_{\star})^{-1} \right\| \end{split}$$

so that

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

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



Theorem (Local Convergence of Newton method)

Let $\mathbf{F}(\boldsymbol{x})$ satisfying standard assumptions, and \boldsymbol{x}_{\star} a simple root (i.e. $\nabla \mathbf{F}(\boldsymbol{x}_{\star})$ non singular). Then, if $\|\boldsymbol{x}_{0} - \boldsymbol{x}_{\star}\| \leq \delta$ 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:

•
$$\|\boldsymbol{x}_k - \boldsymbol{x}_\star\| \le \delta$$
 for $k = 0, 1, 2, 3, ...$
• $\|\boldsymbol{x}_{k+1} - \boldsymbol{x}_\star\| \le C \|\boldsymbol{x}_k - \boldsymbol{x}_\star\|^2$ for $k = 0, 1, 2, 3, ...$
• $\lim_{k \to \infty} \boldsymbol{x}_k = \boldsymbol{x}_\star.$

• The point 2 of the theorem is the second *q*-order of convergence of Newton method.

Proof.

Consider a Newton step with $\|oldsymbol{x}_k - oldsymbol{x}_\star\| \leq \delta$ and

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

taking the norm and using Taylor like lemma

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

from Jacobian norm control lemma there exist a δ such that $2 \|\nabla \mathbf{F}(\boldsymbol{x}_k)^{-1}\| \geq \|\nabla \mathbf{F}(\boldsymbol{x}_\star)^{-1}\|$ for all $\|\boldsymbol{x}_k - \boldsymbol{x}_\star\| \leq \delta$. Reducing eventually δ such that $\gamma \delta \|\nabla \mathbf{F}(\boldsymbol{x}_\star)^{-1}\| \leq 1$ we have

$$\|oldsymbol{x}_{k+1} - oldsymbol{x}_{\star}\| \leq \gamma \left\|
abla \mathbf{F}(oldsymbol{x}_{\star})^{-1}
ight\| \delta \|oldsymbol{x}_{k} - oldsymbol{x}_{\star}\| \leq \|oldsymbol{x}_{k} - oldsymbol{x}_{\star}\| \, ,$$

So that by induction we prove point $1. \ \mbox{Point} \ 2 \ \mbox{and} \ 3 \ \mbox{follows}$ trivially.

- The problem of Newton method is that it 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

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

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

- **1** x_{\star} is a global minimum of f(x);
- 2 $\mathbf{F}(x_{\star}) = \mathbf{0}$, i.e. is a solution of the nonlinear system $\mathbf{F}(x)$.
- So that finding a global minimum of the merit function f(x) is the same of finding a solution of the nonlinear system F(x).



- 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
 - Compute the search direction by solving $\nabla \mathbf{F}(\boldsymbol{x}_k)\boldsymbol{d}_k + \mathbf{F}(\boldsymbol{x}_k) = \mathbf{0};$
 - Find an approximate solution of the problem $\alpha_k = \arg \min_{\alpha \ge 0} \|\mathbf{F}(\boldsymbol{x}_k + \alpha \boldsymbol{d}_k)\|^2;$
 - 3 Update the solution $x_{k+1} = x_k + \alpha_k d_k$.
- The previous algorithm work if the direction d_k is a descent direction.

Is d_k a descent direction?

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

$$\begin{split} \frac{\partial}{\partial x_k} \mathsf{f}(\boldsymbol{x}) &= \frac{1}{2} \frac{\partial}{\partial x_k} \, \|\mathbf{F}(\boldsymbol{x})\|^2 = \frac{1}{2} \frac{\partial}{\partial x_k} \sum_{i=1}^n F_i(\boldsymbol{x})^2 \\ &= \sum_{i=1}^n \frac{\partial F_i(\boldsymbol{x})}{\partial x_k} F_i(\boldsymbol{x}) \end{split}$$

this can be written as

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



The Newton Raphson

Is d_k a descent direction?

Now we check $-\nabla f(\boldsymbol{x}_k)\boldsymbol{d}_k$:

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

= $\mathbf{F}(\boldsymbol{x}_k)^T \nabla \mathbf{F}(\boldsymbol{x}_k) \nabla \mathbf{F}(\boldsymbol{x}_k)^{-1} \mathbf{F}(\boldsymbol{x}_k)$
= $\mathbf{F}(\boldsymbol{x}_k)^T \mathbf{F}(\boldsymbol{x}_k)$
= $\|\mathbf{F}(\boldsymbol{x}_k)\|^2 > 0$

so that Newton direction is a descent direction.

< 3 > < 3 >

Is the angle from d_k and $-\nabla f(\boldsymbol{x}_k)$ bounded from $\pi/2$?

Let $heta_k$ the angle form $abla \mathsf{f}(oldsymbol{x}_k)$ and $oldsymbol{d}_k$, then we have

$$\cos \theta_k = \frac{-\nabla f(\boldsymbol{x}_k) \boldsymbol{d}_k}{\|\nabla f(\boldsymbol{x}_k)\| \| \boldsymbol{d}_k\|} = \frac{\|\mathbf{F}(\boldsymbol{x}_k)\|^2}{\|\mathbf{F}(\boldsymbol{x}_k)^T \nabla \mathbf{F}(\boldsymbol{x}_k)\| \| \nabla \mathbf{F}(\boldsymbol{x}_k)^{-1} \mathbf{F}(\boldsymbol{x}_k)\|} \\ \geq \frac{\|\mathbf{F}(\boldsymbol{x}_k)\|^2}{\|\nabla \mathbf{F}(\boldsymbol{x}_k)\| \| \nabla \mathbf{F}(\boldsymbol{x}_k)^{-1}\| \| \mathbf{F}(\boldsymbol{x}_k)\|^2} \\ = \frac{1}{\|\nabla \mathbf{F}(\boldsymbol{x}_k)\| \| \nabla \mathbf{F}(\boldsymbol{x}_k)^{-1}\|}$$

so that, if for example $\|\nabla \mathbf{F}(\boldsymbol{x})\| \|\nabla \mathbf{F}(\boldsymbol{x})^{-1}\|$ 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.

Algorithm (The globalized Newton method)

 $k \leftarrow 0$; x assigned; $f \leftarrow \mathbf{F}(\boldsymbol{x});$ while $||f|| > \epsilon$ do — Evaluate search direction Solve $\nabla \mathbf{F}(\mathbf{x})\mathbf{d} = \mathbf{F}(\mathbf{x})$; — Evaluate dumping factor λ Approximate $\lambda = \arg \min_{\alpha > 0} \|\mathbf{F}(\boldsymbol{x} - \alpha \boldsymbol{d}_k)\|^2$ by line-search; — perform step $x \leftarrow x - \lambda d$: $f \leftarrow \mathbf{F}(\boldsymbol{x})$: $k \leftarrow k+1$: end while



- 4 伺 ト 4 ヨ ト 4 ヨ ト



The Newton Raphson



3 The dumped Broyden method

4 Stopping criteria and q-order estimation



Non-linear problems in n variable

- Newton method is a fast (q-order 2) numerical scheme to approximate the root of a function F(x) but needs the knowledge of the Jacobian ∇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

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

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

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

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

(2/5)

Algorithm (Generic Secant iterative scheme)

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

• Solve for p_k :

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

2 Update the root approximation

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

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

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

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

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

- 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 others conditions to determine M_{k+1}.



イロト イ団ト イヨト イヨト

- The way an update of $M_k \to M_{k+1}$ in Broyden scheme is the following:
 - **1** $M_{k+1}(x_k) = \mathbf{F}(x_k);$
 - 2 $M_{k+1}(x) M_k(x)$ is small in some sense;
- 2 The first condition imply

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

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

The second condition become

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

$$\left\| M_{k+1}(\boldsymbol{x}) - M_{k}(\boldsymbol{x}) \right\| \leq \left\| \boldsymbol{A}_{k+1} - \boldsymbol{A}_{k} \right\| \left\| \boldsymbol{x} - \boldsymbol{x}_{k} \right\|$$

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.

(5/5)

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 $oldsymbol{A}_{k+1}$ which satisfy:

2 If we choose for the norm $\|\!|\!|\cdot|\!|\!|$ the Frobenius norm $\|\!|\cdot|\!|_F$

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

then the problem admits a unique solution.



(1/4)

The Frobenius norm $\left\|\cdot\right\|_{F}$

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

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

3
$$\|A + B\|_F \le \|A\|_F + \|B\|_F;$$

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

The Frobenius matrix norm

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, ..., 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, ..., n.

(3/4)

 $\begin{array}{l} \mathsf{Proof of} \ \| \boldsymbol{A} + \boldsymbol{B} \|_F \leq \| \boldsymbol{A} \|_F + \| \boldsymbol{B} \|_F. \\ \mathsf{By using triangular inequality} \end{array}$

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

- ∢ 臣 ► ∢ 臣 ►

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

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

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

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

 $= \left\| \boldsymbol{A} \right\|_{F} \left\| \boldsymbol{B} \right\|_{F}.$



E

個 と く ヨ と く ヨ と

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

Lemma

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

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

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

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

moreover $oldsymbol{B}$ has the following form

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

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

(1/4).

Proof.

First of all notice that \mathcal{B} is not empty, in fact

$$rac{1}{oldsymbol{s}^Toldsymbol{s}}oldsymbol{y}oldsymbol{s}^T\in\mathcal{B}\qquadiggl[rac{1}{oldsymbol{s}^Toldsymbol{s}}oldsymbol{s}=oldsymbol{y}$$

So that the problem is not empty. Next we reformulate the problem as a constrained minimum problem:

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

The solution is a stationary point of the Lagrangian:

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

Proof.

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

$$oldsymbol{B} = oldsymbol{A} + oldsymbol{\lambda} oldsymbol{s}^T \qquad oldsymbol{B} oldsymbol{s} = oldsymbol{y}$$

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.



(3/4).

Proof.

The matrix \boldsymbol{B} is a minimum, in fact

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

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

$$egin{aligned} &\|m{B}-m{A}\|_F = \left\|rac{(m{C}m{s}-m{A}m{s})m{s}^T}{m{s}^Tm{s}}
ight\|_F = \left\|(m{C}-m{A})rac{m{s}m{s}^T}{m{s}^Tm{s}}
ight\|_F \ &\leq \|m{C}-m{A}\|_F \left\|rac{m{s}m{s}^T}{m{s}^Tm{s}}
ight\|_F = \|m{C}-m{A}\|_F \end{aligned}$$

because in general

$$\left\| \boldsymbol{u}\boldsymbol{v}^{T} \right\|_{F} = \left(\sum_{i,j=1}^{n} u_{i}^{2} v_{j}^{2}\right)^{\frac{1}{2}} = \left(\sum_{i=1}^{n} u_{i}^{2} \sum_{j=1}^{n} v_{j}^{2}\right)^{\frac{1}{2}} = \left\| \boldsymbol{u} \right\| \left\| \boldsymbol{v} \right\|$$

(4/4).

Proof.

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

$$\left\|\boldsymbol{A} - \frac{1}{2}(\boldsymbol{B}' + \boldsymbol{B}'')\right\|_{F} \le \frac{1}{2} \left\|\boldsymbol{A} - \boldsymbol{B}'\right\|_{F} + \frac{1}{2} \left\|\boldsymbol{A} - \boldsymbol{B}''\right\|_{F}$$

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

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

and

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

but this is true only when $\lambda = 1$, i.e. B' = B''.



イロト イ団ト イヨト イヨト

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: $A_{k+1}s_k = y_k$ and A_{k+1} is the nearest matrix in the Frobenius norm that satisfy the secant condition.

Changing the norm we can have different results and in general you can loose uniqueness of the update.

(1/2)

Algorithm (The Broyden method)

 $k \leftarrow 0$; x_0 and A_0 assigned; $f_0 \leftarrow \mathbf{F}(\boldsymbol{x}_0);$ while $||f_k|| > \epsilon$ do Solve for s_k the linear system $A_k s_k + f_k = 0$; $x_{k+1} \leftarrow x_k + s_k;$ $f_{k+1} \leftarrow \mathbf{F}(x_{k+1});$ $y_k \leftarrow f_{k+1} - f_k;$ Update: $oldsymbol{A}_{k+1} \leftarrow oldsymbol{A}_k + rac{(oldsymbol{y}_k - oldsymbol{A}_k s_k) oldsymbol{s}_k^T}{oldsymbol{s}_k^T oldsymbol{s}_k}$; $k \leftarrow k+1$: end while



40 / 68

- 4 伺 ト 4 ヨ ト 4 ヨ ト

The Broyden method

(2/2)

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; $f \leftarrow \mathbf{F}(x)$; 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

- 4 3 6 4 3 6

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 $\|\mathbf{x}_0 - \mathbf{x}_*\| \le \epsilon$ and $\|\mathbf{A}_0 - \nabla \mathbf{F}(\mathbf{x}_*)\| \le \delta$, then the sequence $\{\mathbf{x}_k\}$ generated by the Broyden method is well defined and converge q-superlinearly to \mathbf{x}_* , i.e.

$$\lim_{k \to \infty} \frac{\|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\|}{\|\boldsymbol{x}_k - \boldsymbol{x}_\star\|} = 0$$

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.



42 / 68

Broyden algorithm properties

(2/2)

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 $\nabla \mathbf{F}(\mathbf{x}_{\star})$ nonsingular. Then there exists positive constants ϵ , δ such that if $\|\mathbf{x}_0 - \mathbf{x}_{\star}\| \leq \epsilon$ and $\|\mathbf{A}_0 - \nabla \mathbf{F}(\mathbf{x}_{\star})\| \leq \delta$, then the sequence $\{\mathbf{x}_k\}$ generated by the Broyden method satisfy

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

D.M.Gay

Some convergence properties of Broyden's method. SIAM J. Numer. Anal., **16** 623–630, 1979.



Reorganizing Broyden update

- Broyden method needs to solve a linear system for $oldsymbol{A}_k$ 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 $oldsymbol{A}_k$ 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.

Sherman-Morrison formula

Sherman-Morrison formula permit to explicit write the inverse of a matrix changed with a rank $1\ {\rm perturbation}$

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.



The Broyden method

Application of Sherman-Morrison formula

• From the Broyden update formula

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

• By using Sherman–Morrison formula

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

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

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

(1/2)

イロン イ理 とく ヨン ト ヨン・

Application of Sherman-Morrison formula

• The update formula for H_k :

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

• Can be reorganized as follows

$$\begin{array}{l} \bullet \quad \text{Compute } \boldsymbol{z}_{k+1} = \boldsymbol{H}_k \boldsymbol{f}_{k+1};\\ \bullet \quad \text{Compute } \beta_k = \boldsymbol{s}_k^T \boldsymbol{s}_k + \boldsymbol{s}_k^T \boldsymbol{z}_{k+1};\\ \bullet \quad \text{Compute } \boldsymbol{H}_{k+1} = (\boldsymbol{I} - \beta_k^{-1} \boldsymbol{z}_{k+1} \boldsymbol{s}_k^T) \boldsymbol{H}_k; \end{array}$$

.

(2/2)

The Broyden method with inverse updated

Algorithm (The Broyden method (updating inverse))

$$\begin{split} k \leftarrow 0; x_0 \text{ assigned}; \\ f_0 \leftarrow \mathbf{F}(x_0); \\ \mathbf{H}_0 \leftarrow \mathbf{I} \text{ or better } \mathbf{H}_0 \leftarrow \nabla \mathbf{F}(x_0)^{-1}; \\ \text{while } \|f_k\| > \epsilon \text{ do} \\ \hline - \text{ perform step} \\ s_k \leftarrow -\mathbf{H}_k f_k; \\ x_{k+1} \leftarrow x_k + s_k; \\ f_{k+1} \leftarrow \mathbf{F}(x_{k+1}); \\ \hline - \text{ update } \mathbf{H} \\ \mathbf{z}_{k+1} \leftarrow \mathbf{H}_k f_{k+1}; \\ \beta_k \leftarrow s_k^T s_k + s_k^T \mathbf{z}_{k+1}; \\ \mathbf{H}_{k+1} \leftarrow (\mathbf{I} - \beta_k^{-1} \mathbf{z}_{k+1} s_k^T) \mathbf{H}_k; \\ k \leftarrow k+1; \\ \text{end while} \end{split}$$



< (¹¹) ▶

- 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

$$\boldsymbol{H}_{k} = \left(\boldsymbol{I} - \beta_{k-1}\boldsymbol{z}_{k}\boldsymbol{s}_{k-1}^{T}\right)\cdots\left(\boldsymbol{I} - \beta_{1}\boldsymbol{z}_{2}\boldsymbol{s}_{1}^{T}\right)\left(\boldsymbol{I} - \beta_{0}\boldsymbol{z}_{1}\boldsymbol{s}_{0}^{T}\right)\boldsymbol{H}_{0}$$

- 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.

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

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

- 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.

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

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

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

2

< Ξ > < Ξ >

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

 $egin{aligned} oldsymbol{x}_{k+1} &\leftarrow oldsymbol{x}_k - oldsymbol{d}_k \ oldsymbol{f}_{k+1} &\leftarrow oldsymbol{F}(oldsymbol{x}_{k+1}) \ oldsymbol{z}_{k+1} &\leftarrow oldsymbol{H}_k oldsymbol{f}_k^T oldsymbol{d}_k \ oldsymbol{d}_k^T oldsymbol{d}_k - oldsymbol{d}_k^T oldsymbol{z}_{k+1} \ oldsymbol{e}_{k+1} &\leftarrow oldsymbol{\left(oldsymbol{I} + oldsymbol{d}_k^T oldsymbol{d}_k - oldsymbol{d}_k^T oldsymbol{z}_{k+1} \ oldsymbol{H}_{k+1} &\leftarrow oldsymbol{\left(oldsymbol{I} + oldsymbol{d}_k^T oldsymbol{d}_k - oldsymbol{d}_k^T oldsymbol{z}_{k+1} \ oldsymbol{z}_{k+1} &\leftarrow oldsymbol{\left(oldsymbol{I} + oldsymbol{d}_k^T oldsymbol{d}_k - oldsymbol{d}_k^T oldsymbol{z}_{k+1} \ oldsymbol{H}_{k+1} &\leftarrow oldsymbol{\left(oldsymbol{I} + oldsymbol{d}_k^T oldsymbol{d}_k - oldsymbol{d}_k^T oldsymbol{d}_k \ oldsymbol{H}_k \ oldsymbol{H}_{k+1} &\leftarrow oldsymbol{\left(oldsymbol{I} + oldsymbol{d}_k^T oldsymbol{d}_k - oldsymbol{d}_k^T oldsymbol{d}_k \ oldsymbol{d}_k \ oldsymbol{d}_k \ oldsymbol{d}_k \ oldsymbol{H}_k \ oldsymbol{$

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.



(3/3)

・ロト ・聞ト ・ ヨト ・ ヨト

Algorithm (The Broyden method (low memory usage))

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



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



Non-linear problems in n variable

Algorithm (The dumped Broyden method)

```
k \leftarrow 0; \boldsymbol{x}_0 assigned;
f_0 \leftarrow \mathbf{F}(\boldsymbol{x}_0); \ \boldsymbol{H}_0 \leftarrow \nabla \mathbf{F}(\boldsymbol{x}_0)^{-1};
while ||f_k|| > \epsilon do
     - compute search direction
    d_{l} \leftarrow H_{l} f_{l}
    Approximate \arg \min_{\lambda > 0} \|\mathbf{F}(\boldsymbol{x}_k - \lambda \boldsymbol{d}_k)\|^2 by line-search;
     — perform step
     s_k \leftarrow -\lambda_k d_k;
     x_{k+1} \leftarrow x_k + s_k;
     f_{k+1} \leftarrow \mathbf{F}(x_{k+1})
     y_k \leftarrow f_{k+1} - f_k;
     — update H_{k+1}
    oldsymbol{H}_{k+1} \leftarrow oldsymbol{H}_k + rac{(oldsymbol{s}_k - oldsymbol{H}_k oldsymbol{y}_k) oldsymbol{s}_k^T}{oldsymbol{s}_k^T oldsymbol{H}_k oldsymbol{y}_k} oldsymbol{H}_k;
     k \leftarrow k+1:
end while
```

Notice that

$$oldsymbol{H}_koldsymbol{y}_k = oldsymbol{H}_koldsymbol{f}_{k+1} - oldsymbol{H}_koldsymbol{f}_k = oldsymbol{z}_{k+1} - oldsymbol{d}_k, \quad ext{and} \quad oldsymbol{s}_k = -\lambda_koldsymbol{d}_k$$

and

$$egin{aligned} m{H}_{k+1} &\leftarrow m{H}_k + rac{(m{s}_k - m{H}_k m{y}_k) m{s}_k^T}{m{s}_k^T m{H}_k m{y}_k} m{H}_k \ &\leftarrow m{H}_k + rac{(-\lambda_k m{d}_k - m{z}_{k+1} + m{d}_k)(-\lambda_k m{d}_k^T)}{-\lambda_k m{d}_k^T (m{z}_{k+1} - m{d}_k)} m{H}_k \ &\leftarrow igg(m{I} + rac{(-\lambda_k m{d}_k - m{z}_{k+1} + m{d}_k) m{d}_k^T}{m{d}_k^T (m{z}_{k+1} - m{d}_k)}igg) m{H}_k \ &\leftarrow igg(m{I} + rac{(-\lambda_k m{d}_k - m{z}_{k+1} + m{d}_k) m{d}_k^T}{m{d}_k^T (m{z}_{k+1} - m{d}_k)}igg) m{H}_k \end{aligned}$$

æ

イロト イヨト イヨト イヨト

A step of the broyden iterative scheme can be rewritten as

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



(2/5)

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

æ

イロト イ団ト イヨト イヨト



Solving for z_{k+1}

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

and substituting in $oldsymbol{H}_{k+1}$ we have

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



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

$$egin{aligned} & m{x}_{k+1} \leftarrow m{x}_k - \lambda_k m{d}_k \ & m{f}_{k+1} \leftarrow m{F}(m{x}_{k+1}) \ & m{z}_{k+1} \leftarrow m{H}_k m{f}_{k+1} \ & m{d}_{k+1} \leftarrow m{d}_k^T m{d}_k) m{z}_{k+1} + (\lambda_k - 1)(m{d}_k^T m{z}_{k+1}) m{d}_k \ & m{d}_k^T m{d}_k - m{d}_k^T m{z}_{k+1}) \ & m{d}_k \ & m{H}_{k+1} \leftarrow m{\left(m{I} + m{(m{d}_{k+1} + (\lambda_k - 1)m{d}_k)m{d}_k^T \ & m{d}_k^T m{d}_k \ & m{d}_k^T m{d}_k \ & m{d}_k^T m{d}_k \ & m{d}_k \ &$$

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.



(5/5)

伺下 イヨト イヨト

Algorithm (The dumped Broyden method)

 $k \leftarrow 0$; x assigned; $f \leftarrow \mathbf{F}(\boldsymbol{x}); \ \boldsymbol{H}_0 \leftarrow \nabla \mathbf{F}(\boldsymbol{x})^{-1}; \ \boldsymbol{d}_0 \leftarrow \boldsymbol{H}_0 \boldsymbol{f}; \ \ell_0 \leftarrow \boldsymbol{d}_0^T \boldsymbol{d}_0;$ while $||f_k|| > \epsilon$ do Approximate $\arg \min_{\lambda > 0} \|\mathbf{F}(\boldsymbol{x} - \lambda \boldsymbol{d}_k)\|^2$ by line-search: — perform step $\boldsymbol{x} \leftarrow \boldsymbol{x} - \lambda_k \boldsymbol{d}_k$: $\boldsymbol{f} \leftarrow \mathbf{F}(\boldsymbol{x})$: —- evaluate $H_k f$ $z \leftarrow H_0 f$: for $j = 0, 1, \dots, k - 1$ do $z \leftarrow z + \left[(\boldsymbol{d}_{i}^{T} \boldsymbol{z}) / \ell_{i} \right] (\boldsymbol{d}_{i+1} + (\lambda_{i} - 1) \boldsymbol{d}_{i});$ end for — update H_{k+1} $\boldsymbol{d}_{k+1} \leftarrow [\ell_k \boldsymbol{z} + (\lambda_k - 1)(\boldsymbol{d}_k^T \boldsymbol{z}) \boldsymbol{d}_k] / (\ell_k - \boldsymbol{d}_k^T \boldsymbol{z});$ $\ell_{k+1} \leftarrow d_{k+1}^T d_{k+1};$ $k \leftarrow k+1$: end while

Outline

- The Newton Raphson
- 2 The Broyden method
- 3 The dumped Broyden method
- 4 Stopping criteria and *q*-order estimation



Stopping criteria for *q*-convergent sequences

- Consider an iterative scheme that produce a sequence {xk} 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. Stopping criteria and q-order estimation

Stopping criteria q-convergent sequences

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

3 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



(2/2)

Stopping criteria and q-order estimation

Estimation of the q-order

Consider an iterative scheme that produce a sequence {xk} 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

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

 $\ensuremath{\mathfrak{O}}$ From this two ratio we can deduce p as

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

Estimation of the *q*-order

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.

② 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 $q\mbox{-order}$ of a sequence.

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!



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.