

Move on Newton's method

Considerations for optimization

- Suppose we want to minimize $f: \mathbb{R}^n \rightarrow \mathbb{R}^1$

Adopt Newton's method by applying it to

$$\nabla f(x) = 0 \quad (*)$$

Note: solution of (*) may be maximizers, saddle points as well as minimizers or none of those

Later: develop procedures that make convergence to anything but a minimizer unlikely

Note: Newton Iteration

$$x_+ = x - \nabla^2 f(x)^{-1} \nabla f(x)$$

where

$$\nabla^2 = \nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Essentially always have

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$$

$\Rightarrow \nabla^2 f(x)$ symmetric

Usually, $\nabla^2 f(x)$ is also positive definite near a local minimum.

Then solving $\nabla^2 f(x) \cdot s = -\nabla f(x)$

In free-matrix case, can solve with Cholesky decomposition

- In sparse case, can use PCG
- Nice property of Newton's method on discretized differential or integral equations: Mesh independence, # of iterations required for given level of accuracy does not increase much as mesh is refined. Often seen in practice. Few theoretical results (Aegommer ref.)
- Another property of Newton's method: Newton's method enjoys a certain scale independence

Suppose we rescale F :

$$F \leftarrow \bar{F} = BF \quad , \quad B \in \mathbb{R}^{n \times n} \text{ nonsingular}$$

then

$$\begin{aligned} x_+ &= x - \bar{F}'(x)^{-1} \bar{F}(x) \\ &= x - (BF'(x))^{-1} BF(x) \\ &= x - F'(x)^{-1} F(x) \end{aligned}$$

Suppose we rescale x :

$$\hat{x} = Ax ; \quad A \in \mathbb{R}^{n \times n} \text{ nonsingular}$$

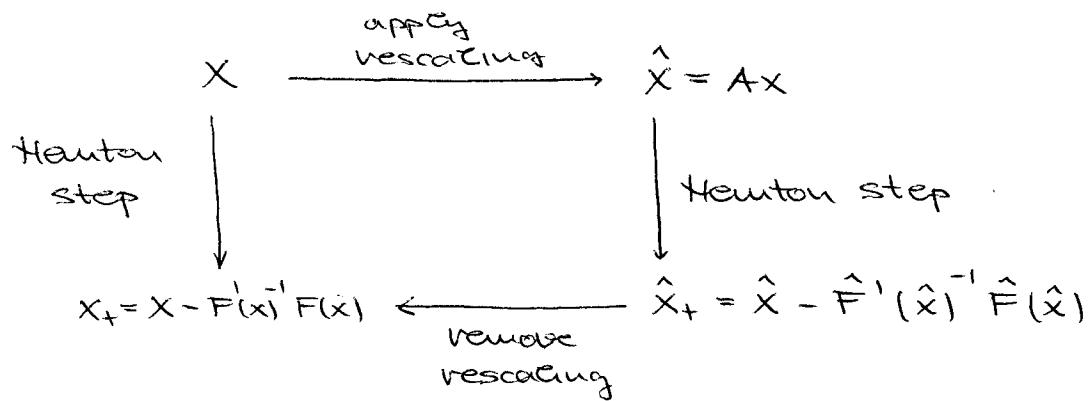
$$\begin{aligned} \text{set } \hat{F}(\hat{x}) &= F(A^{-1}\hat{x}) \\ \Rightarrow \hat{F}'(\hat{x}) &= F'(A^{-1}\hat{x}) A^{-1} \end{aligned}$$

$$\begin{aligned} \hat{x}_+ &= \hat{x} - \hat{F}'(\hat{x})^{-1} \hat{F}(\hat{x}) \\ &= \hat{x} - AF'(A^{-1}\hat{x})^{-1} \hat{F}(A^{-1}\hat{x}) \end{aligned}$$

$$Ax_+ = Ax - AF'(x)^{-1} F(x)$$

$$x_+ = x - F'(x)^{-1} F(x)$$

Schematically



- But scaling may be important to various aspects of practical implementations
- For example, rescale stopping tests

$$\|F(x)\| \leq \text{tol}_F \quad \|S\| \leq \text{tol}_X, \quad S = -F'(x)^{-1} F(x)$$

- If components of F or S are badly scaled they may not be effective as desired
- Also, bad scaling may cause problems in solving

$$F'(x) \cdot S = -F(x)$$

Also, bad scaling may cause problems with "globalization"
(Carten)

What to do about bad scaling?

Idea: Want components of F and components of x to be about the same in magnitude

May be able to find diagonal matrices D_x, D_F with diagonal entries that are typical of scales of x - and F -components

Then replace x, F by $D_x^{-1}x$ and $D_F^{-1}F$

Some variants of Newton's method

- Finite difference Newton's method
- Often evaluating $F'(x)$ is undesirable or even impossible

Idea: approximate $F'(x)$ with finite differences

Forward difference approximation:

j^{th} column of $F'(x)$

$$F'(x) = \frac{\partial F(x)}{\partial x_j} = \frac{1}{h} \left\{ F(x + h e_j) - F(x) \right\} + O(h)$$

Centrale difference approximation

$$F'(x) = \frac{\partial F(x)}{\partial x_j} = \frac{1}{2h} \left\{ F(x + h e_j) - F(x - h e_j) \right\} + O(h^2)$$

- Forward diff. is cheaper, centrale diff. is more accurate for small h
- see: approximating $F'(x)$ requires n F -evaluations (forward diff.) or $2n$ F -evaluations (centrale diff.)
- If F' is sparse, can often greatly reduce F -evaluations with Curtis - Powell - Reid "trick"

Illustrate: suppose

$$F'(x) = \begin{pmatrix} x & x & & \\ x & \ddots & & \\ & \ddots & x & \\ & & x & x \end{pmatrix} \quad \text{tridiagonal}$$

$\Leftrightarrow F'_i(x)$ depends only on x_{i-1}, x_i, x_{i+1}

see: columns 1, 4, 7, 10, ... of F' have no common nonzero components

columns 2, 5, 8, 11, ... || —

columns 3, 6, 9, 12, ... — || —

So forming

$$\frac{1}{h} \left\{ F(x + h(e_1 + e_4 + e_7 + \dots)) - F(x) \right\}$$

gives non-zero components of

1st column of F' in components 1-2

4th column ————— || ————— 3-5

7th column ————— || ————— 6-8

Farming

$$\frac{1}{h} \left\{ F(x + h(e_2 + e_5 + e_8 + \dots)) - F(x) \right\}$$

gives non-zero components of

2nd column of F' in components 1-3

5th ————— || ————— 4-6

8th ————— || ————— 7-9

Finally form

$$\frac{1}{h} \left\{ F(x + h(e_3 + e_6 + e_9 + \dots)) - F(x) \right\} \text{ to get}$$

non-zero entries of columns 3, 6, 9

So get F.D. approximation of all $F'(x)$ with just 3 F-eval.

Aside: These days automatic differentiation may be useful in evaluating $F'(x)$ when it is otherwise unavailable

What about convergence of FD Newton?

Look at convergence of very general Newton-Like Method

given x_0 and B_0

for $k = 0, 1, \dots$

$$x_{k+1} = x_k - B_k^{-1} F(x_k)$$

determine B_{k+1}

Standard Assumption: Suppose that F is continuously differentiable near x_* such that $F(x_*) = 0$ and $F'(x_*)$ is nonsingular.

Theorem: Suppose the standard assumption holds and that $\{x_k\}$ is produced by NLM with $\{B_k\}$ such that for all k

$$\textcircled{1} \quad \|B_k^{-1}\| \leq m$$

$$\textcircled{2} \quad \|I - B_k^{-1}F'(x_k)\| \leq \epsilon_{\max} < 1$$

Then for any ϵ such that $\epsilon_{\max} < \epsilon < 1$ there is a $\delta > 0$ such that if

$$\|x_0 - x_*\| \leq \delta$$

then $x_k \rightarrow x_*$ with

$$\|x_{k+1} - x_*\| \leq \epsilon \|x_k - x_*\|$$

Moreover if

$$\lim_{k \rightarrow \infty} \|I - B_k^{-1}F'(x_k)\| = 0 \quad \text{then}$$

the convergence is superlinear.

Proof idea:

$$\begin{aligned} \text{Set } e_k &= x_k - x_*, \text{ for each } k, \quad J_* = F'(x_*); \quad J(x) = F'(x) \\ x_{k+1} &= x_k - B_k^{-1}F(x_k) \\ \Rightarrow e_{k+1} &= e_k - B_k^{-1} \left\{ F(x_*) + \underbrace{J_* e_k}_{[J(x_k) + J_* - J(x)]} + \Theta(e_k) \right\} \\ &= [I - B_k^{-1}J(x_k)]e_k + B_k^{-1}[J(x_k) - J_*]e_k + B_k^{-1}\Theta(e_k) \end{aligned}$$

$$\|e_{k+1}\| \leq \underbrace{\|[I - B_k^{-1}J(x_k)]e_k\|}_{\epsilon_{\max}\|e_k\|} + \underbrace{\|B_k^{-1}[J(x_k) - J_*]e_k\|}_{M\Theta(\|e_k\|)\|e_k\|} + \underbrace{\|B_k^{-1}\Theta(e_k)\|}_{M \frac{\|\Theta(e_k)\|}{\|e_k\|}\|e_k\|} \dots (*)$$

so get

$$\|e_{k+1}\| \leq \left\{ E_{\max} + M \delta (\|e_k\|) + M \frac{\|f(e_k)\|}{\|e_k\|} \right\} \|e_k\|$$

Pick $\delta > 0$ so small that $\{ \dots \} \leq \epsilon$ whenever

$$\|x_k - x_*\| \leq \delta$$

Easy induction: if $\|x_0 - x_*\| \leq \delta$ then $\|x_k - x_*\| \leq \delta$
for all k and $\|x_{k+1} - x_*\| \leq \epsilon \|x_k - x_*\|$

from (*), see $\lim_{k \rightarrow \infty} \frac{\|e_{k+1}\|}{\|e_k\|} = \lim_{k \rightarrow \infty} \|\mathbb{I} - B_k^{-1} F'(x_k)\|$

In particular if $\lim_{k \rightarrow \infty} \|\mathbb{I} - B_k^{-1} F'(x_k)\| = 0$

then

$$\lim_{k \rightarrow \infty} \frac{\|e_{k+1}\|}{\|e_k\|} = 0 \quad (\text{superlinear convergence})$$

Apply Theorem to convergence of FD Newton

Have $B_k =$ FD approximation of $F'(x_k)$

Disregarding rounding and function evaluation errors
can make $\|\mathbb{I} - B_k^{-1} F'(x_k)\|$ arbitrarily small by
taking h sufficiently small.

Theorem: Suppose standard assumptions holds. Then
for $\epsilon > 0$ $\exists \delta > 0$ and $\eta > 0$ such that if
 $\|x_0 - x_*\| < \delta$ and $|h| < \eta$ in all FD
approximations of F' , then the FD Newton
iterates x_k converge to x_* with

$$\|x_{k+1} - x_*\| \leq \epsilon \|x_k - x_*\|$$

for each k , if $h_k =$ difference step at k^{th}
iteration $\rightarrow 0$ as $k \rightarrow \infty$ then convergence
is superlinear.

In practice, pick some small h and use it for all iterations

What should h be?

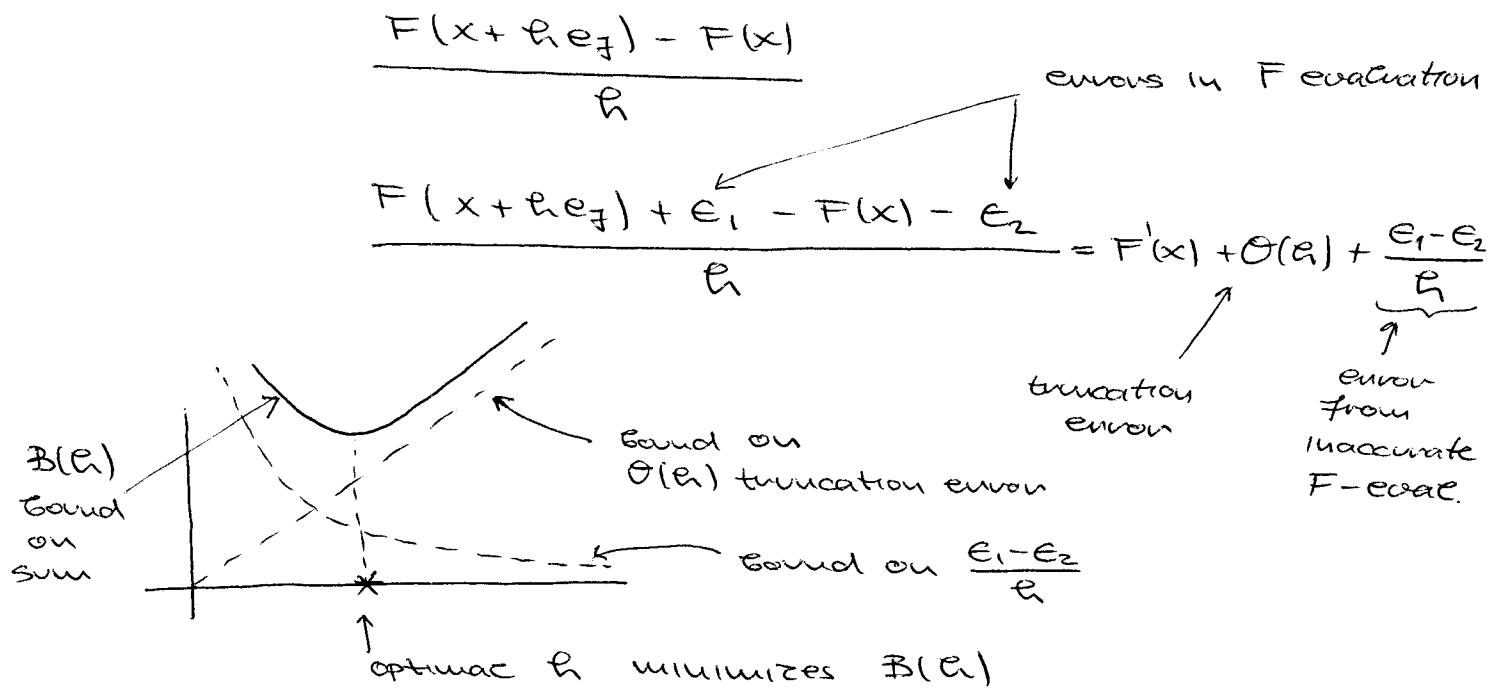
Taking h smaller makes truncation error smaller
(truncation error

$\Theta(h)$ - for forward diff.

$\Theta(h^2)$ - for central diff.)

But taking h too small makes larger the errors from F-evaluations and roundoff.

To illustrate this



For scalar function can derive a bound on error that is minimized by "optimal h "

Extending to vector-valued function is "sketchy"

I like the following (from Penrice-Wacker 1997)

ϵ_F - bound on relative error in computed F-values (function epsilon)

$$F_{\text{computed}} = F + \epsilon$$

$$\frac{\|\epsilon\|}{\|F\|} \leq \epsilon_F$$

$$h = \sqrt{(1 + \|x\|) \epsilon_F} \quad \text{for forward differences}$$

$$h = \sqrt[3]{(1 + \|x\|) \epsilon_F} \quad \text{for central differences}$$

main assumption F, F' and F'' are have the same scale

If $\epsilon_F \approx 10^{-k}$ $\approx k$ accurate digits in computed F
 expect $\approx \frac{k}{2}$ accurate digits in forward diff. approx.
 $\approx \frac{2k}{3}$ —————— || —————— in central ——————

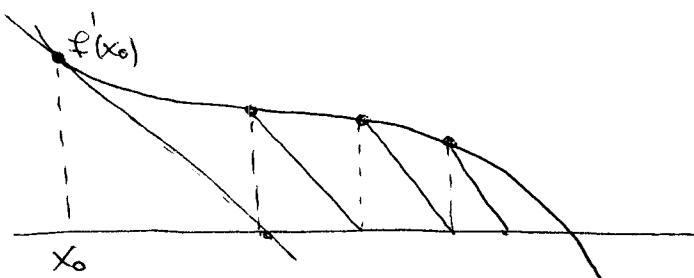
Often just take

$$h = \begin{cases} \epsilon^{1/2} & \text{for forward. diff.} \\ \epsilon^{1/2} & \text{for central diff.} \end{cases}$$

If F evaluation is unlikely to incur much error

- Usually, FD Newton converges about as rapidly as Newton (but not always)

Cord method



given initial x

evaluate $F(x)$ stop?

evaluate $J = F'(x)$ (*)

until termination do

$$x \leftarrow x - J^{-1} F(x)$$

evaluate $F(x) \rightarrow$ stop

Other's advantage:

- only one F' evaluation
- also if using direct linear algebra method
then only factor J once at $(*)$ and re-use
factors at $(**)$

Theorem: Suppose standard assumption holds. Then for
any $\epsilon > 0 \exists \delta > 0$ such that if $\|x_0 - x_*\| \leq \delta$
then Chord method iterates converge to x_* with

$$\|x_{k+1} - x_*\| \leq \epsilon \|x - x_*\|$$

for each k

proof: Follows from convergence theorem for NLM

Chord method (a.k.a. "modified Newton" method)
very successful in contexts with good initial
guess (e.g. stiff ODE solvers)

Shamanski's method

given x

evaluate $F(x)$ stop

until termination do

 evaluate $J = F'(x)$

 for $j = 1, \dots, k$ do

$x \leftarrow x - J^{-1} F(x)$

 evaluate $F(x)$ stop

- get local superlinear convergence from Theorem
for NLM