

DIRECTIONAL NEWTON METHODS IN n VARIABLES

YURI LEVIN AND ADI BEN-ISRAEL

ABSTRACT. Directional Newton methods for functions f of n variables are shown to converge, under standard assumptions, to a solution of $f(\mathbf{x}) = 0$. The rate of convergence is quadratic, for near-gradient directions, and directions along components of the gradient of f with maximal modulus. These methods are applied to solving systems of equations without inversion of the Jacobian matrix.

1. INTRODUCTION

Consider a single equation in n unknowns,

$$f(\mathbf{x}) = 0, \quad \text{or} \quad f(x_1, x_2, \dots, x_n) = 0. \quad (1)$$

Given a point \mathbf{x}^0 where f is differentiable and a direction vector \mathbf{d} , we restrict f to the line

$$L := \{\mathbf{x}^0 + t\mathbf{d} : t \in \mathbb{R}\},$$

where it is a function of one variable

$$F(t) := f(\mathbf{x}^0 + t\mathbf{d}).$$

The Newton iteration for F at $t^0 = 0$ gives the next point

$$t^1 := -\frac{F(0)}{F'(0)},$$

and the corresponding iteration for f is

$$\mathbf{x}^1 := \mathbf{x}^0 - \frac{f(\mathbf{x}^0)}{\nabla f(\mathbf{x}^0) \cdot \mathbf{d}} \mathbf{d} \quad (2)$$

since $F(0) = f(\mathbf{x}^0)$ and $F'(0)$ is the directional derivative

$$F'(0) = f'(\mathbf{x}^0, \mathbf{d}) = \nabla f(\mathbf{x}^0) \cdot \mathbf{d}.$$

Continuing in this fashion we get the iterations

$$\mathbf{x}^{k+1} := \mathbf{x}^k - \frac{f(\mathbf{x}^k)}{\nabla f(\mathbf{x}^k) \cdot \mathbf{d}^k} \mathbf{d}^k, \quad k = 0, 1, \dots \quad (3)$$

that we call a **directional Newton method**. If $n = 1$, (3) is the classical Newton method.

Date: September 14, 1999. Revised May 9, 2000.

1991 Mathematics Subject Classification. Primary 65H05, 65H10; Secondary 49M15.

Key words and phrases. Newton Method, Single equations, Systems of equations.

The first author was supported by the Center for Discrete Mathematics and Theoretical Computer science (DIMACS), Rutgers University.

For directions \mathbf{d}^k sufficiently close to the gradients $\nabla f(\mathbf{x}^k)$, Theorems 1–2 establish quadratic convergence of the method (3) under standard assumptions, namely

- the gradient of f not “too small”, see (14) and (23b), and
- the second derivative (Hessian matrix) of f not “too large”, see (8a).

A special case is the gradient method (28) studied in Corollary 1.

Another choice of \mathbf{d}^k is the unit vector $\mathbf{e}^{m(k)}$ where $m(k)$ is the index of a component of $\nabla f(\mathbf{x}^k)$ of maximal modulus

$$\left| \nabla f(\mathbf{x}^k)[m(k)] \right| := \max_{j=1,2,\dots,n} \left| \nabla f(\mathbf{x}^k)[j] \right|. \quad (4)$$

For this choice of \mathbf{d}^k , the directional Newton method (3) becomes

$$\mathbf{x}^{k+1} := \mathbf{x}^k - \frac{f(\mathbf{x}^k)}{\nabla f(\mathbf{x}^k)[m(k)]} \mathbf{e}^{m(k)}, \quad k = 0, 1, \dots \quad (5)$$

a method whose quadratic convergence is established in Theorem 3. This method is suitable for parallel implementations.

These results are applied in § 4 to general systems of m equations in n unknowns.

MAPLE programs for these and related methods can be downloaded from [6].

Notation: We use the Euclidean norm $\|\mathbf{x}\|$, and the corresponding matrix norm $\|A\|$, except in § 3 where the ∞ -norm is used for vectors and matrices, denoted by $\|\mathbf{x}\|_\infty$ and $\|A\|_\infty$ respectively.

The angle between two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ is denoted $\angle(\mathbf{u}, \mathbf{v})$ and given by

$$\angle(\mathbf{u}, \mathbf{v}) := \arccos \left(\frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \right). \quad (6)$$

2. GRADIENT AND NEAR GRADIENT METHODS

In this section we study the convergence of the directional Newton method (3), along gradient and near gradient directions. The proofs use standard arguments, see e.g. [8, Chapter 7].

Theorem 1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable, let \mathbf{x}^0 be a point where*

$$f(\mathbf{x}^0) \neq 0, \quad \nabla f(\mathbf{x}^0) \neq \mathbf{0},$$

let $\mathbf{d}^0 \in \mathbb{R}^n$ be such that $\|\mathbf{d}^0\| = 1$ and let

$$\mathbf{h}^0 := -\frac{f(\mathbf{x}^0)}{\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0} \mathbf{d}^0, \quad (7a)$$

$$\mathbf{x}^1 := \mathbf{x}^0 + \mathbf{h}^0. \quad (7b)$$

Consider the ball

$$X_0 := \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^1\| \leq \|\mathbf{h}^0\|\}$$

and assume that $f \in C^2[X_0]$ and

$$\sup_{\mathbf{x} \in X_0} \|f''(\mathbf{x})\| = M, \quad (8a)$$

$$|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|^2 \geq 2 |f(\mathbf{x}^0)| M, \quad (8b)$$

where f'' is the Hessian matrix of f and $\|f''\|$ is its matrix norm corresponding to the Euclidean norm.

Consider the sequence $\{\mathbf{x}^i : i = 1, 2, \dots\}$ defined recursively by

$$\mathbf{x}^{i+1} := \mathbf{x}^i + \mathbf{h}^i, \quad (9a)$$

$$\text{where } \mathbf{h}^i := -\frac{f(\mathbf{x}^i)}{\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i} \mathbf{d}^i, \quad (9b)$$

$\mathbf{d}^i \in \mathbb{R}^n$ is such that $\|\mathbf{d}^i\| = 1$ and

$$\angle(\mathbf{d}^{i+1}, \nabla f(\mathbf{x}^{i+1})) \leq \angle(\mathbf{d}^i, \nabla f(\mathbf{x}^i)), \quad i = 0, 1, \dots \quad (10)$$

Then:

- (a) All \mathbf{x}^i lie in X_0 .
- (b) The sequence $\{\mathbf{x}^i\}$ converges, as $i \rightarrow \infty$, to a point $\mathbf{x}^* \in X_0$ that is a zero of f .
- (c) Moreover¹, $\nabla f(\mathbf{x}^*) \neq \mathbf{0}$ unless $\|\mathbf{x}^* - \mathbf{x}^0\| = 2 \|\mathbf{h}^0\|$.
- (d) Further, for $i = 1, 2, \dots$

$$\|\mathbf{x}^{i+1} - \mathbf{x}^i\| \leq \frac{M}{2 |\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i|} \|\mathbf{x}^i - \mathbf{x}^{i-1}\|^2, \quad (11a)$$

$$\|\mathbf{x}^* - \mathbf{x}^{i+1}\| \leq \frac{M}{2 |\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i|} \|\mathbf{x}^i - \mathbf{x}^{i-1}\|^2. \quad (11b)$$

Remark. Since $\|\mathbf{d}^i\| = 1$ for all i , condition (10) is equivalent to:

$$\frac{|\nabla f(\mathbf{x}^{i+1}) \cdot \mathbf{d}^{i+1}|}{\|\nabla f(\mathbf{x}^{i+1})\|} \geq \frac{|\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i|}{\|\nabla f(\mathbf{x}^i)\|}, \quad i = 0, 1, \dots \quad (12)$$

These conditions state that the direction \mathbf{d}^i does not have to be precisely along the gradient $\nabla f(\mathbf{x}^i)$ (the most common choice). Small perturbations in the angle $\angle(\mathbf{d}^i, \nabla f(\mathbf{x}^i))$ are allowed, if they do not increase with i .

Proof.

Part 1. Proof that $\|\nabla f(\mathbf{x}^1)\| \geq \frac{1}{2} \|\nabla f(\mathbf{x}^0)\|$. (13)

We rewrite condition (8b) as

$$|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0| \geq 2 \|\mathbf{h}^0\| M. \quad (14)$$

Each component of $f''(\mathbf{x})d\mathbf{x}$ is an exact 1-form, and by [4, Theorem 6.1],

$$\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^0) = \int_{\mathbf{x}^0}^{\mathbf{x}} f''(\mathbf{x})d\mathbf{x}.$$

¹The example: $n = 1$, $f(x) = x^2$, $x_0 = 1$ illustrates the case $\|\mathbf{x}^* - \mathbf{x}^0\| = 2 \|\mathbf{h}^0\|$ and $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Therefore, by (8a),

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^0)\| \leq M \|\mathbf{x} - \mathbf{x}^0\| \quad (15)$$

for all $\mathbf{x} \in X_0$. In particular, for \mathbf{x}^1 ,

$$\begin{aligned} \|\nabla f(\mathbf{x}^1) - \nabla f(\mathbf{x}^0)\| &\leq M \|\mathbf{x}^1 - \mathbf{x}^0\| = M \|\mathbf{h}^0\| \\ &\leq \frac{1}{2} |\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|, \text{ by (14) ,} \\ &\leq \frac{1}{2} \|\nabla f(\mathbf{x}^0)\|, \text{ since } \|\mathbf{d}^0\| = 1. \quad (16) \\ \therefore \|\nabla f(\mathbf{x}^1)\| &\geq \|\nabla f(\mathbf{x}^0)\| - \|\nabla f(\mathbf{x}^1) - \nabla f(\mathbf{x}^0)\|, \\ &\geq \|\nabla f(\mathbf{x}^0)\| - \frac{1}{2} \|\nabla f(\mathbf{x}^0)\|, \text{ by (16) ,} \\ &= \frac{1}{2} \|\nabla f(\mathbf{x}^0)\|, \text{ proving (13) .} \end{aligned}$$

Part 2. Proof that

$$|f(\mathbf{x}^1)| \leq \frac{1}{2} \|\mathbf{h}^0\|^2 M, \quad (17a)$$

$$|\nabla f(\mathbf{x}^1) \cdot \mathbf{d}^1| \geq 2 \|\mathbf{h}^1\| M, \quad (17b)$$

$$\text{and } \|\mathbf{h}^1\| \leq \frac{1}{2} \|\mathbf{h}^0\|. \quad (17c)$$

First we prove that

$$f(\mathbf{x}^1) = \int_{\mathbf{x}^0}^{\mathbf{x}^1} (\mathbf{x}^1 - \mathbf{x}) f''(\mathbf{x}) d\mathbf{x}, \quad (18)$$

using integration by parts:

$$\begin{aligned} \int_{\mathbf{x}^0}^{\mathbf{x}^1} (\mathbf{x}^1 - \mathbf{x}) \cdot f''(\mathbf{x}) d\mathbf{x} &= -(\mathbf{x}^1 - \mathbf{x}^0) \cdot \nabla f(\mathbf{x}^0) + f(\mathbf{x}^1) - f(\mathbf{x}^0) \\ &= -\mathbf{h}^0 \cdot \nabla f(\mathbf{x}^0) + f(\mathbf{x}^1) - f(\mathbf{x}^0) \\ &= \frac{f(\mathbf{x}^0)}{(\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0)} (\mathbf{d}^0 \cdot \nabla f(\mathbf{x}^0)) + f(\mathbf{x}^1) - f(\mathbf{x}^0) \\ &= f(\mathbf{x}^1), \text{ by (7a) .} \end{aligned}$$

Next we prove (17a) by integrating (18) using a variable t defined by $\mathbf{x} = \mathbf{x}^0 + t\mathbf{h}^0$

$$\therefore \mathbf{x}^1 - \mathbf{x} = \mathbf{x}^1 - \mathbf{x}^0 - t\mathbf{h}^0 = \mathbf{h}^0 - t\mathbf{h}^0 = (1-t)\mathbf{h}^0, \quad d\mathbf{x} = \mathbf{h}^0 dt.$$

Thus (18) becomes

$$f(\mathbf{x}^1) = \int_0^1 (1-t) \mathbf{h}^0 \cdot f''(\mathbf{x}^0 + t\mathbf{h}^0) \mathbf{h}^0 dt.$$

Since $1 - t \geq 0$, it follows that

$$\begin{aligned} |f(\mathbf{x}^1)| &\leq \int_0^1 (1-t) \frac{|f(\mathbf{x}^0)|^2}{|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|^2} \mathbf{d}^0 \cdot f''(\mathbf{x}^0 + t\mathbf{h}^0) \mathbf{d}^0 dt, \\ &\leq \frac{M |f(\mathbf{x}^0)|^2}{2 |\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|^2} \|\mathbf{d}^0\|^2, \text{ by (8a)}, \\ &= \frac{1}{2} \|\mathbf{h}^0\|^2 M, \text{ by (7a), proving (17a)}. \end{aligned}$$

For \mathbf{h}^1 given by (9b) we have

$$\begin{aligned} \|\mathbf{h}^1\| &= \frac{|f(\mathbf{x}^1)|}{|\nabla f(\mathbf{x}^1) \cdot \mathbf{d}^1|} \leq \frac{\frac{1}{2} \|\mathbf{h}^0\|^2 M}{|\nabla f(\mathbf{x}^1) \cdot \mathbf{d}^1|}, \text{ by (17a)}, \\ &\leq \frac{\frac{1}{2} \|\mathbf{h}^0\|^2 M}{|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|} \frac{\|\nabla f(\mathbf{x}^0)\|}{\|\nabla f(\mathbf{x}^1)\|}, \text{ by (12)} \\ &\leq \frac{\|\mathbf{h}^0\|^2 M}{|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|}, \text{ by (13)}. \tag{19} \\ \therefore \frac{2M \|\mathbf{h}^1\|}{\|\nabla f(\mathbf{x}^1)\|} &\leq \frac{2M \|\mathbf{h}^1\|}{|\nabla f(\mathbf{x}^1) \cdot \mathbf{d}^1|}, \text{ since } \|\mathbf{d}^1\| = 1, \\ &\leq \frac{2M^2 \|\mathbf{h}^0\|^2}{|\nabla f(\mathbf{x}^1) \cdot \mathbf{d}^1| |\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|}, \text{ by (19)}, \\ &\leq \frac{2M^2 \|\mathbf{h}^0\|^2}{|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|^2} \frac{\|\nabla f(\mathbf{x}^0)\|}{\|\nabla f(\mathbf{x}^1)\|}, \text{ by (12)}, \\ &\leq \frac{2^2 M^2 \|\mathbf{h}^0\|^2}{|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|^2}, \text{ by (13)}, \\ &= \left(\frac{2M \|\mathbf{h}^0\|}{|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|} \right)^2 \leq 1, \text{ by (14), proving (17b)}. \end{aligned}$$

From (19) we also get

$$\frac{\|\mathbf{h}^1\|}{\|\mathbf{h}^0\|} \leq \frac{1}{2} \left(\frac{2 \|\mathbf{h}^0\| M}{|\nabla f(\mathbf{x}^0) \cdot \mathbf{d}^0|} \right) \leq \frac{1}{2}, \text{ by (14), proving (17c)}.$$

Part 3. Proof of claims (a) and (b).

It follows from (17c) that

$$\|\mathbf{x}^2 - \mathbf{x}^1\| = \|\mathbf{h}^1\| \leq \frac{1}{2} \|\mathbf{h}^0\|$$

showing that $\mathbf{x}^2 \in X_0$. Further, the ball

$$X_1 := \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^2\| \leq \|\mathbf{h}^1\|\}$$

is contained in X_0 . The inequality (17b) shows that the hypotheses of our theorem remain true if we replace \mathbf{x}^0 and \mathbf{h}^0 by \mathbf{x}^1 and \mathbf{h}^1 , respectively. The same argument can be repeated for all \mathbf{x}^i and \mathbf{h}^i given by (9), $i = 0, 1, \dots$

The following analogs of (17) therefore hold for all i ,

$$|f(\mathbf{x}^i)| \leq \frac{1}{2} \|\mathbf{h}^{i-1}\|^2 M, \quad (20a)$$

$$|\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i| \geq 2 \|\mathbf{h}^i\| M, \quad (20b)$$

$$\text{and } \|\mathbf{h}^{i+1}\| \leq \frac{1}{2} \|\mathbf{h}^i\|, \quad (20c)$$

showing that the nested balls

$$X_i := \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^{i+1}\| \leq \|\mathbf{h}^i\|\}$$

have radii converging to zero. Therefore there is a unique point \mathbf{x}^* contained in all the balls X_i , and $\mathbf{x}^i \rightarrow \mathbf{x}^*$ since each \mathbf{x}^i is the center of the ball X_{i+1} .

We prove now that \mathbf{x}^* is a zero of $f(\mathbf{x})$.

$$\begin{aligned} |f(\mathbf{x}^i)| &\leq \frac{1}{2} \|\mathbf{h}^{i-1}\|^2 M, \text{ by (20a),} \\ &\leq \left(\frac{1}{2}\right)^i \|\mathbf{h}^0\|^2 M, \text{ by (20c).} \end{aligned}$$

$$\therefore f(\mathbf{x}^*) = 0.$$

Part 4. Proof of claim (c).

We prove that $\nabla f(\mathbf{x}^*) \neq \mathbf{0}$ except if \mathbf{x}^* is diametrically opposite to \mathbf{x}^0 . For any $\mathbf{x} \in X_0$ we have, by (15),

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^0)\| \leq M \|\mathbf{x} - \mathbf{x}^0\| \leq 2M \|\mathbf{h}^0\|.$$

If \mathbf{x} is not diametrically opposite to \mathbf{x}^0 , i.e. if $\|\mathbf{x} - \mathbf{x}^0\| < 2 \|\mathbf{h}^0\|$, then

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^0)\| < 2M \|\mathbf{h}^0\| \leq \|\nabla f(\mathbf{x}^0)\|,$$

or

$$\|\nabla f(\mathbf{x}^0)\| > \|\nabla f(\mathbf{x}^0) - \nabla f(\mathbf{x})\|,$$

showing that $\nabla f(\mathbf{x}) \neq \mathbf{0}$.

Part 5. Proof of claim (d).

The inequality (11a) is equivalent to

$$\|\mathbf{h}^i\| \leq \frac{M \|\mathbf{h}^{i-1}\|^2}{2 \|\nabla f(\mathbf{x}^i)\|}, \quad i = 1, 2, \dots \quad (21)$$

Indeed,

$$\|\mathbf{h}^i\| = \frac{|f(\mathbf{x}^i)|}{|\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i|} \leq \frac{M \|\mathbf{h}^{i-1}\|^2}{2 |\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i|}, \text{ by (20a).}$$

To prove (11b), we note

$$\|\mathbf{x}^* - \mathbf{x}^{i+1}\| \leq \|\mathbf{x}^{i+1} - \mathbf{x}^i\| = \|\mathbf{h}^i\|,$$

since \mathbf{x}^{i+1} is the center of the ball X_i of radius $\|\mathbf{h}^i\|$, \mathbf{x}^i is on the boundary of X_i and \mathbf{x}^* is in X_i . Then (11b) follows from (11a). \square

Lemma 1. *If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable for all \mathbf{x} in a convex region $X_0 \in \mathbb{R}^n$, and if there exists a constant M such that*

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq M \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in X_0, \quad (22a)$$

then

$$|f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y})| \leq \frac{M}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X_0. \quad (22b)$$

Proof. Follows from [10, Lemma 5.3.1]. \square

Remark. The Lipschitz bound M in (22a) can be used in (8a) if f is twice differentiable.

Theorem 2. *(Quadratic convergence of the directional Newton method)*

Let the assumptions of Theorem 1 hold, and in addition assume that there exist positive constants L, M such that

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq M \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in X_0, \quad (23a)$$

$$|\nabla f(\mathbf{x}) \cdot \mathbf{d}| \geq \frac{1}{L}, \quad \forall \mathbf{x} \in X_0, \mathbf{d} \in \mathbb{R}^n, \|\mathbf{d}\| = 1, \quad (23b)$$

$$\text{and } q := \frac{ML \|\mathbf{h}^0\|}{2} < 1. \quad (23c)$$

Then

$$\|\mathbf{x}^* - \mathbf{x}^i\| \leq \|\mathbf{h}^0\| \frac{q^{2^i - 1}}{1 - q^{2^i}}, \quad i = 1, 2, \dots \quad (24)$$

Note: Since $0 < q < 1$, the inequalities (24) show that the directional Newton method is at least quadratically convergent.

Proof. By Theorem 1, part (a), all iterates \mathbf{x}^i lie in X_0 .

$$\begin{aligned} \therefore \|\mathbf{x}^{i+1} - \mathbf{x}^i\| &= \frac{|f(\mathbf{x}^i)|}{|\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i|} \\ &\leq L |f(\mathbf{x}^i)|, \quad \text{by (23b)}, \\ &= L |f(\mathbf{x}^i) - f(\mathbf{x}^{i-1}) - \nabla f(\mathbf{x}^{i-1}) \cdot (\mathbf{x}^i - \mathbf{x}^{i-1})|, \end{aligned}$$

since

$$f(\mathbf{x}^{i-1}) + \nabla f(\mathbf{x}^{i-1}) \cdot (\mathbf{x}^i - \mathbf{x}^{i-1}) = 0$$

follows by multiplying both sides of

$$\mathbf{x}^i - \mathbf{x}^{i-1} = -\frac{f(\mathbf{x}^{i-1})}{\nabla f(\mathbf{x}^{i-1}) \cdot \mathbf{d}^{i-1}} \mathbf{d}^{i-1} \quad \text{with } \nabla f(\mathbf{x}^{i-1}).$$

$$\therefore \|\mathbf{x}^{i+1} - \mathbf{x}^i\| \leq \frac{LM}{2} \|\mathbf{x}^i - \mathbf{x}^{i-1}\|^2, \quad \text{by (22b)}. \quad (25)$$

We prove by induction that

$$\|\mathbf{x}^{i+1} - \mathbf{x}^i\| \leq \|\mathbf{h}^0\| q^{2^i - 1}, \quad i = 0, 1, \dots \quad (26)$$

By definition this inequality holds for $i = 0$. Assume it is correct for $i - 1$.

$$\begin{aligned} \therefore \|\mathbf{x}^{i+1} - \mathbf{x}^i\| &\leq \frac{LM}{2} \|\mathbf{x}^i - \mathbf{x}^{i-1}\|^2, \text{ by (25),} \\ &\leq \frac{LM}{2} \|\mathbf{h}^0\|^2 q^{2^i-2}, \text{ by (26) for } i-1, \\ &= \|\mathbf{h}^0\| q^{2^i-1}, \text{ proving (26) for } i. \end{aligned}$$

For $m > n$ we therefore have

$$\begin{aligned} \|\mathbf{x}^m - \mathbf{x}^n\| &\leq \|\mathbf{x}^m - \mathbf{x}^{m-1}\| + \|\mathbf{x}^{m-1} - \mathbf{x}^{m-2}\| + \dots + \|\mathbf{x}^{n+1} - \mathbf{x}^n\| \\ &\leq \|\mathbf{h}^0\| q^{2^n-1} (1 + q^{2^n} + (q^{2^n})^2 + \dots + (q^{2^n})^{m-n-1}), \\ &\quad \text{by (26),} \\ &\leq \|\mathbf{h}^0\| q^{2^n-1} (1 + q^{2^n} + (q^{2^n})^2 + \dots), \\ &= \frac{\|\mathbf{h}^0\| q^{2^n-1}}{1 - q^{2^n}}. \end{aligned}$$

$$\therefore \lim_{m \rightarrow \infty} \|\mathbf{x}^m - \mathbf{x}^n\| = \|\mathbf{x}^* - \mathbf{x}^n\| \leq \frac{\|\mathbf{h}^0\| q^{2^n-1}}{1 - q^{2^n}}, \text{ proving (24).}$$

□

The main special case of Theorem 1 is the **gradient method** discussed next.

Corollary 1. For $i = 0, 1, \dots$ define the direction \mathbf{d}^i and step \mathbf{h}^i by

$$\mathbf{d}^i := \frac{\nabla f(\mathbf{x}^i)}{\|\nabla f(\mathbf{x}^i)\|}, \quad (27a)$$

$$\mathbf{h}^i := -\frac{f(\mathbf{x}^i)}{\|\nabla f(\mathbf{x}^i)\|^2} \nabla f(\mathbf{x}^i), \quad (27b)$$

and let $f, \mathbf{x}^0, X_0, M, L, q$ be as in Theorems 1–2. Then the iterations

$$\mathbf{x}^{i+1} := \mathbf{x}^i - \frac{f(\mathbf{x}^i)}{\|\nabla f(\mathbf{x}^i)\|^2} \nabla f(\mathbf{x}^i), \quad i = 0, 1, \dots, \quad (28)$$

satisfy all conclusions of Theorems 1–2. □

3. NEWTON DIRECTIONS ALONG MAXIMAL MODULUS COMPONENTS OF THE GRADIENT

A directional Newton method not covered by Theorem 1 is where the direction \mathbf{d} in each iteration is chosen as the unit vector along the maximal absolute value $|\frac{\partial f}{\partial x_j}|$. However, the proof of Theorem 1 applies here, with small changes.

Theorem 3. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable, let \mathbf{x}^0 be a point where

$$f(\mathbf{x}^0) \neq 0, \quad \nabla f(\mathbf{x}^0) \neq \mathbf{0}.$$

Let $m(0)$ be an index such that

$$|f'_*(\mathbf{x}^0)| := \left| \frac{\partial f}{\partial x_{m(0)}}(\mathbf{x}^0) \right| = \max_{j=1,\dots,n} \left| \frac{\partial f}{\partial x_j}(\mathbf{x}^0) \right|$$

and define the vectors \mathbf{h}^0 and \mathbf{x}^1 by

$$\mathbf{h}^0[k] := \begin{cases} -\frac{f(\mathbf{x}^0)}{|f'_*(\mathbf{x}^0)|} & , k = m(0) , \\ 0 & , k \neq m(0) , \end{cases} \quad (29a)$$

$$\mathbf{x}^1 := \mathbf{x}^0 + \mathbf{h}^0 . \quad (29b)$$

Consider the interval

$$X_0 := \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^1\|_\infty \leq \|\mathbf{h}^0\|_\infty\}$$

and assume that $f \in C^2[X_0]$, and

$$\sup_{\mathbf{x} \in X_0} \|f''(\mathbf{x})\|_\infty = M , \quad (30a)$$

$$\left| \frac{\partial f(\mathbf{x}^0)}{\partial x_{m(0)}} \right|^2 \geq 2 |f(\mathbf{x}^0)| M . \quad (30b)$$

Define sequences $\{\mathbf{x}^i\}$, $\{\mathbf{h}^i\}$ recursively as follows. Let $m(i)$ be an index of the maximal modulus of $\frac{\partial f(\mathbf{x}^i)}{\partial x_j}$,

$$|f'_*(\mathbf{x}^i)| := \left| \frac{\partial f}{\partial x_{m(i)}}(\mathbf{x}^i) \right| = \max_{j=1,\dots,n} \left| \frac{\partial f}{\partial x_j}(\mathbf{x}^i) \right| ,$$

$$\mathbf{h}^i[k] := \begin{cases} -\frac{f(\mathbf{x}^i)}{|f'_*(\mathbf{x}^i)|} & , k = m(i) , \\ 0 & , k \neq m(i) , \end{cases} \quad (31a)$$

$$\mathbf{x}^{i+1} := \mathbf{x}^i + \mathbf{h}^i . \quad (31b)$$

Then

- (a) All \mathbf{x}^i lie in X_0 .
- (b) The sequence $\{\mathbf{x}^i\}$ converges, as $i \rightarrow \infty$, to a point $\mathbf{x}^* \in X_0$ that is a zero of f .
- (c) Moreover, $\nabla f(\mathbf{x}^*) \neq \mathbf{0}$ unless $\|\mathbf{x}^* - \mathbf{x}^0\|_\infty = 2 \|\mathbf{h}^0\|_\infty$.
- (d) Further, for $i = 1, 2, \dots$

$$\|\mathbf{x}^{i+1} - \mathbf{x}^i\|_\infty \leq \frac{M}{2 \|\nabla f(\mathbf{x}^i)\|_\infty} \|\mathbf{x}^i - \mathbf{x}^{i-1}\|_\infty^2 , \quad (32a)$$

$$\|\mathbf{x}^* - \mathbf{x}^{i+1}\|_\infty \leq \frac{M}{2 \|\nabla f(\mathbf{x}^i)\|_\infty} \|\mathbf{x}^i - \mathbf{x}^{i-1}\|_\infty^2 . \quad (32b)$$

Proof. The proof of Theorem 1 can be adapted here, by replacing each occurrence of $|\nabla f(\mathbf{x}^i) \cdot \mathbf{d}^i|$ by $\|\nabla f(\mathbf{x}^i)\|_\infty$, and by using the ∞ -norm

instead of the Euclidean norm. For example, condition (30b) is analogous to (8b). It can be rewritten as

$$\|\nabla f(\mathbf{x}^0)\|_\infty \geq 2 \|\mathbf{h}^0\|_\infty M, \quad (33)$$

the analog of (14), etc. \square

The convergence rate of method (5) is quadratic under conditions analogous to Theorem 2. The proof is analogous to that of Theorem 2.

4. SYSTEMS OF EQUATIONS

Consider an arbitrary system of m equations in n unknowns:

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}, \text{ or } f_i(x_1, x_2, \dots, x_n) = 0, \quad i = 1, \dots, m. \quad (34)$$

If $m = n$ the usual Newton method for solving (34) uses the iterations

$$\mathbf{x}^{k+1} := \mathbf{x}^k - J_{\mathbf{f}}(\mathbf{x}^k)^{-1} \mathbf{f}(\mathbf{x}^k), \quad k = 0, 1, \dots, \quad (35)$$

where the Jacobian matrix

$$J_{\mathbf{f}}(\mathbf{x}) := \begin{pmatrix} \frac{\partial f_i}{\partial x_j} \end{pmatrix},$$

is assumed nonsingular. If $J_{\mathbf{f}}(\mathbf{x})$ is singular, or if $m \neq n$, the Moore-Penrose inverse can be used in (35),

$$\mathbf{x}^{k+1} := \mathbf{x}^k - J_{\mathbf{f}}(\mathbf{x}^k)^\dagger \mathbf{f}(\mathbf{x}^k), \quad k = 0, 1, \dots, \quad (36)$$

see [1]. Note that the gradient method (28) is a special case of (36), corresponding to $m = 1$.

Alternatively, several authors (e.g. [3, Vol II, p. 165], [9, p. 362]) suggest minimizing the sum of squares $\sum_{i=1, \dots, m} f_i^2(\mathbf{x})$, using a suitable method, such as the steepest descent method. Following this idea, consider the single equation,

$$F(\mathbf{x}) := \sum_{i=1, \dots, m} f_i^2(\mathbf{x}) = 0, \quad (37)$$

that is equivalent to the system (34), in the sense of having the same solutions. There are now two (Newton) approaches to solving (34) or (37):

Approach 1: Apply a vector Newton method, e.g. (35) or (36), to the system (34).

Approach 2: Apply a directional Newton method, e.g. (28) or (5), to equation (37).

An advantage of Approach 2 is that it avoids the inversion of the Jacobian matrix, as in (35) or (36).

Another advantage is suitability for parallel implementation.

Disadvantages of Approach 2 include slower convergence. Indeed, the gradient of (37) is

$$\nabla F(\mathbf{x}) = 2 \sum_{i=1, \dots, m} f_i(\mathbf{x}) \nabla f_i(\mathbf{x})$$

and as the values $f_i(\mathbf{x})$ get close to zero, the condition (14) may no longer hold. Thus quadratic convergence is lost.

Example 1. Another example of what can go wrong in Approach 2 is illustrated by the system of n equations in n unknowns,

$$\begin{aligned} f_1(\mathbf{x}) &:= 1 - x_1 = 0, \\ f_i(\mathbf{x}) &:= 10(x_{i-1} - x_i^2) = 0, \quad i = 2, \dots, n, \end{aligned}$$

that is solved quickly by the Newton method (35), for all n that can be handled. However, the sum of squares $\sum f_i^2$ is the notorious Rosenbrock function

$$R_n(\mathbf{x}) := (1 - x_1)^2 + 100 \sum_{i=2}^n (x_{i-1} - x_i^2)^2, \quad (38)$$

for which no directional Newton method works well, even for small n . \square

To salvage Approach 2 we can apply, in iteration k , the Newton method to the modified function

$$\frac{F(\mathbf{x})}{\|\mathbf{x} - \mathbf{x}^{k-1}\|^\alpha}, \quad \text{for suitably chosen } \alpha > 0, \quad (39)$$

to get a next point \mathbf{x}^{k+1} , see e.g. [2]. The denominator in (39) creates a barrier at \mathbf{x}^{k-1} , and the resulting method may be called a **Newton barrier method**.

The modified function (39) is “steeper” than the original F , making it more likely to satisfy (14). However, the modified function is less likely to satisfy the second derivative condition (8a). The choice of α will be studied elsewhere.

Example 2. (Intersection of surfaces) In computer graphics it is often required to compute and display the intersection $\mathcal{C} := \mathcal{S}_1 \cap \mathcal{S}_2$ of two surfaces $\mathcal{S}_1, \mathcal{S}_2$ in \mathbb{R}^3 . The surfaces in question may be represented explicitly

$$\mathcal{S}_1 := \{(x, y, z) : z = f(x, y)\}, \quad \mathcal{S}_2 := \{(x, y, z) : z = g(x, y)\}$$

in which case their intersection \mathcal{C} is the set of points $\mathbf{r} = (x, y, z)$ satisfying

$$f(x, y) = g(x, y) \quad (40)$$

and then $z := f(x, y)$. If the surfaces are given implicitly

$$\mathcal{S}_1 := \{(x, y, z) : F(x, y, z) = 0\}, \quad \mathcal{S}_2 := \{(x, y, z) : G(x, y, z) = 0\}$$

then the intersection points satisfy

$$F(x, y, z) = 0, \quad G(x, y, z) = 0. \quad (41)$$

The so-called **marching methods** compute a **starting solution** $\mathbf{r}^0 = (x^0, y^0, z^0) \in \mathcal{C}$ and from it compute other intersection points sufficiently close to each other to enable displaying the curve \mathcal{C} , see [7].

A starting solution requires solving the single equation (40) in 2 unknowns, or the system (41) of 2 equations in 3 unknowns. Either way, the

Newton method (35) cannot be used. The remaining alternatives include the method (36) using the generalized inverse of the Jacobian matrix, (see [1], [7]), or the directional Newton methods of this paper.

For any point $\mathbf{r}^k = (x^k, y^k, z^k) \in \mathcal{C}$, a next point \mathbf{r}^{k+1} is found by solving (40) or (41) with an additional equation, such as

$$x = x^k + \Delta, \quad \Delta \text{ sufficiently small,}$$

to give movement in the x direction. One can solve (in parallel) for movements in the y and z directions, then select the point closest to \mathbf{r}^k . \square

Example 3. (Complex roots of non-analytic functions) If $f(z)$ is analytic, its roots $z = x + iy$ can be approximated by the complex Newton method

$$z^{k+1} := z^k - \frac{f(z^k)}{f'(z^k)}, \quad k = 0, 1, \dots \quad (42)$$

or equivalently, by solving the system of 2 equations in 2 unknowns

$$\Re f(x, y) = 0, \quad \Im f(x, y) = 0. \quad (43)$$

If f is not analytic then (42) cannot be used, however the system (43) may still be solved by the real Newton method (35). If f is a non-analytic function of several complex variables, say $f(z_1, \bar{z}_1, z_2, \bar{z}_2, \dots, z_m, \bar{z}_m)$, then (43) is a system of 2 equations in $2m$ variables, and its surrogate equation $\Re^2 f + \Im^2 f = 0$ can be solved by the above directional Newton methods. \square

ACKNOWLEDGMENT

We thank Professor Arkadi Nemirovski and the referees for their help and constructive suggestions.

REFERENCES

- [1] A. Ben-Israel, *A Newton-Raphson method for the solution of systems of equations*, J. Math. Anal. Appl. **15**(1966), 243–252
- [2] A. Ben-Israel, *Newton's method with modified functions*, Contemp. Math. **204**(1997), 39–50
- [3] I.S. Berezin and N.P. Zhidkov, *Computing Methods*, Pergamon Press, 1965
- [4] W. Fleming, *Functions of Several Variables*, 2nd Edition, Springer, 1977
- [5] C.-E. Fröberg, *Numerical Mathematics: Theory and Computer Applications*, Benjamin, 1985
- [6] Y. Levin and A. Ben-Israel, MAPLE programs for directional Newton methods are available at:
ftp://rutcor.rutgers.edu/pub/bisrael/Newton-Dir.mws
- [7] G. Lukács, *The generalized inverse matrix and the surface-surface intersection problem*, pp. 167–185 in *Theory and Practice of Geometric Modeling* (W. Strasser and H.-P. Seidel, editors), Springer-Verlag, 1989

- [8] A.M. Ostrowski, *Solution of Equations in Euclidean and Banach Spaces*, 3rd Edition, Academic Press, 1973
- [9] A. Ralston and P. Rabinowitz, *A First Course in Numerical Analysis*, 2nd edition, McGraw-Hill, 1978
- [10] J. Stoer and K. Bulirsch, *Introduction to Numerical Analysis*, Springer-Verlag, 1976

YURI LEVIN, RUTCOR–RUTGERS CENTER FOR OPERATIONS RESEARCH, RUTGERS UNIVERSITY, 640 BARTHOLOMEW RD, PISCATAWAY, NJ 08854-8003, USA
E-mail address: `ylevin@rutcor.rutgers.edu`

ADI BEN-ISRAEL, RUTCOR–RUTGERS CENTER FOR OPERATIONS RESEARCH, RUTGERS UNIVERSITY, 640 BARTHOLOMEW RD, PISCATAWAY, NJ 08854-8003, USA
E-mail address: `bisrael@rutcor.rutgers.edu`