Proc. Estonian Acad. Sci. Phys. Math., 1994, **43**, 2, 49 – 63 https://doi.org/10.3176/phys.math.1994.2.01

## ON SOLVING ILL-CONDITIONED SYSTEMS OF NON-LINEAR EQUATIONS

#### Otu VAARMANN

Tallinna Tehnikaülikool (Tallinn Technical University), Akadeemia tee 1, EE-0026 Tallinn, Eesti (Estonia)

Presented by G. Vainikko

Received October 7, 1993; revised version received December 22, 1993 and February 15, 1994; accepted February 15, 1994

Abstract. Methods under consideration are intended to solve systems of nonlinear equations in n variables where the number of equations m is greater or equal to n and the Jacobian may be illconditioned. Two approaches for solving such problems are examined. The first one is based on the computation of the weighted pseudoinverse being, as a rule, concurrent with great computational efforts and therefore the corresponding algorithms are appropriate for solving only small dimensional problems, especially if m >> n. The second approach rests on solving approximately corresponding preconditioned linear equations by taking finitely many steps of an iterative procedure. Thus they are applicable to large scale nonlinear problems as well.

**Key words:** nonlinear equations, least squares solution, weighted pseudoinverse, preconditioner, Neumann's series, regularization.

Giving up some extra work for the numerical stability is particularly justified for problems in which the evaluation of function values and their derivatives are expensive and dominate the other costs. For instance, inverse problems are usually nonlinear and often involve very complicated subproblems with partially unknown properties. As a rule, they are illconditioned and numerically unstable. In this case it is reasonable to combine the algorithm with a preconditioner to reduce the condition number.

#### METHODS

The problem is to solve the nonlinear equation,

$$F(x) = 0,$$

where  $F: D \subset \mathbb{R}^n \to \mathbb{R}^m$ ,  $m \neq n$  and F is continuously differentiable. If the equation (1) has any solution in the classical sense, then it is sought in the

Teaduste Akadeemia Raamatukogu

49

(1)

50.6.31

least squares sense, i.e. minimizing the functional

$$T(x) = \frac{1}{2} \|F(x)\|_{2}^{2}.$$
 (2)

**1.** For solving (1) or for seeking a local minimum of (2), one might use the iterative method

$$x_{k+1} = x_k - [WF'(x_k)] + WF(x_k), \qquad (3)$$

where  $A_k$ : =  $[WF'(x_k)]^+W = [F'(x_k)]^+_{WI}$ ,  $W \in \mathbb{R}^{n \times m}$  and  $[F'(x)]^+_{WI}$ means the WI – weighted pseudoinverse of F'(x) [<sup>1</sup>].

If  $m \ge n$  and F'(x) and W have a full rank, then one gets

$$x_{k+1} = x_k - [WF'(x_k)]^{-1}WF(x_k), \qquad (4)$$

which for m=n coinsides with the standard Newton method, and taking  $W = W_k = \alpha_k [F'(x_k)]^T, \alpha_k > 0$ , one obtains the Gauss-Newton method for  $m \ge n$ . Here  $[\cdot]^T$  denotes the dual mapping.

Let  $P_{R(x)}$  denote the orthogonal projector of  $R^{m}$  on to the range of F'(x), and let  $P_{R(x)}$  be defined by the relation  $P_{R(x)} = P_{R(x)} \cdot P_{W}$ , where  $P_{W} = W^{+}W$ , and W is an  $n \times m$  matrix approximating  $[F'(x)]^{+}$ . Let  $\lambda$  denote a positive scalar satisfying the inequality  $||A_{k}|| \le \lambda_{k} = \lambda \le \infty$  and  $P_{k} = P_{R(F'(x_{k}))}$ .

Lemma 1 [<sup>2, 3</sup>]. Let on some set  $S = \{x \in \mathbb{R}^n : ||x - x_0|| \le \rho\}$  the following conditions be fulfilled

$$\| [F'(x)]^+ \| \le C, \quad \forall x \in S.$$
$$\| F'(x) - F'(y) \| \le L \| x - y \|, \quad \forall x, y \in S,$$

then

$$\|P_{R(x)} - P_{R(y)}\| \le L_0 \|x - y\|, \quad \forall x, y \in S,$$

where

$$L_0 = CL.$$

R e m a r k 1. From Lemma 1 it follows that there exist positive scalars N and N such that

$$\left\| \left( P_{R(y)} - P_{R(y)} P_{R(x)} \right) F(x) \right\| = \left\| \left( P_{R(y)} - P_{R(x)} \right) \left( I - P_{R(x)} \right) F(x) \right\| =$$
$$= N' \|x - y\|, \qquad \forall x, y \in S,$$

and for  $x_k, x_{k+1} \in S$ 

$$\left\| (P_{k+1} - P_{k+1}P_k) F(x_k) \right\| \le N \left\| P_k F(x_k) \right\|$$

where

$$\sup_{x \in S} \| (I - P_{R(x)}) F(x) \| \le N'', N' = N'' L_0 \text{ and } N = \lambda N'.$$

Since, in general, the equality  $\overline{P}_{R(x)} = \overline{P}_{R(x)}^2$  does not hold, then

 $\overline{P}_{R(y)} - \overline{P}_{R(y)}\overline{P}_{R(x)} = (\overline{P}_{R(y)} - \overline{P}_{R(x)}) (I - \overline{P}_{R(x)}) + \overline{P}_{R(x)} - \overline{P}_{R(x)}^{2}.$ For  $W = [F'(x_{0})]^{+}$  we have

 $\left\| \overline{P}_{R(y)} - \overline{P}_{R(x)} \right\| = \left\| (\overline{P}_{R(y)} - \overline{P}_{R(x)}) P_0 \right\| \le \left\| P_{R(y)} - P_{R(x)} \right\|,$  $\overline{P}_{R(x)} - \overline{P}_{R(x)}^2 = P_{R(x)} \overline{P}_{R(x)} - P_{R(x)} P_0 \overline{P}_{R(x)} = (P_{R(x)} - P_{R(x)} P_0) \overline{P}_{R(x)}.$ 

By Lemma 1 and Remark 1

$$\left\| \overline{P}_{R(x)} - \overline{P}_{R(x)}^{2} \right\| \leq \tilde{N} \|x - x_{0}\| \leq \tilde{N}\rho.$$

$$\left\| (\overline{P}_{k+1} - \overline{P}_{k+1}\overline{P}_{k})F(x_{k}) \right\| \leq (\overline{N} + \tilde{N}\rho) \|P_{k}F(x_{k})\|$$

where

$$\sup_{x \in S} \left\| (I - \overline{P}_{R(x)}) F(x) \right\| \le \overline{N}^{"}, \quad \overline{N}^{'} = \overline{N}^{"} L_{0}, \quad \overline{N} = L_{0} \left\| I - P_{0} \right\|$$

and  $N = \lambda N'$ .

First, we will examine convergence properties of (4) with  $W = [F'(x_0)]^+$ , and thus  $\overline{P}_{R(x)} = P_{R(x)}P_0$ .

Theorem 1. Let  $m \ge n$ ,  $x_0 \in \mathbb{R}^n$ ,  $S = \{x \in \mathbb{R}^n : ||x - x_0|| \le \rho\}$  and the following conditions are valid in S:

1° operator F has Frechet-derivative and  $|| [F'(x)]^+ || \le C$ ;

2° derivative F' satisfies Lipschitz condition

 $||F'(x) - F'(y)|| \le L||x - y||;$ 

 $3^{\circ}$  derivative F'(x) has full rank.

If  $\delta = \overline{N} + \overline{N}r + \frac{LC^2}{2} \|P_0F(x_0)\| \le 1$  and  $r = C \|P_0F(x_0)\|/(1-\delta) \le \rho$ then the sequence generated by (4) with  $W = [F'(x_0)]^+$  has a limit  $x^*$  which appears to be a solution of the equation  $P_{R(x)}F(x) = 0$  with  $\|x_k - x_0\| \le r$  and

$$\left\|x_k - x^*\right\| \le r\delta^k.$$

Proof. If A and B are  $m \times n$  and  $n \times m$  matrices, respectively, then  $(AB)^{+}=B^{+}A^{+}$ , provided rank(A) = rank(B) = n [<sup>4</sup>]. On the basis of this relation one concludes that

$$A_{k} = \{ [F'(x_{0})]^{+} F'(x_{k}) \}^{+} [F'(x_{0})]^{+} = [F'(x_{k})]^{+} F'(x_{0}) [F'(x_{0})]^{+} =$$

$$= [F'(x_k)]^+ P_0 = [F'(x_k)]^+ P_k P_0 = [F'(x_k)]^+ P_k.$$

According to Taylor expansion, one obtains

$$\begin{split} \bar{P}_{k+1}F(x_{k+1}) &= (\bar{P}_{k+1} - \bar{P}_{k+1}\bar{P}_k)F(x_k) + \bar{P}_{k+1} \left\{ (\bar{P}_k - F'(x_k)A_k)F(x_k) + \right. \\ &+ \int_0^1 \left[ F'(x_k) - F'(x_k + t(x_{k+1} - x_k)) \right] A_k F(x_k) dt \left. \right\}. \end{split}$$

Further

 $\overline{P}_{k} - F'(x_{k})A_{k} = P_{k}\overline{P}_{k} - F'(x_{k})A_{k}\overline{P}_{k} = (P_{k} - F'(x_{k})[F'(x_{k})]^{+})\overline{P}_{k} = 0$ since  $P_{k} = F'(x_{k})[F'(x_{k})]^{+}$  and

$$\begin{split} \left\| \overline{P}_{k+1} F(x_{k+1}) \right\| &\leq \left\| \overline{N} + \widetilde{N}r + \frac{LC^2}{2} \right\| \overline{P}_k F(x_k) \| \right\| \cdot \left\| \overline{P}_k F(x_k) \right\| \leq \\ &\leq \delta \left\| \overline{P}_k F(x_k) \right\| \leq \left\| \overline{P}_0 F(x_0) \right\| \delta^{k+1}, \\ &\left\| x_{k+1} - x_k \right\| \leq C \left\| \overline{P}_k F(x_k) \right\| \leq C \left\| \overline{P}_0 F(x_0) \right\| \delta^k, \\ &\left\| x_p - x_k \right\| \leq r \left( \delta^k - \delta^p \right), \quad p \geq k, \\ &\left\| x_p - x_0 \right\| \leq r \left( 1 - \delta^p \right) \leq r \leq \rho, \quad \left\| x^* - x_k \right\| \leq r \delta^k, \\ &x^* = \lim_{k \to \infty} x_k, \quad \left\| P_{R(x^*)} \cdot P_o \cdot F(x^*) \right\| = \lim_{k \to \infty} \left\| \overline{P}_k F(x_k) \right\| = 0. \end{split}$$

In particular, if W is an  $(m \times m)$ -unitary matrix then it can be shown that  $[WF'(x)]^+=[F'(x)]^+W^{-1}$ , and therefore  $x^*$  appears to be a solution of the nonlinear normal equation [<sup>5</sup>]. The problems of global convergence for (3) and the concept of appropriate level functions are discussed in the next section.

2. When large-scale nonlinear problems are to be solved, then the exact solution of associated linear equations at each iteration can be very expensive, not to speak of matrix inversion. In this case, the use of an inexact standard iterative method (e.g. truncated-Newton, truncated

Gauss-Newton, etc.) may be justified, which is based on solving corresponding equations approximately, frequently via the conjugate gradient method. Iterative methods solve Kx = g sometimes faster than elimination methods. Moreover, finding a current approximate solution  $x_k$  with high accuracy may not be necessary when  $x_k$  is far from the solution  $x^*$ .

In particular, if the Newton equation  $F'(x_k)(x_{k+1}-x_k) = -F(x_k)$  is to be solved within the tolerance  $\eta_k \| F(x_k) \| (\eta_k > 0)$ , then the quantity  $\eta_k$ , accurate to a factor, coincides with the quantity  $\gamma_k$ , where  $\| I - F'(x_k)A_k \| = \gamma_k$  and  $A_k \approx [F'(x_k)]^{-1}$ . The sequence  $\{\eta_k\}$  is known as forcing sequence, and it can be used for controlling the computational process; specifically, the approximation rate  $\gamma_k = 0 \| F(x_k) \|$  ensures that the sequence  $\{x_k\}$  converges quadratically (e.g. [6, 7]). The search direction is usually assessed by the norm of the residual. However, the norm of the residual happens to be not always a good predictor of a good search direction, especially if the problem to be solved is ill-conditioned  $[^{8, 9}]$  and hence preconditioning is needed. In particular, for symmetric and positive definite linear systems of equations Kx = g, a combination of iterative methods and a preconditioner based on incomplete Cholesky factorization (LU factorization) is applicable. That is, one solves iteratively the equivalent system  $N^{-1}Kx = N^{-1}g$ , where N is the matrix representation of the incomplete factorization. The convergence properties of the iterative process can be further improved by using a polynomical preconditioner

$$C(N^{-1}K)N^{-1}Kx = C(N^{-1}K)N^{-1}g,$$

where  $C(\lambda)$  is a preconditioning polynomial (usually  $C(\lambda) \approx \lambda^{-1}$ ) and  $C(N^{-1}K)$  is the associated polynomical preconditioner [<sup>10</sup>]. If N and K are symmetric, then so is the preconditioner  $C(N^{-1}K)N^{-1}$  and several variants of the conjugate gradient methods can be put to use.

The execution of one step of the iteration method (e.g. see [11])

$$x_{k+1} = x_k - \varepsilon_k \alpha_k C(\alpha_k [F'(x_k)]^T F'(x_k)) [F'(x_k)]^T F(x_k), \ \alpha_k > 0,$$

where the parameter  $\varepsilon_{\nu}$  determines the stepsize and

$$C(\alpha_{k}, x_{k}) = C(\alpha_{k}[F'(x_{k})]^{T}F'(x_{k})) = I + \Psi_{k} + \dots + \Psi_{k}^{q-1}$$

 $\Psi_k = I - \alpha_k [F'(x_k)]^T F'(x_k)$  with  $q(q \ge 1)$  an integer for  $\Psi^{q-1}$  can be regarded as the execution of one step of the ordinary iteration method as applied to the preconditioned system

$$\alpha_k C(\alpha_k, x_k) [F'(x_k)]^{\mathrm{T}} F(x) = 0, \qquad (5)$$

where  $\alpha [F'(x)]^T$  plays the role of W.

Let us take  $A_k := \alpha_k C(\alpha_k, x_k) [F'(x_k)]^T$  and rewrite it as  $A_k = \alpha_k [F'(x_k)]^T [P_k + E_k + \dots + E_k^{q-1}]$ , where  $E_k = P_k - \alpha_k F'(x_k) \times [F'(x_k)]^T$ . If  $0 < \alpha_k < 2/K_0^2$ , where  $||F'(x)|| \le K_0$  for all  $x \in S$ , then there exists a positive scalar  $\mu$  such that  $||E_k|| \le \mu < 1$  and  $\lim_{q \to \infty} \alpha_k [F'(x_k)]^T [P_k + E_k + \dots + E_k^{q-1}] = [F'(x_k)]^+$ ,

and for the iterative process

$$x_{k+1} = x_k - \varepsilon_k A_k F(x_k), \qquad (6)$$

where the parameter  $\varepsilon_{\nu}$  determines the stepsize one obtains.

The orem 2. Let, in addition to the conditions  $1^{\circ}-2^{\circ}$  of Theorem 1, the following conditions be fulfilled in S.

$$1^{\circ} \delta = \delta_0 = 1 - \varepsilon_0 + \varepsilon_0 (N + \mu^q) + \frac{1}{2} \varepsilon^2 \lambda^2 L \| P_0 F(x_0) \| < 1,$$

$$\lambda \equiv \lambda_{q-1} = \min\left\{\frac{2}{K_0} \left(1 + \mu + \dots + \mu^{q-1}\right), C\left(1 + \mu^q\right)\right\};$$

$$2^{\circ} 0 < \varepsilon_0 \le \varepsilon_{k-1} \le \varepsilon_k = \min\{1, \varepsilon_{k-1}, \delta^{-1/2}\}.$$

If  $r = \lambda (\|P_0F(x_0)\|/(1-\delta) \le \rho)$ , then the sequence generated by (6) has a limit which appears to be a solution of the equation  $[F'(x)]^T F(x) = 0$  with  $\|x_k - x_0\| \le r$  and

$$\left\|x_{k}-x^{*}\right\|\leq r\delta^{k}.$$

(cf. also Theorem 2 [<sup>2</sup>]).

Proof. Since

$$\|P_k - F'(x_k)A_k\| \le \|P_k - \alpha_k F'(x_k)[F'(x_k)]^T\|^q \le \mu^q$$

and  $A_k = A_k P_k$ , then

 $||A_k|| = ||[F'(x_k)]^+ (F'(x_k)A_k - P_k) + [F'(x_k)]^+ || \le C(1+\mu)^q.$ 

On the other hand,

$$||A_k|| \le \frac{2}{K_0} (1 + \mu + \dots + \mu^{q-1})$$

polied to the predot

and in the capacity of  $\lambda_{q-1}$  one can take

$$\lambda_{q-1} = \min\left\{\frac{2}{K_0} (1 + \mu + \dots + \mu^{q-1}), C(1 + \mu)^q\right\}$$

Taking U = I in Theorem 1 [<sup>3</sup>] and bearing in mind that  $\alpha_k$  stands in it for the regularization parameter, the straightforward application of this theorem yields Theorem 2 under consideration here.

If F'(x) has full rank, then another possibility to obtain the next iterate  $x_{k+1}$  is to solve at each iteration the linearized preconditioned equation (5), i.e. the equation  $W_k F'(x_k) (x - x_k) = W_k F(x_k)$  with  $W_k = A_k = \alpha_k C(\alpha_k, x_k) [F'(x_k)]^T$  or, equivalently, to compute  $x_{k+1} = x_k - [W_k F'(x_k)]^{-1} W_k F(x_k)$ , while  $[W_k F'(x_k)]^{-1} W_k = B^{-1}(x_k) \times [F'(x_k)]^T = [F'(x_k)]^+$ ,  $B(x) = [F'(x)]^T F'(x)$ . The preconditioner  $W_k$  suffers from the disadvantage that for ill-conditioned Jacobian a rank reduction may occur when the matrix B(x) is constructed. To overcome this drawback, a regularization technique is necessary.

Let H be an arbitrary symmetric positive definite matrix. One possibility to handle safely the rank-deficient (ill posed) problems is to use instead of B(x) the matrix M(x): =  $B(x) + \beta H$  with  $\beta > 0$ , not trying necessarily to find an optimal value for  $\beta$  but choosing it simply large enough to eliminate the singularity of the matrix B(x). Doing so, one obtains the following regularized iterative process

$$x_{k+1} = x_k - \varepsilon_k D_k \left[ F'(x_k) \right]^{\mathrm{T}} F(x_k) , \qquad (7)$$

where  $D_k = D_k^0 [I + ... + Q_k^{q-1}]$ ,  $Q_k = I - M_k D_k^0$ ,  $M_k = M(\beta_k, x_k)$  and  $D_k^0$ is an approximation to  $M_k^{-1}$  such that  $||Q_k|| < 1$ , while  $\lim_{q \to \infty} D_k = M_k^{-1}$  and  $\lim_{q \to \infty} D_k [F'(x_k)]^T = [F'(x_k)]^+$ . In particular, if  $0 < \alpha_k < \beta_{k \to 0}$ 

< 2/ ( $K_0^2 + \alpha ||H||$ ),  $\alpha_k \le \alpha$ , then  $D_k^0 = \alpha_k I$  is a proper choice.

On the basis of Theorem 1  $[^3]$ , one can prove for (7) analogous theorem as Theorem 2 for (6).

Alternatively, for solving (1) one might at each iteration to solve the linearized preconditioned equation  $W_k F'(x_k) (x - x_k) = W_k F(x_k)$  with  $W_k = D_k [F(x_k)]^T$ , or to compute  $x_{k+1} - x_k = -[W_k F'(x_k)]_R^+ W_k F(x_k)$ , where  $[\cdot]^+$  denote the reflexive generalized inverse  $[W_k F'(x_k)]_R^+ = B_k^+ D_k^{-1}$ , and

$$\lim_{\beta_k \to 0} [W_k F'(x_k)]_R^+ W_k = [F'(x_k)]^+.$$

R e m a r k 2. For improving the accuracy of the approximation for  $[F'(x)]^+$  based on the Neumann series, one can use the Chebyshev acceleration procedure (e.g. see [<sup>11</sup>]). When  $x_i$  (i = 1,...,n) vary widely in magnitude, it is reasonable to transform the variable space  $\bar{x} = Dx$  where D is a diagonal matrix.

### MOTIVATION OF METHODS

The usual convergence analysis is based on the monotonicity test of the objective function

$$T(x_{k+1}) \le T(x_k), \tag{8}$$

where  $\{x_k\}$  is a sequence generated by a particular local method. But frequently  $||F(x_k)||$  is not monotonically decreasing. For badly nonlinear F(x) monotonicity in T(x) may only ultimately occur when iterates  $\{x_k\}$ are sufficiently close to the solution  $x^*$ . Sometimes the use of relaxation parameter determining the appropriate step size produces a decrease of ||F(x)|| but at a cost of a lower speed of convergence.

That is why one is interested in solving a weighted least squares problem

$$\Phi(x): = \frac{1}{2} \|WF(x)\|_{2}^{2}, \qquad (9)$$

where W is a weighting matrix. In some case a simple weighting strategy is not sufficient and a dynamic weighting strategy is needed, i.e. one recalculates the matrix W at each iteration once again or periodically after several iterations. Thus, in general, W=W(x). The introduction of a weighting matrix means that one changes the notion of the distance, i.e. instead of  $||y||_2^2 = y^T y$  one uses  $||y||_W^2 = y^T W^T W y$ , where  $||\cdot||_W$  means, in general, a seminorm. In particular, if W is a symmetric positive definite matrix, it is usually called a scaling matrix or preconditioner. The basic idea of preconditioning is to introduce a preliminary scaling on the vector of independent variables and dependent variables. For example, instead of solving a linear equation Kx=g explicitly one often solves the equivalent preconditioned system

$$MKx = Mg, \tag{10}$$

where M is a scaling matrix. Specifically, using a diagonal weighting matrix D as a preconditioner rescales the dependent variables, i.e. one changes their units. If M is a proper approximation to  $K^{-1}$  or  $K^+$ , it will be a lot more efficient to solve (10) than Kx=g, because the spectral condition number cond<sub>2</sub>(MK) =  $\kappa$  (MK) is considerably less than  $\kappa(K)$  or than that

 $\kappa(K^{T}K) = \kappa^{2}(K)$  which is the case when the Gauss transformation is used.

For a given matrix A the objective function of nonlinear least squares (NLSQ) problems may be written as

$$T(x|A): = \frac{1}{2} \|AF(x)\|_{2}^{2}.$$
(11)

The solution (11) defines a level function T(x|A) and an associated with him level set in  $x_k$  [<sup>8</sup>]

$$G_k(A) = \{x \in S | T(x|A) \le T(x_k|A)\}, x \in S \subseteq \mathbb{R}^n.$$

The condition

$$T(x_{k+1}|A) \le T(x_k|A)$$

is equivalent to the requirement

 $x_{k+1} \in G_k(A)$  for some A.

Obviously,

$$0 = T(x^* | J(x^*)^+) \leq T(x | J(x^*)^+)$$

for all  $x \in S$ , where  $x^*$  denotes a solution of the NLSQ problem,  $A = J(x^*)^+$ and J(x) = F'(x). Thus the monotonicity test

$$x_{k+1} \in G_k(J(x^*)^{+})$$
(13)

would be as natural as (12). Note that this test uses a norm in the space of independent variables  $\mathbb{R}^n$ , which in many cases is more preferable than the usual monotonicity test (8) specifying a norm in the space  $\mathbb{R}^m$ . For instance, when for expanding the domain of convergence of a Gauss-Newton type method one uses a suitably chosen relaxation (damping) parameter  $\lambda$  to guarantee the convergence of it from a poor starting point, then for steep valley it is desirable to avoid the bottom of the valley as long as the full step of the method is unacceptable [<sup>12</sup>]. The property of global convergence is a criterion of robustness of the algorithm and therefore finding an appropriate globalization procedure is highly desirable for solving ill-conditioned NLSQ-s. Usually the line search or the trust region strategy is used for globalization of a method. Another efficient way for choosing damping parameter for ill-conditioned problems based on test (12) has been proposed by Deuflhard [<sup>8</sup>]. As pointed out in [<sup>8</sup>], the choice  $A = W(x) = J(x)^+$  turns out to be a most

As pointed out in [°], the choice  $A = W(x) = J(x)^+$  turns out to be a most natural choice, since

$$x_{k+1} = x_k - \text{grad} \left( T(x \mid J(x_k)^+) \mid_{x = x_k} \right)$$

i.e. the Gauss-Newton direction for the problem (2) is the steepest descent direction of  $T(x|J(x_{\nu})^{+})$ . Indeed, since

 $J(x)^* J(x)^+ = (J(x^+) J(x))^* = J(x^+) J(x)$  and  $J(x^+) J(x) J(x^+) = J(x^+)$ , one obtains

grad 
$$(T(x|J(x_k)^+)|_{x=x_k} = J(x)^*J(x)^{+}J(x_k)^+F(x)|_{x=x_k}$$
  
=  $J(x_k)^+F(x_k)$ .

Therefore  $T(x|J(x_k)^+)$  is called the natural level function. But, in general,

 $x^* \notin G_k(J(x_k)^+)$ 

except for a special case  $x_k = x^*$  (unfortunately, the element  $x^*$  is unknown) and therefore a global convergence proof cannot be based on the natural monotonocity test (13) and other convergence criteria are necessary. P. Deuflhard introduced the concept of appropriate level functions [13] and showed that the natural level function  $T(x|J(x_k)^+)$  is locally appropriate in  $x_k$ , and there exists some neighbourhood  $U(x_k)$  with  $J(y)^+ \in A_0(x_k)$  for all  $y \in U(x_k)$ , where  $A_0(x_k)$  denotes the set of matrices for each of them, the level function T(x|A) is locally appropriate in  $x_k$  [<sup>8, 13</sup>).

Let  $r(x) = (I - J(x)J(x)^+)F(x)$ ,  $\Delta x = -J(x)^+F(x)$ , and for fixed  $x_k$  let A(x) be some matrix set such that

 $(\Delta x_k)^{\mathrm{T}}$ grad  $T(x_k|A) < 0$ , if  $\Delta x_k \neq 0$ ,  $\forall x_k \in S$ , for all  $A \in A(x_k)$ .

Definition [8]. A level function T(x|A) is said to be locally appropriate in  $x_k$  if and only if  $A \in A_0(x_k)$ , where a special sub-set  $A_0(x) \subset A(x)$  is defined by

 $\|Ar(x_k)\| < \|AF(x_k)\|$ , if  $\Delta x_k \neq 0$ , for all  $A \in A_0(x_k)$ .

In general, the definition yields a not easily controllable condition  $\sigma(x) := \|WW^+ - WJ(x)J^+(x)W^+\| < 1$  for T(x|W) to be a locally appropriate level function. In particular, for  $W = [F'(x_0)]^+$  and for all  $x \in S$ , one obtains

$$\sup_{x \in S} \| [F'(x_0)]^+ (P_0 - P_0 P_{R(x)}) F'(x_0) \| < 1.$$

i.e.  $\sigma(x) \leq \kappa(J(x_0)) L_0 \rho < 1$ .

#### **ASPECTS OF NUMERICAL REALIZATION**

A common complaint about normal equation methods is that they are less accurate than orthogonal methods, since the latter ones preserve the condition number  $\kappa$  of the original problem, while during the solution of the normal system the condition number is squared. This argument is partially true since it ignores the presence of the term  $\kappa^2$  in the bound of the approximate pseudosolution error for the perturbed problem. Let us consider the system

$$Kx = g, \tag{14}$$

and a perturbed system

$$(K+\delta K)x = g+\delta g. \tag{15}$$

A sufficient condition of rank conservation of the problem (12) is the inequality  $||K^+|| || \delta K|| < 1$ . On the basis of the results by Wedin [<sup>14</sup>] and Stewart [<sup>15</sup>], one obtains

$$\|\delta x\| / \|x\| \le \kappa \{ (1 + \kappa \rho) \alpha + \gamma \beta \}, \tag{16}$$

$$\kappa = \kappa(K) = ||K|| ||K^+||, \ \kappa = \frac{\kappa}{1 - \kappa\alpha} = \frac{||K|| ||K^+||}{1 - ||\delta K|| ||K^+||}, \ \alpha = \frac{||\delta K||}{||K||},$$

$$\beta = \frac{\|\delta g\|}{\|g\|}, \ \gamma = \frac{\|g\|}{\|K\|\|\bar{x}\|}, \ \rho = \frac{\|\bar{r}\|}{\|K\|\|\bar{x}\|},$$

where  $\bar{x}$  is the pseudosolution of the problem (14), r is residual corresponding to  $\bar{x}$ , and  $\hat{x} = \bar{x} + \delta \bar{x}$  is the pseudosolution of the perturbed problem (15). Hence, for problems with large residual r the right-hand side of (16) includes the term with  $\kappa^2$  whose adverse influence is unavoidable regardless of the successful choice of the numerical method as applied to the problem. Thus, the use of a preconditioner might be fruitful.

Moreover, if J(x) is highly ill-conditioned, even a very small residual cannot guarantee an accurate solution. In other words, the rounded solution of an equation may have a large residual. The situation is more complicated when an accurate solution has a large residual, specifically when  $S(x) := [F''(x)]^T F(x)$  is larger in some sense than  $J(x)^T J(x)$ (obviously, if ||F(x)|| is large, then it may cause ||S(x)|| to be also large). It is known that in the latter case the sequence  $\{x_k\}$  generated by the Gauss-Newton method is always divergent even if  $J(x)^T J(x)$  is invertible [<sup>16</sup>]. Then Newton-type methods (e.g. augmented Gauss-Newton methods) as applied to (8) or (11) might be useful since they exploit more completely information on  $\nabla^2 T(x)$  or methods with the convergence order  $p \ge 3$ , as they are based on a quadratic model. Besides, sometimes, the latter permit to reduce the run times in seconds on the computer. Yet one reason for using methods with the convergence order higher than that of the Newton method is that variants of the Gauss-Newton method may experience great difficulty or fail completely if  $\nabla T(x)$  is singular or ill-conditioned at the solution point  $x^*$ , since they are based on a linear model. On the other hand, using even very rough approximations, the second-order derivatives in methods based on a quadratic model may provide their numerical stability [<sup>17, 18</sup>].

## NUMERICAL EXPERIMENTS

A computer experimentation which is aimed at comparing the relative performance of methods was carried out for:

1. The Fridman method

$$x_{k+1} = x_k - \alpha_k g(x_k), g(x) = J(x)^{\mathrm{T}} F(x), \ \alpha_k = \|F(x_k)\|^2 / \|g(x_k)\|^2.$$

2. The two-step gradient method of Măruşter [<sup>19</sup>]

$$x_{k+1} = x_k - \alpha_k g(x_k) + \beta_k J(x_k)^{-1} J(x_k) g(x_k),$$

where the parameters  $\alpha_k$  and  $\beta_k$  are determined by the system

$$\beta_{k} \| J(x_{k}) g(x_{k}) \|^{2} + \alpha_{k} \| g(x_{k}) \|^{2} + \| F(x_{k}) \| = 0,$$

$$\beta_{k} \left\| J(x_{k})^{T} J(x_{k}) g(x_{k}) \right\|^{2} + \alpha_{k} \left\| J(x_{k}) g(x_{k}) \right\|^{2} + \left\| g(x_{k}) \right\|^{2} = 0.$$

3. The modified Newton method.

4. The Newton method a) with derivatives b) with finite-differences.

- 5. The method defined by (4) with  $W_k = A_{k-1}$ .
- 6. The Gauss-Newton method.

7. The cubically convergent method [<sup>20</sup>]

$$\begin{aligned} x_{k+1} &= y_k - \left[F\left(2y_k - x_k; x_k\right)\right]^{-1} F\left(y_k\right), \\ y_k &= x_k - \left[F\left(2y_{k-1} - x_{k-1}; x_{k-1}\right)\right]^{-1} F\left(x_k\right), \end{aligned}$$

where f(.;.) denotes a finite-difference approximation to J(x).

8. The combination of the latter with the Newton method with finitedifferences.

9. The Kogan method [<sup>21</sup>]

$$x_{k+1} = x_k - \left[ F'\left(x_k - \frac{1}{2}\Gamma_k(x_k)\right) \right]^{-1} F(x_k),$$

where  $\Gamma_{k} = [F'(x_{k})]^{-1}$ .

This paper is not intended to give a numerical comparative study of methods under consideration with the other related methods. For the sake of illustration only some results on their numerical behavior are presented in the table.

Method		1	2	3	4		5	6	7	8	9
Test	Di- men- sion		-	Philip:	а	Ь		5xerr			
1	2	k>300	4	3	3	3	3	2	2	2	2
2	4	k>300	k>300	10>300	32	31	32	k>300	-	38	21
3	2	<i>k</i> >300	26	-	13	19	13	14	-	12	9
4	4	k>300	k>300	k>300	15	47	15	15	-	17	28
5	3	k>300	134		11	25055	11	11	-	10	6
6	6	k>300	<i>k</i> >300	k>300	13	45	13	13	-	14	9
7	5	25	6	210 Bend	6	6	6	6	149	5	4
8	10	k>300	k>300	1. Secon	91	80	-1-	-		92	<i>k</i> >300
9	10	k>300	<i>k</i> >300	6	4	3	4	4	3	3	3
10	10	16	7	6	4	4	4	4	3	3	3
11	10	<i>k</i> >300	k>300	<i>k</i> >300	8	6	8	8.0.9	32.9	6	7
12	10	15	27	k>300	15	20	15	15	-	16	10
13	10	114	33	22	5	5	5	6		5	4
14	10	35	15	60	7	6	7	7	14	6	5

The set of test problems containing 14 problems for systems of nonlinear equations of Argonne National Laboratory plus Freudenstein and Roth function and Box three-dimensional function for nonlinear least squares was taken from  $[^{22}]$ , and the same starting points were used as in  $[^{22}]$ .

All the problems were run on a EC 1060 computer under a FORTRAN-IV compiler. The calculus was performed in double precision and stopped  $||x_{k+1} - x_k|| \le 10^{-9}$ . Termination of the routine also occurred when the number of iterations exceeded the given maximum value 300. The table gives iteration number k at which the presigned accuracy was achieved and "-" denotes the failure (nonconvergence).

In conclusion, we discussed the methods mainly on the theoretical basis. One reason for doing so is that mathematical properties exhibit indisputable features of the algorithms; in contrast to the computer experiments their numerical results presented in papers and books often show contradictory aspects. Also, mathematical aspects fix limits of what can be expected from the use of an algorithm.

## ACKNOWLEDGEMENT

The author wishes to express his special thanks to Prof. G. Vainikko for very helpful discussions.

#### REFERENCES

- Elden, L. A note on weighted pseudoinverses with application to the regularization of Fredholm integral equations of the first kind: Report Lith-MAT-R-1975-11, Department of Mathematics, Linköping University (Sweden), 1975.
- 2. Ваарманн О., Ломп М. Изв. АН ЭССР. Физ. Матем., 1982, 31, 4, 410-417.
- 3. Vaarmann, O. Eesti NSV TA Toim. Füüs.-Matem., 1989, 38, 2, 146-153.
- Albert, A. Regression and the Moor-Penrose pseudoinverse. Academic Press, New York-London, 1972.
- 5. Deuflhard, P. & Heindl, G. SIAM J. Numer. Anal., 1979, 16, 1-10.
- 6. Vaarmann, 0. Proc. Estonian Acad. Sci. Phys. Math., 1991, 40, 2, 99-104.
- Vaarmann, O. In: Tikhonov, A. N. (ed.). Ill-Posed Problems in Natural Sciences. Utrecht: VSP/ Moscow: TVP Sci. Publ. 1992, 191-201.
- Deuflhard, P. & Apostolescu, V. Preprint N° 51. Universität Heidelberg, Institute f
  ür Angewandte Mathematik. January, 1980.
- 9. Nash, S & Sofer, A. Oper. Res. Lett., 1990, 4, 4, 219-221.
- 10. Ashby, S. F., Manteuffel, T. A. & Otto, J. S. SIAM J. Sci. Comput., 1992, 13, 1, 1-29.
- 11. Ваарманн О. Изв. АН ЭССР. Физ. Матем., 1971, 20, 4, 386-394.
- Lindqvist, S.-g. Programs for nonlinear least squares problems a user's guide. Computing centre, the Swedish University of Åbo (Finland), Report N
  <sup>o</sup> 9, Jan. 1981, 19.
- 13. Deuflhard, P. Numer. Math., 1974, 22, 289-315.
- 14. Wedin, P.-A. Bit, 1973, 13, 217-232.
- 15. Stewart, G. W. SIAM Rev., 1977, 19, 4, 634-662.
- 16. Пугачев Б. П. ЖВМ и МФ, 1978, 18, 6, 1593-1595.
- 17. Ehle, G. & Schwetlick, H. SIAM J. Numer. Anal., 1976, 13, 3, 433-443.
- 18. Schnabel, R. B. & Frank, P. D. SIAM J. Numer. Anal., 1984, 21, 5, 815-843.

 Măruşter, St. On the two-step gradient method for nonlinear equations. Seminarul De Informatic Si Analiz Numerica. Timişoara, 1985, 20.

20. Ваарманн О., Полль В. Изв. АН ЭССР. Физ. Матем., 1977, 26, 2, 123-127.

21. Коган Т. И. Сиб. мат. журн., 1967, 8, 4, 958-960.

22. Moré, J. J., Garbow, B. S. & Hillstrom, K. E. ACM Trans. Math. Soft., 1987, 7, 1, 17-41.

## EBASTABIILSETE MITTELINEAARSETE VÕRRANDISÜSTEEMIDE LAHENDAMISEST

#### Otu VAARMANN

On vaadeldud *m* võrrandi ja *n* otsitavaga  $(m \ge n)$  mittelineaarse võrrandisüsteemi lahendamist juhul, kui selle Jacobi maatriksi konditsiooniarv võib olla väga suur. On uuritud sellise ülesande lahendusmeetodeid, mis põhinevad: 1) Jacobi maatriksi kaalutud pseudopöördmaatriksi arvutamisel; 2) vastavate lineaarsete võrrandite ligikaudsel lahendamisel mõne iteratsioonimeetodi abil. On esitatud saadud teoreetilisi uurimistulemusi illustreerivad arvutustulemused.

## О РЕШЕНИИ ПЛОХО ОБУСЛОВЛЕННЫХ НЕЛИНЕЙНЫХ СИСТЕМ УРАВНЕНИЙ

# Оту ВААРМАНН

Для решения плохо обусловленных систем нелинейных уравнений с m уравнениями и n неизвестными ( $m \ge n$ ) применяются итерационные методы, основанные на вычислении взвешенной псевдоинверсии для якобиана и на приближенном решении соответствущих линейных уравнений с помощью некоторого итерационного метода. Доказаны теоремы о сходимости рассматриваемых методов. Для иллюстрации полученных теоретических выводов приводятся результаты численных экспериментов.

usually, employed is that of needecting the enor due to samp