

The use of inner preconditioned conjugate gradient iterations in large sparse nonlinear optimization problems

©I. E. Kaporin

Computing Center of Russian Academy of Sciences, Vavilova 40,
Moscow 119991, Russia
E-mail: kaporin@sci.kun.nl, kaporin@ccas.ru

1 Introduction

In our presentation, we consider the problem of the (sub)optimum termination of inner linear Preconditioned Conjugate Gradient (PCG) iterations used within Inexact Newton nonlinear solvers [2, 4, 3].

In general, we show that for certain solution error measure ϵ_k , where k is the number of nonlinear (outer) iteration, one has

$$\epsilon_{k+1} \leq (1 - \tau\vartheta_k(s_k))\epsilon_k \quad (1)$$

where $0 < \tau \leq 1$ is a certain characterization of the problem nonlinearity ($\tau = 1$ when the problem is linear), and $0 < \vartheta_k(s_k) \leq 1$ is a special error measure for the solution of some inner linear equation using s_k iterations of the PCG method (this linear error is zero whenever $\vartheta_k = 1$). We present a simple practical rule for choosing s_k which does not assume the knowledge of τ and makes it possible to avoid making redundant PCG iterations within each nonlinear step.

The general outline is as follows. Let the nonlinear convergence is measured by

$$\epsilon_m \leq \varepsilon\epsilon_0, \quad 0 < \varepsilon \ll 1.$$

Then (1) yields the following sufficient condition for the nonlinear iterations to be converged:

$$\sum_{j=0}^{m-1} \vartheta_j(s_j) \approx \tau^{-1} \log(\varepsilon^{-1}).$$

Let then \mathcal{P} and \mathcal{I} be the computational costs of one outer iteration (including the linear solver initialization cost) and the computational cost per one PCG

iteration, respectively. Hence, the total computational cost is estimated as

$$\begin{aligned} \sum_{j=0}^{m-1} (\mathcal{P} + \mathcal{I}s_j) &\approx \left(\sum_{j=0}^{m-1} (\mathcal{P} + \mathcal{I}s_j) / \sum_{j=0}^{m-1} \vartheta_j(s_j) \right) \tau^{-1} \log(\varepsilon^{-1}) \\ &\leq \left(\max_{0 \leq j \leq m-1} (\mathcal{P} + \mathcal{I}s_j) / \vartheta_j(s_j) \right) \tau^{-1} \log(\varepsilon^{-1}). \end{aligned}$$

Hence, the stopping criterion can now be formulated in terms of choosing the inner iteration numbers s_j providing for a reasonably small value of each ratio

$$\psi(s) = (\mathcal{P} + \mathcal{I}s) / \vartheta_j(s).$$

Since the quantity $\vartheta_j(s)$ can be readily evaluated at any s -th inner iteration, cf.[3], one can stop at $s = s_j$ where the first local minimum of $\psi_j(s)$ is attained.

2 Nonlinear Least Squares Problem

Let F be a differentiable mapping

$$F : R^n \rightarrow R^N, N \geq n,$$

and let the nonlinear least squares problem

$$u_* = \arg \inf \|Fu\|$$

be approximately solved using inexact Gauss-Newton iterations. Each iteration of the method can be presented as follows: given an approximation u , one calculates the matrix

$$A = F'(u)^T F'(u),$$

the right hand side

$$b = -F'(u)^T Fu,$$

then performs some inexact solution of related linear system, i.e.,

$$Ax \approx b$$

and then obtains the next approximation as

$$u_+ = u + \tau x$$

choosing some proper steplength $\tau > 0$. It was shown in [4] that if the approximate Gauss-Newton direction x is properly scaled, that is,

$$b^T x = x^T A x, \quad (2)$$

then the estimate

$$\|F(u_+)\| \leq \sqrt{1 - \tau\theta^2} \|Fu\|, \quad (3)$$

where

$$\theta = |b^T x| / (\|Fu\| \sqrt{x^T A x})$$

holds for some positive $\tau \leq 1$, the maximum admissible value of which is considered as the key characterization of the local nonlinearity of the mapping F . Note that

$$\theta = \cos[Fu, F'(u)x] = \sqrt{x^T A x} / \|Fu\|$$

where $\cos[p, q]$ is the cosine of the acute angle between the m -vectors p and q , and the last equality holds by the scaling condition (2). Clearly, the equation (3) can readily be rewritten in the form (1).

In [4], the quantity τ was called *the limiting stepsize from u along the normalized direction x* and defined as the largest positive number satisfying

$$2\|Fu\| \|F(u + \alpha x) - Fu - \alpha F'(u)x\| \leq \alpha(1 - \alpha) \|F'(u)x\|$$

for all $0 < \alpha \leq \tau$.

If F is a linear mapping, one can simply take $\tau = 1$, and it makes sense to solve the linear problem for x to full precision. However, for the nonlinear problems the case of $\tau < 1$ always takes place, so one should specify a proper termination rule for the PCG iterations when solving for x .

We will describe here a criterion for stopping PCG iterations which is independent of any characterization of nonlinearity and is expressed in terms of scalar coefficients involved in PCG recursions and the ratio of the PCG startup cost to a regular PCG iteration cost. Our strategy is aimed towards *maximization of θ* rather than (standard) minimization of the residual $\|b - Ax\|$. It worth noting here that, in general, the relative residual $\|b - Ax\|/\|b\|$ is rather loosely related to the above quantity θ .

3 The conjugate gradient iterations

Let us recall the PCG algorithm. The PCG iterations [1] for the solution of the problem $Ax = b$ can be written as follows:

$$r_0 = b - Ax_0,$$

```

 $p_0 = Cr_0;$ 
for  $i = 0, 1, \dots :$ 
 $\alpha_i = r_i^T Cr_i / p_i^T Ap_i,$ 
 $x_{i+1} = x_i + p_i \alpha_i,$ 
 $r_{i+1} = r_i - Ap_i \alpha_i,$ 
 $\beta_i = r_{i+1}^T Cr_{i+1} / r_i^T Cr_i,$ 
 $p_{i+1} = Cr_{i+1} + p_i \beta_i.$ 

```

Here C is a properly chosen SPD preconditioning matrix, which should approximate, in some sense, the matrix A^{-1} . The choice of the matrix C is subject to the requirement that a vector $w = Cr$ be easily calculated for any r . For instance, one of the best choices is the approximate Cholesky preconditioning, where $C = (U^T U)^{-1}$ and $U^T U \approx A$ with the upper triangular matrix U being much sparser than the exact Cholesky factor of A , cf.[5] and references therein.

Fortunately, if $x_0 = 0$ in the above method, then one can easily see (e.g., using the A -orthogonality property of p_i , cf.[3]) that each iterate x_i obtained in the PCG iterations applied to $Ax = b$ satisfies the above scaling condition (1) and, moreover, the equality

$$x_i^T Ax_i = \sum_{j=0}^{i-1} \omega_j$$

holds, where the quantities

$$\omega_i = (r_i^T H r_i)^2 / p_i^T A p_i$$

are readily available from the scalar products calculated in the course of the PCG iterations. Therefore, the above mentioned quantity $\theta = \sqrt{x^T Ax} / \|Fu\|$ can be calculated directly using

$$\theta = \|Fu\|^{-1} \left(\sum_{j=0}^{i-1} \omega_j \right)^{1/2}, \quad i = 0, 1, \dots, n-1,$$

which makes it possible to develop our new PCG stopping criterion.

The scaling relationship (2) can be proved as follows. Let x be obtained after s iterations of the PCG method with zero initial guess. Therefore,

$$x \in K_s = \text{span}\{Cb, CACb, \dots, (CA)^{s-1}Cb\},$$

and, by the PCG optimality property,

$$x = \arg \min_{x \in K_s} (b - Ax)^T A^{-1} (b - Ax).$$

Since $\alpha x \in K_s$ for any scalar α , one gets

$$(b - \alpha Ax)^T A^{-1} (b - \alpha Ax) \geq (b - Ax)^T A^{-1} (b - Ax),$$

which readily gives, with $\alpha = x^T b / x^T Ax$, the inequality

$$0 \geq (-x^T b + x^T Ax)^2,$$

which readily yields the required scaling condition.

Furthermore, using the well known techniques developed for the estimation of the PCG iteration error [1], one gets

$$(b - Ax)^T A^{-1} (b - Ax) / b^T A^{-1} b \leq 1 / \cosh^2 \left(\frac{2s}{\sqrt{\kappa}} \right)$$

where

$$\kappa = \text{cond}(B^{-1}A)$$

which, by the scaling condition, gives

$$\theta^2 = x^T Ax / b^T A^{-1} b \geq \tanh^2 \left(\frac{2s}{\sqrt{\kappa}} \right). \quad (4)$$

Hence, $0 < \theta < 1$ and $\theta \rightarrow 1$ as the PCG iteration number s grows, the faster the better the preconditioner C .

4 Using the new inner iteration stopping criterion

Relying on the above expression for θ , one can try to find a proper balance between the costs of initializing and performing the inner PCG iterations and the acceleration obtained at the outer nonlinear iterations due to larger values of θ .

The inner iteration stopping criterion based on maximization of θ can be constructed as follows. Let m nonlinear iterations be performed, and let τ be the lower bound for the stepsizes used. Assuming that the nonlinear convergence is measured by

$$\|F(u_m)\| \leq \varepsilon \|F(u_0)\|,$$

one can therefore find that a sufficient condition for the nonlinear iterations to be converged can be taken as

$$\sum_{j=1}^m \theta_j^2 \approx \tau^{-1} \log(\varepsilon^{-2}).$$

Let then \mathcal{P} be the costs of the outer iteration (including at least the generation of b , A , and the preconditioner H) and \mathcal{I} be the costs per one PCG iteration (determined mainly by the costs of matrix-vector multiplications with A and H). The total computational cost (i.e., the running time) is estimated as

$$\begin{aligned} \sum_{j=1}^m (\mathcal{P} + \mathcal{I}s_j) &\approx \left(\sum_{j=1}^m (\mathcal{P} + \mathcal{I}s_j) / \sum_{j=1}^m \theta_j^2 \right) \tau^{-1} \log(\varepsilon^{-2}) \\ &\leq \left(\max_{1 \leq j \leq m} (\mathcal{P} + \mathcal{I}s_j) \theta_j^{-2} \right) \tau^{-1} \log(\varepsilon^{-2}), \end{aligned}$$

where s_j is the number of inner iterations at the j -th outer iteration step. The stopping criterion can now be formulated in terms of choosing the inner iteration numbers s_j providing for a reasonably small value of each ratio

$$(\mathcal{P} + \mathcal{I}s_j) \theta_j^{-2} = \|Fu_j\|^2 (\mathcal{P} + \mathcal{I}s_j) / (\omega_0 + \dots + \omega_{s_j}).$$

In our experiments we choose the value s_j for which the increase of this ratio occurred for the first time, i.e. the iterations were performed until the condition

$$s \geq -(\mathcal{P}/\mathcal{I}) + (\omega_0 + \dots + \omega_s) / \omega_s \quad (5)$$

holds true. In view that $\omega_0 + \dots + \omega_s$ is bounded from above by the squared A -norm of the solution x , and ω_s tend to decrease as the PCG iterations progress, one can expect rather early termination of the inner iterations. The above inner iteration stopping criterion was successfully used in numerical experiments reported in [2].

5 The case of unconstrained minimization problem

Quite similarly, the new PCG stopping criterion can be applied to the construction of algorithms for unconstrained minimization. In particular, rather hard-to-solve minimization problems arising in global untangling of computational

grids via continuation technique, see [6], were successfully solved there by the application of related inexact Newton-like minimization procedure. Next we present some theoretical analysis which has been omitted in [6].

Let a differentiable functional

$$\varphi(u) : R^n \rightarrow R^1$$

be bounded from below, have gradient $g(u) \in R^n$, and

$$|\varphi(u + h) - \varphi(u) - h^T g(u)| \leq \frac{\gamma}{2} h^T A h \quad (6)$$

for certain symmetric positive definite $n \times n$ matrix $A = A(u)$. Here, A is supposed to be given explicitly, while the exact knowledge of γ is not necessary.

The minimizer u_* of φ is approximated by the iterates u which are updated as follows:

Step 1. Compute $g = g(u)$ and check the convergence, i.e., if

$$\|g\| \leq \varepsilon,$$

then quit;

Step 2. Compute $A = A(u)$ and find $x \approx A^{-1}g$ such that the scaling condition

$$-x^T g = x^T Ax \quad (7)$$

holds and the quantity

$$\theta^2 = x^T Ax / g^T A^{-1}g, \quad 0 \leq \theta \leq 1,$$

is sufficiently large;

Step 3. Find

$$\alpha = \arg \min_{\beta > 0} \varphi(u + \beta x),$$

set

$$u_+ = u + \alpha x$$

then $u := u_+$ and go to Step 1.

In order to specify further implementation details, let us consider this scheme more closely.

Let x at Step 2 be obtained after s iterations of the (preconditioned) conjugate gradients with zero initial guess applied to the linear system $Ax = -g$. As follows from (2) of the required scaling condition (7) holds by $b = -g$.

Next we will estimate the reduction in the functional value attained by the descent along the direction x . For any stepsize $\beta > 0$ one has, by (6) and (7),

$$\begin{aligned}\varphi(u + \beta x) &= \varphi(u) + \beta x^T g + (\varphi(u + \beta x) - \varphi(u) - \beta x^T g) \\ &\leq \varphi(u) - \beta x^T Ax + \frac{\gamma}{2} \beta^2 x^T Ax \\ &= \varphi(u) - \left(\beta - \frac{\gamma}{2} \beta^2\right) \theta^2 g^T A^{-1} g.\end{aligned}$$

One gets then

$$\varphi(u_+) = \varphi(u + \alpha x) = \min_{\beta > 0} \varphi(u + \beta x) \leq \varphi(u + \frac{1}{\gamma} x) \leq \varphi(u) - \frac{\theta^2}{2\gamma} g^T A^{-1} g.$$

Hence we obtain

$$\varphi(u_+) = \varphi(u) - \frac{\theta^2}{2\gamma} g^T A^{-1} g,$$

and therefore $\lim_{i \rightarrow \infty} \|g_i\|_{A_i^{-1}} = 0$, where i is the outer iteration number. Moreover, letting u_* be the minimizer of $\varphi(u)$, one has

$$\varphi(u_+) - \varphi(u_*) \leq \left(1 - \frac{\theta^2 g^T A^{-1} g}{2\gamma(\varphi(u) - \varphi(u_*))}\right) (\varphi(u) - \varphi(u_*)). \quad (8)$$

Choosing τ as the lower bound for $g^T A^{-1} g / (2\gamma(\varphi(u) - \varphi(u_*)))$ over all outer iterations, one can write the error equation (8) in the form (1). Therefore, in order to provide the suboptimum decrease in $\varphi(u) - \varphi(u_*)$ at each outer iteration, one can use the same stopping criterion (3) for inner PCG iterations, where $s = s_i$ is the number of PCG iterations, ω_k are calculated from the PCG scalar products in the same way as earlier, and \mathcal{P}/\mathcal{I} is the ratio of the computational costs implied by evaluation of $\varphi(u)$, $g(u)$, $A(u)$ and the preconditioner C to the computational costs of one inner PCG iteration.

6 Acknowledgement

This work has been partially supported by the Netherlands Technology Foundation (STW) grant NNS.4683 and by the National Science Foundation grant NWO 047.008.007.

References

- [1] O. Axelsson. A class of iterative methods for finite element equations. *Computer Meth. Appl. Mech. Engrg.*, 9, 123-137, 1976.
- [2] O. Axelsson and I. E. Kaporin. Minimum residual adaptive multilevel procedure for the finite element solution of nonlinear stationary problems. *SIAM J. Numer. Anal.*, 35, 1213-1229, 1998.
- [3] O. Axelsson and I. E. Kaporin. Error norm estimation and stopping criteria in preconditioned conjugate gradient iterations, *Numer. Linear Algebra Appl.* 8, 265-286, 2001.
- [4] I. E. Kaporin and O. Axelsson. On a class of nonlinear equation solvers based on the residual norm reduction over a sequence of affine subspaces. *SIAM J. Sci. Comput.*, 16, 228-249, 1994.
- [5] I. E. Kaporin. High quality preconditioning of a general symmetric positive definite matrix based on its $U^T U + U^T R + R^T U$ -decomposition. *Numer. Linear Algebra Appl.* 5, 483-509, 1998.
- [6] V. A. Garanzha and I. E. Kaporin. Regularization of barrier variational method for constructing computational grids. (*Russian*) *Zh. Vychisl. Mat. i Mat. Fiz.* 39, 1489-1503, 1999.