Harlan Crowder Philip Wolfe

Linear Convergence of the Conjugate Gradient Method

Abstract: There are two procedures for applying the method of conjugate gradients to the problem of minimizing a convex, nonlinear function: the "continued" method, and the "restarted" method in which all the data except the best previous point are discarded, and the procedure is begun anew from that point. It is demonstrated by example that in the absence of the standard initial starting condition on a quadratic function, the continued conjugate gradient method will converge to the solution no better than linearly. Furthermore, it is shown that for a general nonlinear function, the nonrestarted conjugate gradient method converges no worse than linearly.

Introduction

The method of conjugate gradients was originally introduced as an alternative to Gaussian elimination for the solution of large systems of linear equations on a computer [1]. Although elimination remains the preferred method for the standard linear equation problem, the conjugate gradient method does find application in such areas as the solution of certain large, sparse systems of linear equations, certain infinite-dimensional linear problems, and the minimization of general nonlinear functions. The latter point is the topic of this study.

In particular, suppose we try to minimize some realvalued, nonlinear function f(x), where $x \in E^n$. We want f(x) to be differentiable, and preferably convex. By applying an algorithm which employs the conjugate gradient method (see the following section), a sequence is generated that converges to the minimum of f(x) [the minimum being defined as $f(x^*)$ such that $f(x^*) \le f(x)$ for all x sufficiently close to x^* ; if the function is convex, then $f(x^*) \le f(x)$ for all $x \in E^n$]. For a particular subclass of unconstrained minimization problems, the conjugate gradient method not only converges to the solution but also terminates at the solution in a finite number of steps. Specifically, let f(x) be a quadratic function of the form $f(x) = c + p^{T}x + \frac{1}{2}x^{T}Qx$, where c is a scalar constant, p is an nth order column vector, and Q is an nth order symmetric and positive semidefinite matrix. Thus f(x) is convex. Then the conjugate gradient method will terminate at the minimum in at most n steps, provided the standard initial starting condition, which is defined in the second section, is used.

It should be noted that this standard initial starting condition (or *standard start*) is of primary importance,

for if this initial condition is not applied to the quadratic problem mentioned above, then the conjugate gradient method will not terminate at the solution, although the method will converge to the solution as it does when applied to a nonquadratic or general nonlinear function. We can thus conclude that the application of the conjugate gradient method to a quadratic function in the absence of the standard start behaves analogously to the application of the method to a general nonlinear function with the standard start.

There are two main ways which the conjugate gradient method can be applied to a general nonlinear function (or a quadratic function without the standard start): the continued method, in which the same recursive procedure is used throughout; and the restarted method, in which, after a certain number of steps, most of the iterative data are discarded and the procedure is begun afresh from the best solution thus far obtained. In the present study, we find that the speed of convergence to the solution of the continued conjugate gradient method is linear, as is the case in a geometric series. Since the restarted conjugate method has superlinear convergence [2] for a wide class of problems, it is evident that the continued method should not be used for general nonlinear functions.

Method of conjugate gradients

Let f(x) be a quadratic function as defined in the previous section. A statement of the conjugate gradient method for this function is:

Given
$$x_0$$
, let $d_0 = -\nabla f(x_0)$.
For $k \ge 0$, given x_k and d_k , let $x_{k+1} = x_k + t_k d_k$,

431

where t_k is the value of t minimizing $f(x_k + td_k)$.

If
$$\nabla f(x_{k+1}) = 0$$
, stop. Otherwise, let

$$d_{k+1} = -\nabla f(x_{k+1}) + s_k d_k,$$

where s_k is chosen so that $d_{k+1}^T Q d_k = 0$.

Let $g_k = \nabla f(x_k)$ for all k. Noting that

$$g_{k+1} = \nabla f(x_{k+1})$$

$$= p^{T} + Q(x_k + t_k d_k)$$

$$= p^{T} + Qx_k + t_k Q d_k$$

$$= g_k + t_k Q d_k,$$

we can write the recursion as

$$d_0 = -g_0; (1)$$

For $k \geq 0$,

$$g_{k+1} = g_k + t_k Q d_k$$

$$d_{k+1} = -g_{k+1} + s_k d_k \tag{2}$$

where

$$t_k = -g_k^T d_k / d_k^T Q d_k \tag{3}$$

and

$$s_k = g_{k+1}^T Q d_k / d_k^T Q d_k \tag{4}$$

or

$$s_k = g_{k+1}^T (g_{k+1} - g_k) / d_k^T (g_{k+1} - g_k).$$
 (5)

The formula (5) for s_k is essentially formula (3:2b) of Hestenes and Stiefel, rather than the more commonly used formula (3:1e), which in our notation is $s_k = g_{k+1}^2/g_k^2$. As they subsequently show, the former gives better protection against the accumulation of roundoff error. More importantly, it ensures that $d_{k+1}^TQd_k=0$ for each k, independently of whether the other steps have been carried out accurately, which the latter formula does not. If all the needed relations do hold accurately, it can be shown that the successive directions d_0 , d_1 , \cdots are all linearly independent and conjugate (that is, $d_j^TQd_k=0$ for $j\neq k$), and that x_k minimizes the function f on the affine set passing through x_0 and spanned by d_0 , d_1 , \cdots , d_{k-1} . Consequently the procedure must terminate with $g_k=0$ for some $k\leq n$.

The problem

It is important to note that both the starting condition (1) and the determinations (3) and (4) of the coefficients t_k and s_k must be observed precisely in order that the above termination ensues; it cannot be shown otherwise. Indeed, failure to choose a standard start—one in which

 d_0 is parallel to g_0 – makes it impossible to retain the conjugacy relation $d_i^T Q d_k = 0$ for $|j - k| \neq 1$ using formulas of the type of (2). (We have seen this fact overlooked in some reports in the literature, leading to an overestimate of the convergence rate of the method.) Since, however, the procedure is almost invariably used under circumstances in which the condition $g_k = 0$ cannot be precisely met-with the quadratic problem in an environment of roundoff error, and, more significantly, in extensions of the method to nonquadratic problems, such as that due to Fletcher and Reeves [3] – provision for continuing after the nth step must be made. It has generally been recognized as good practice to restart the procedure after n(or possibly n + 1 or n + 2) iterations; that is, to begin all over again, using the latest point x_k found as the new x_0 and thus rebuild a new set of conjugate directions.

The purpose of the present study was to determine whether restarting was, in fact, necessary, or whether the procedure could be continued indefinitely without restarting and not suffer. We have concluded that restarting is necessary for quick convergence. Indeed, we have an example (for n = 3) of a quadratic problem which shows that convergence can be no better than linear when a nonstandard start is used. A standard start or restart would, of course, cause termination in at most three iterations.

Example: Convergence is at best linear

We have run about 50 steps of the continued conjugate gradient method as defined by Eqs. (2) to (4) on each of some 100 quadratic, three-variable problems, examining graphically the ratios $f(x_{k+1})/f(x_k)$ of successive values of the function $f(x) = \frac{1}{2} x^T Q x$. In about half of the trials, Q was the diagonal matrix, the eigenvalues of which are (0.1, 1, 1); the starting vectors g_0 and d_0 were chosen randomly. In every case the ratios, which at first seemed randomly scattered between 0 and 1, were found to lie in a rather definitely marked interval [a,b] with 0 < a <b < 1. In many cases it appeared that something very much like a sine curve having a period between three and five steps could be fitted to the set of successive ratios. After considerable experimentation with the starting data, we found an example in which the ratios were constant. The other data of the procedure then exhibited a remarkable periodicity, and the discovery of simple relationships among these led to the following example:

$$let Q = \begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$g_0 = (1, -\sqrt{5}, 0)^T / \sqrt{6}$$

$$d_0 = (-10\sqrt{5}, 14, -3\sqrt{6})^T/4\sqrt{30}$$

One step of the method given by Eqs. (2) to (4) is

$$g_{k+1} = g_k + t_k Q d_k$$

$$d_{k+1} = -g_{k+1} + s_k d_k$$
.

In our case $t_k = 8/5$, and $s_k = 9/25$ for all k. Furthermore, the relations $g_{k+1} = rRg_k$ and $d_{k+1} = rRd_k$ hold for all k where r = 3/5 and R is the orthogonal matrix

$$R = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1/5 & -(2\sqrt{6})/5 \\ 0 & (2\sqrt{6})/5 & -1/5 \end{bmatrix}.$$

Thus $g_k = (rR)^k g_0$ and $d_k = (rR)^k d_0$ for all k. Each successive application of the matrix rR rotates the gradient and the direction through an angle $\arccos{(-1/5)}$ around the long axis of the three-dimensional ellipsoid $x^TQx = 1$, and diminishes both of these vectors in magnitude by the factor r = 3/5. Thus the ratio $f(x_{k+1})/f(x_k)$ is 9/25 for all k.

Theorem: Convergence is at worst linear

To bound the rate of convergence of the nonrestarted conjugate gradient method from both sides, we will show that its convergence is at worst linear.

Let $f(x) = \frac{1}{2} x^T Q x$. We can always transform the original problem so that it has this form. Since $g = \nabla f(x) = Qx$, $f(x) = \frac{1}{2} g^T Q^{-1} g$. The minimal of f along any line x + td is given by $t = \hat{t} = -g^T d/d^T Q d$ (compare with formula (3), suppressing k). Setting $x_+ = x + \hat{t} d$ and $g_+ = \nabla f(x_+)$, we have

$$2f(x_{+}) = g_{+}^{T}Qg_{+} = (g + \hat{t}Qd)^{T}Q^{-1}(g + \hat{t}Qd)$$
$$= gQ^{-1}g - (g^{T}d)^{2}/d^{T}Qd.$$

We consider two cases:

1) d = -g; that is, the step is an ordinary steepest descent step.

Then

$$2f(x_{+}) = g^{T}Q^{-1}g - (g^{T}g)^{2}/g^{T}Qg.$$

2) The point x was obtained by minimizing f along some line having the direction c, whence $g^Tc=0$, and then the direction d was obtained as in formulas (2) and (4), so that d=-g+sc and $d^TQc=0$.

Then $g^T d = 0 - g^T g$ and

$$d^{T}Qd = -g^{T}Qd = -g^{T}(-Qg + sQc) = g^{T}Qg - sg^{T}Qc$$
$$= g^{T}Qg - (g^{T}Qc)^{2}/c^{T}Qc.$$

We see that $d^TQd \leq g^TQg$.

Since in case 2) $2f(x_+) = g^T Q^{-1}g - (g^T g)^2/d^T Q d$, the resulting value of f is no higher in case 2) than in case 1);

the fact that the direction d was obtained by conjugating -g with respect to the previous direction, rather than taking it to be -g itself, has not hurt. Thus each step of the continued conjugate gradient method decreases the function at least as much as would one step of steepest descent taken at the same point. The inequality

$$f(x_{k+1})/f(x_k) \le [(A-1)/(A+1)]^2$$

is known to hold for steepest descent, where A is the condition number of the matrix Q, namely, the ratio of the largest to smallest eigenvalue. It follows that the inequality also holds for the conjugate gradient method, so that its convergence is at worst linear.

Conclusion

We have used a small, rather simple quadratic function to show that, in the absence of the standard initial starting condition for the conjugate gradient method, the convergence of the generated sequence will be no better than linear. We have also shown that if the conjugate gradient method is applied to a general nonlinear function, and if the procedure is not restarted every n or so steps, then the convergence rate to the solution is no worse than linear. This characterization of the rate of convergence of the conjugate gradient procedure should be used when considering whether or not to apply the method to a specific nonlinear optimization problem. In addition, our result should benefit anyone attempting to understand the fundamental aspects of the conjugate gradient method.

Acknowledgment

This work was supported in part by the Office of Naval Research under Contract No. N00014-71-C-0112.

References

- M. R. Hestenes and E. Stiefel, "Methods of Conjugate Gradients for Solving Linear Systems," J. Res. NBS 49, 409-436 (1952).
- Garth P. McCormick, and Klaus Ritter, "On the Convergence and Rate of Convergence of the Conjugate Gradient Method," MRC Technical Summary Report No. 1118, June, 1971, Mathematics Research Center, The University of Wisconsin.
- 3. R. Fletcher and C. M. Reeves, "Function Minimization by Conjugate Gradients," *Computer J.* 7, 149-154 (1964-5).

Received November 3, 1971

The authors are located at the IBM Thomas J. Watson Research Center, Yorktown Heights, New York 10598.