

On the Barzilai and Borwein choice of steplength for the gradient method

MARCOS RAYDAN

Department of Mathematics, University of Kentucky, Lexington, Kentucky, 40506

[Received 30 October 1991 and in revised form 17 June 1992]

In a recent paper, Barzilai and Borwein presented a new choice of steplength for the gradient method. Their choice does not guarantee descent in the objective function and greatly speeds up the convergence of the method. They presented a convergence analysis of their method only in the two-dimensional quadratic case. We establish the convergence of the Barzilai and Borwein gradient method when applied to the minimization of a strictly convex quadratic function of any number of variables.

1. Introduction

In a recent paper, Barzilai & Borwein (1988) presented a new choice of steplength for the gradient method. Their choice does not guarantee descent in the objective function, requires less computational work than the optimum choice of the steepest descent method, and greatly speeds up the convergence of the method (see also Fletcher (1990)). Barzilai & Borwein (1988) presented a convergence analysis of their method only in the two-dimensional quadratic case. In this work, we establish the convergence of the Barzilai and Borwein gradient method when applied to the minimization of a strictly convex quadratic function of any number of variables.

The Barzilai and Borwein method for the unconstrained minimization of a differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is defined by

$$x_{k+1} = x_k - \frac{1}{\alpha_k} g_k, \quad (1)$$

where g_k is the gradient vector of f at x_k and the scalar α_k is given by either

$$\alpha_k = \frac{s'_{k-1} y_{k-1}}{s'_{k-1} s_{k-1}} \quad (2)$$

or

$$\alpha_k = \frac{y'_{k-1} y_{k-1}}{s'_{k-1} y_{k-1}}, \quad (3)$$

where $s_{k-1} = x_k - x_{k-1}$ and $y_{k-1} = g_k - g_{k-1}$. When $f(x) = \frac{1}{2}x'Ax - b'x + c$ is a quadratic function and A is a symmetric positive definite (SPD) matrix, then α_k in

(2) and (3) becomes

$$\alpha_k = \frac{s'_{k-1} A s_{k-1}}{s'_{k-1} s_{k-1}}, \quad (4)$$

$$\alpha_k = \frac{s'_{k-1} A^2 s_{k-1}}{s'_{k-1} A s_{k-1}} \quad (5)$$

respectively. In this case, α_k is the Rayleigh quotient of A at either the vector s_{k-1} or at the vector $\sqrt{A} s_{k-1}$. Since A is SPD,

$$0 < \lambda_{\min} \leq \alpha_k \leq \lambda_{\max} \quad \text{for all } k, \quad (6)$$

where λ_{\min} and λ_{\max} are respectively the smallest and largest eigenvalues of A . Hence, there is no danger of dividing by zero in (1).

In the rest of this paper, we will only consider the Barzilai and Borwein method with the choice of α_k defined by (4). The reason for this is that all results established for this method with the choice of α_k given by (4) also hold with the choice of α_k given by (5).

2. Convergence analysis

We will establish the convergence of the Barzilai and Borwein method applied to any quadratic function

$$f(x) = \frac{1}{2} x' A x - b' x + c$$

with an SPD Hessian matrix A .

Let x_* be the unique minimizer of f , $\{x_k\}$ the sequence generated by the Barzilai and Borwein method from a given vector x_0 , and $e_k = x_* - x_k$ for all k . Then, using (1) and the fact that $g_k = A x_k - b$, where $b = A x_*$, we have

$$A e_k = \alpha_k s_k \quad \text{for all } k. \quad (7)$$

Substituting $s_k = e_k - e_{k+1}$ in (7) we obtain for any k

$$e_{k+1} = \frac{1}{\alpha_k} (\alpha_k I - A) e_k. \quad (8)$$

Now for any initial error e_0 , there exist constants $d_1^0, d_2^0, \dots, d_n^0$ such that

$$e_0 = \sum_{i=1}^n d_i^0 v_i,$$

where $\{v_1, v_2, \dots, v_n\}$ are orthonormal eigenvectors of A associated with the eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$. Using (8) we obtain for any integer k ,

$$e_{k+1} = \sum_{i=1}^n d_i^{k+1} v_i, \quad (9)$$

where

$$d_i^{k+1} = \left(\frac{\alpha_k - \lambda_i}{\alpha_k} \right) d_i^k = \prod_{j=0}^k \left(\frac{\alpha_j - \lambda_i}{\alpha_j} \right) d_i^0. \quad (10)$$

We observe that the convergence properties of the sequence $\{e_k\}$ will depend on the behaviour of each one of the sequences $\{d_i^k\}$, $1 \leq i \leq n$. In general, these sequences will increase at some iterations. However, the following lemma shows that the sequence $\{d_1^k\}$ will decrease Q-linearly to zero.

LEMMA 1 The sequence $\{d_1^k\}$ converges to zero Q-linearly with convergence factor $\hat{c} = 1 - (\lambda_{\min}/\lambda_{\max})$.

Proof. For any positive integer k ,

$$d_1^{k+1} = \left(\frac{\alpha_k - \lambda_{\min}}{\alpha_k} \right) d_1^k.$$

Since α_k satisfies (6), we have

$$|d_1^{k+1}| = \left(1 - \frac{\lambda_{\min}}{\alpha_k} \right) |d_1^k| \leq \hat{c} |d_1^k|,$$

where

$$\hat{c} = 1 - \frac{\lambda_{\min}}{\lambda_{\max}} < 1. \quad \square$$

In the proof of our convergence theorem, we will use the following result.

LEMMA 2 If the sequences $\{d_1^k\}$, $\{d_2^k\}$, ..., $\{d_l^k\}$ all converge to zero for a fixed integer l , $1 \leq l < n$. Then,

$$\liminf_{k \rightarrow \infty} |d_{l+1}^k| = 0.$$

Proof. Suppose, by way of contradiction, that there exists a constant $\varepsilon > 0$ such that

$$(d_{l+1}^k)^2 \lambda_{l+1}^2 \geq \varepsilon \quad \text{for all } k. \quad (11)$$

By (4) and (7) it follows that the Rayleigh quotient α_{k+1} can be written as

$$\alpha_{k+1} = \frac{e_k' A^3 e_k}{e_k' A^2 e_k}.$$

Then, by (9) and the orthonormality of the eigenvectors $\{v_1, v_2, \dots, v_n\}$, we obtain

$$\alpha_{k+1} = \frac{\sum_{i=1}^n (d_i^k)^2 \lambda_i^3}{\sum_{i=1}^n (d_i^k)^2 \lambda_i^2}. \quad (12)$$

Since the sequences $\{d_1^k\}$, ..., $\{d_l^k\}$ all converge to zero, there exists \hat{k} sufficiently large such that

$$\sum_{i=1}^l (d_i^k)^2 \lambda_i^2 \leq \frac{1}{2} \varepsilon \quad \text{for all } k \geq \hat{k}. \quad (13)$$

By (12) and (13), we obtain, for any $k \geq \hat{k}$,

$$\frac{(\sum_{i=l+1}^n (d_i^k)^2 \lambda_i^2) \lambda_{l+1}}{\frac{1}{2} \varepsilon + (\sum_{i=l+1}^n (d_i^k)^2 \lambda_i^2)} \leq \alpha_{k+1} \leq \lambda_{\max}. \quad (14)$$

Since

$$\sum_{i=l+1}^n (d_i^k)^2 \lambda_i^2 \geq (d_{l+1}^k)^2 \lambda_{l+1}^2 \geq \varepsilon,$$

then it follows from (14) that

$$\frac{2}{3} \lambda_{l+1} \leq \alpha_{k+1} \leq \lambda_{\max} \quad \text{for all } k \geq \hat{k},$$

which implies the bound

$$\left| 1 - \frac{\lambda_{l+1}}{\alpha_k} \right| \leq \max \left(\frac{1}{2}, 1 - \frac{\lambda_{l+1}}{\lambda_{\max}} \right) \quad \text{for all } k \geq \hat{k} + 1. \quad (15)$$

Finally, using (15) and the first part of equation (10), we obtain, for all $k \geq \hat{k} + 1$,

$$|d_{l+1}^{k+1}| = \left| 1 - \frac{\lambda_{l+1}}{\alpha_k} \right| |d_{l+1}^k| \leq \hat{c} |d_{l+1}^k|,$$

where

$$\hat{c} = \max \left(\frac{1}{2}, 1 - \frac{\lambda_{\min}}{\lambda_{\max}} \right) < 1. \quad (16)$$

Because this conclusion contradicts the hypothesis (11), we find that the lemma is true. \square

Theorem 1 establishes the convergence of the Barzilai and Borwein method when applied to a quadratic function with an SPD Hessian.

THEOREM 1 Let $f(x)$ be a strictly convex quadratic function. Let $\{x_k\}$ be the sequence generated by the Barzilai and Borwein gradient method and x_\star the unique minimizer of f . Then, either $x_j = x_\star$ for some finite j , or the sequence $\{x_k\}$ converges to x_\star .

Proof. We need only consider the case in which there is no finite integer j such that $x_j = x_\star$. Hence, it suffices to prove that the sequence $\{e_k\}$ converges to zero. From (9) and the orthonormality of the eigenvectors we have

$$\|e_k\|_2^2 = \sum_{i=1}^n (d_i^k)^2.$$

Therefore, the sequence of errors $\{e_k\}$ converges to zero if and only if each one of the sequences $\{d_i^k\}$ for $i = 1, \dots, n$ converges to zero.

Lemma 1 shows that $\{d_1^k\}$ converges to zero. We prove that $\{d_p^k\}$ converges to zero for $2 \leq p \leq n$ by induction on p . Therefore we let p be any integer from this interval, and we assume that $\{d_1^k\}, \dots, \{d_{p-1}^k\}$ all tend to zero. Then for any given $\varepsilon > 0$ there exists \hat{k} sufficiently large such that

$$\sum_{i=1}^{p-1} (d_i^k)^2 \lambda_i^2 < \frac{1}{2} \varepsilon \quad \text{for all } k \geq \hat{k}. \quad (17)$$

From (12) and (17), we obtain

$$\frac{(\sum_{i=p}^n (d_i^k)^2 \lambda_i^2) \lambda_p}{\frac{1}{2} \varepsilon + (\sum_{i=p}^n (d_i^k)^2 \lambda_i^2)} \leq \alpha_{k+1} \leq \lambda_{\max} \quad (18)$$

for all integers $k \geq \hat{k}$. Moreover, by Lemma 2, there exists $k_p \geq \hat{k}$ such that

$$(d_p^{k_p})^2 \lambda_p^2 < \varepsilon.$$

Now, let us say that $k_0 > k_p$ is any integer for which $(d_p^{k_0-1})^2 \lambda_p^2 < \varepsilon$ and $(d_p^{k_0})^2 \lambda_p^2 \geq \varepsilon$. Clearly,

$$\sum_{i=p}^n (d_i^k)^2 \lambda_i^2 \geq (d_p^k)^2 \lambda_p^2 \geq \varepsilon \quad \text{for } k_0 \leq k \leq j-1, \quad (19)$$

where j is the first integer greater than k_0 for which $(d_p^j)^2 \lambda_p^2 < \varepsilon$. Then, by (18) and (19), we have

$$\frac{2}{3} \lambda_p \leq \alpha_{k+1} \leq \lambda_{\max} \quad \text{for } k_0 \leq k \leq j-1. \quad (20)$$

Thus, using (20) and the first part of equation (10), we obtain

$$|d_p^{k+2}| \leq \hat{c} |d_p^{k+1}| \quad \text{for } k_0 \leq k \leq j-1,$$

where \hat{c} is the constant (16), which satisfies $\hat{c} < 1$. Finally, using the bound

$$|d_p^{k_0+1}| \leq \left(\frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\min}} \right)^2 |d_p^{k_0-1}|,$$

which is implied by expression (6) and the first part of equation (10), we conclude that

$$(d_p^k)^2 \leq \left(\frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\min}} \right)^4 (d_p^{k_0-1})^2 \leq \left(\frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\min}} \right)^4 \frac{\varepsilon}{\lambda_p^2}$$

for all $k_0 + 1 \leq k \leq j + 1$. Further, equation (10) provides the inequality $(d_p^{k_0})^2 \leq [(\lambda_{\max} - \lambda_{\min})/\lambda_{\min}]^2 (d_p^{k_0-1})^2$. It follows from the conditions on k_0 and j that $(d_p^k)^2$ is bounded above by a constant multiple of ε for all $k \geq k_0 - 1$. Hence, since $\varepsilon > 0$ can be chosen arbitrarily small, we deduce $\lim_{k \rightarrow \infty} |d_p^k| = 0$ as required, which completes the proof. \square

Notice that with the choice of $\alpha_{k+1} = s_k' A^2 s_k / s_k' A s_k$, given by (5), equality (12) can be written as

$$\alpha_{k+1} = \frac{\sum_{i=1}^n (d_i^k)^2 \lambda_i^4}{\sum_{i=1}^n (d_i^k)^2 \lambda_i^3}.$$

Then, by a similar argument, the convergence result established in Theorem 1 with the choice of α_k given by (4) also holds with the choice of α_k given by (5).

Acknowledgements

During the course of this research the author was a graduate student in the Department of Mathematical Sciences at Rice University, Houston, Texas. He gratefully acknowledges Universidad Central de Venezuela for support during

this period of time. He also thanks R. Byrd, J. E. Dennis, R. A. Tapia and the referees for helpful comments and suggestions.

REFERENCES

- BARZILAI, J., & BORWEIN, J. M. 1988 Two point step size gradient methods. *IMA Journal of Numerical Analysis* **8**, 141–148.
- FLETCHER, R. 1990 Low storage methods for unconstrained optimization. *Lectures in Applied Mathematics (AMS)* **26**, 165–179.