

An Acceleration of Gradient Descent Algorithm with Backtracking for Unconstrained Optimization

Neculai Andrei

National Institute for Research and Development in Informatics,

Center for Advanced Modeling and Optimization,

8-10, Avereșcu Avenue, Bucharest 1,

ROMANIA

E-mail: nandrei@ici.ro

Abstract. In this paper we introduce an acceleration of gradient descent algorithm with backtracking. The idea is to modify the steplength t_k by means of a positive parameter θ_k , in a multiplicative manner, in such a way to improve the behaviour of the classical gradient algorithm. It is shown that the resulting algorithm remains linear convergent, but the reduction in function value is significantly improved.

Keywords: gradient descent methods, backtracking, acceleration methods

1. Introduction

One of the first and very well known method for unconstrained optimization

$$\min f(x) \tag{1}$$

where $f: R^n \rightarrow R$ is convex and continuously differentiable, is the gradient descent method, designed by Cauchy early in 1847 [3]. In this method the negative gradient direction is used to find the local minimizers. The algorithm starts with an initial point $x_0 \in \text{dom } f$ and generates a sequence of points according to the following iterative procedure:

$$x_{k+1} = x_k + t_k d_k, \tag{2}$$

where t_k is the stepsize and $d_k = -g_k = -\nabla f(x_k)$.

The method proved to be effective for functions very well conditioned, but for functions poorly conditioned it is excessively slow, thus being of no practical value. Even for quadratic functions the gradient descent method with exact line search behave increasingly badly when the condition number of the matrix deteriorates. Early attempts to increase the performance of the method have been considered by Humphrey [9], Forsythe and Motzkin [7] and Schinzinger [14].

As we know, at the current point x_k the direction of the negative gradient is the best direction of search for a minimizer of function f , and this is the direction of gradient descent method. However, as soon as we move in this direction, it ceases to be the best and continue to deteriorate until it becomes orthogonal to $-\nabla f(x_k)$. That is, the method begins to take small steps without making significant progress to minimum. This is the major drawback of the gradient descent method, the steps it takes are too long, i.e. there are some other points z_k on the line segment connecting x_k and x_{k+1} , where $-\nabla f(z_k)$ provides a better new search direction than $-\nabla f(x_{k+1})$.

The purpose of this paper is to present an acceleration of the gradient descent method. The idea is to modify the steplength t_k (computed by backtracking) by means of a positive parameter θ_k in a multiplicative manner in such a way as to improve the behaviour of the classical gradient algorithm. It is shown that the resulting algorithm is linear convergent, but the reducing in function value is significantly improved. The structure of the paper is as follows. In section two we present the line search with backtracking and prove that the Armijo rule in a backtracking scheme generates steplengths bounded away from zero. Section 3 presents the accelerated gradient descent algorithm. It is shown that the accelerated algorithm reduces the function values with an exponent that is smaller than the corresponding exponent of gradient descent algorithm. Some numerical examples and comparisons are given in section 4.

2. Line Search with Backtracking

For implementing the algorithm (2) one of the crucial element is the stepsize computation. Many procedures have been suggested. In the *exact line search* the step t_k is selected as:

$$t_k = \arg \min_{t > 0} f(x_k + td_k). \quad (3)$$

In some special cases (for example quadratic problems) it is possible to compute the step t_k analytically, but in most cases it is computed to approximately minimize f along the ray $\{x_k + td_k : t \geq 0\}$, or at least to reduce f sufficiently. In practice the most used are the *inexact procedures*. A lot of inexact line search methods have been proposed: Goldstein [8], Armijo [1], Wolfe [16], Powell [13], Dennis and Schnabel [4], Fletcher [6], Potra and Shi [12], Lemaréchal [10], Moré and Thuente [11], and many others. In particular, one of the very simple and efficient line search procedure is the backtracking line search. This procedure considers the following scalars $0 < \alpha < 0.5$, $0 < \beta < 1$ and $s_k = -g_k^T d_k / \|d_k\|^2$ and takes the following steps based on the Armijo's rule:

Backtracking procedure

Step 1. Consider the descent direction d_k for f at point x_k . Set $t = s_k$.

Step 2. While $f(x_k + td_k) > f(x_k) + \alpha t \nabla f(x_k)^T d_k$, set $t = t\beta$.

Step 3. Set $t_k = t$.

Typically, $\alpha = 0.0001$ and $\beta = 0.8$, meaning that we accept a small portion of the decrease predicted by linear approximation of f at the current point. Observe that, if $d_k = -g_k$, then $s_k = 1$.

Proposition 1. Suppose that d_k is a descent direction and $\nabla f(x)$ satisfies the Lipschitz condition $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ for all x, y in $\{x: f(x) \leq f(x_0)\}$, where L is a positive constant. If the line search satisfies the Armijo condition, then

$$t_k \geq \min \left\{ 1, \frac{\beta(1-\alpha) - g_k^T d_k}{L \|d_k\|^2} \right\}. \quad (4)$$

Proof. Set $K_1 = \{k: t_k = s_k\}$ and $K_2 = \{k: t_k < s_k\}$. Then for all $k \in K_1$ we have

$$f_k - f_{k+1} \geq -\alpha s_k g_k^T d_k$$

and for all $k \in K_2$:

$$f_k - f_{k+1} \geq -\alpha t_k g_k^T d_k.$$

By Armijo rule since $t_k / \beta \leq s_k$ for all $k \in K_2$, we have

$$f_k - f(x_k + \frac{t_k}{\beta} d_k) < -\alpha \frac{t_k}{\beta} g_k^T d_k,$$

for all $k \in K_2$. Now, using the mean value theorem on the left side of this inequality, it follows that there exists $\xi_k \in [0, 1]$, such that

$$g(x_k + \frac{t_k}{\beta} \xi_k d_k)^T d_k > \alpha g_k^T d_k,$$

for all $k \in K_2$.

Having in view the Lipschitz condition, from the above inequality, by Cauchy-Schwartz inequality, we have:

$$\begin{aligned} \frac{t_k}{\beta} L \|d_k\|^2 &\geq \left\| g(x_k + \frac{t_k}{\beta} \xi_k d_k) - g(x_k) \right\| \|d_k\| \geq \left(g(x_k + \frac{t_k}{\beta} \xi_k d_k) - g(x_k) \right)^T d_k \\ &\geq \alpha g_k^T d_k - g_k^T d_k = -(1-\alpha) g_k^T d_k, \end{aligned}$$

for all $k \in K_2$.

Rearranging this inequality and combining it with the inequality corresponding for $k \in K_1$ we get the bound (5). (See also [15].)

3. Accelerated Gradient Descent Algorithm

In this section we present an accelerated gradient descent algorithm for solving unconstrained optimization problem (1). Considering the initial point x_0 we can compute $f_0 = f(x_0)$, $g_0 = \nabla f(x_0)$ and by backtracking determine $t_0 = \underset{0 < t \leq 1}{\operatorname{argmin}} f(x_0 - t g_0)$. With these, the next

iteration is computed as $x_1 = x_0 - t_0 g_0$ where again f_1 and g_1 can be immediately computed. Now, at the iteration $k = 1, 2, \dots$ we know x_k , f_k and g_k . Using the backtracking procedure the stepsize t_k can be determined with which the following point $z = x_k - t_k g_k$ is computed. By backtracking procedure we get a $t_k \in (0, 1]$ such that:

$$f(z) = f(x_k - t_k g_k) \leq f(x_k) - \alpha t_k g_k^T g_k.$$

With these, let us introduce the accelerated gradient descent algorithm by means of the following iterative scheme:

$$x_{k+1} = x_k - \theta_k t_k g_k, \tag{5}$$

where $\theta_k > 0$ is a parameter which follows to be determined in such a manner to improve the behaviour of the gradient descent algorithm.

Now, we have:

$$f(x_k - t_k g_k) = f(x_k) - t_k g_k^T g_k + \frac{1}{2} t_k^2 g_k^T \nabla^2 f(x_k) g_k + o(\|t_k g_k\|^2).$$

On the other hand, for $\theta > 0$, we have:

$$f(x_k - \theta t_k g_k) = f(x_k) - \theta t_k g_k^T g_k + \frac{1}{2} \theta^2 t_k^2 g_k^T \nabla^2 f(x_k) g_k + o(\|\theta t_k g_k\|^2).$$

Rejecting the high order terms we can write:

$$f(x_k - \theta t_k g_k) = f(x_k - t_k g_k) + \Psi_k(\theta), \tag{6}$$

where

$$\Psi_k(\theta) = (1-\theta) t_k g_k^T g_k - \frac{1}{2} (1-\theta^2) t_k^2 g_k^T \nabla^2 f(x_k) g_k. \tag{7}$$

Let us denote:

$$a_k = t_k g_k^T g_k \geq 0,$$

$$b_k = t_k^2 g_k^T \nabla^2 f(x_k) g_k.$$

Observe that for convex functions $b_k \geq 0$. Therefore

$$\Psi_k(\theta) = (1-\theta) a_k - \frac{1}{2} (1-\theta^2) b_k. \tag{8}$$

Now, we see that $\Psi_k'(\theta) = b_k \theta - a_k$ and $\Psi_k'(\theta_m) = 0$, where $\theta_m = a_k / b_k$. More than this, $\Psi_k(1) = 0$, $\Psi_k(2\theta_m - 1) = 0$ and $\Psi_k'(0) = -a_k < 0$. Therefore, $\Psi_k(\theta)$ is a convex quadratic

function with minimum value in point θ_m and

$$\Psi_k(\theta_m) = -\frac{(a_k - b_k)^2}{2b_k} \leq 0.$$

Considering $\theta = \theta_m$ in (6), we see that for every k

$$f(x_k - \theta_m t_k g_k) = f(x_k - t_k g_k) - \frac{(a_k - b_k)^2}{2b_k} \leq f(x_k - t_k g_k),$$

which is a possible improvement on the values of function f , (when $a_k \neq b_k$).

Therefore, using this simple multiplicative modification of the stepsize t_k as $\theta_k t_k$, where $\theta_k = \theta_m = a_k / b_k$, we get:

$$\begin{aligned} f(x_{k+1}) &= f(x_k - \theta_k t_k g_k) \leq f(x_k) - \left[\alpha t_k g_k^T g_k + \frac{(a_k - b_k)^2}{2b_k} \right] \\ &= f(x_k) - \left[\alpha a_k + \frac{(a_k - b_k)^2}{2b_k} \right] \leq f(x_k), \end{aligned} \quad (9)$$

since

$$\alpha a_k + \frac{(a_k - b_k)^2}{2b_k} \geq 0.$$

Now, in order to establish the algorithm we must determine a way for b_k computation. For this, at point $z = x_k - t_k g_k$ we have:

$$f(z) = f(x_k - t_k g_k) = f(x_k) - t_k g_k^T g_k + \frac{1}{2} t_k^2 g_k^T \nabla^2 f(\tilde{x}_k) g_k,$$

where \tilde{x}_k is a point on the line segment connecting x_k and z . On the other hand, at point $x_k = z + t_k g_k$ we have:

$$f(x_k) = f(z + t_k g_k) = f(z) + t_k g_z^T g_k + \frac{1}{2} t_k^2 g_k^T \nabla^2 f(\bar{x}_k) g_k,$$

where $g_z = \nabla f(z)$ and \bar{x}_k is a point on the line segment connecting x_k and z .

Having in view the local character of searching and that the distance between x_k and z is enough small, we can consider $\tilde{x}_k = \bar{x}_k = x_k$. With these, adding the above equalities we get:

$$b_k = t_k^2 g_k^T \nabla^2 f(x_k) g_k = -t_k y_k^T g_k, \quad (10)$$

where $y_k = g_z - g_k$. Now, we have all the elements to present our algorithm.

Accelerated Gradient Descent Algorithm (AGD)

Step 1. Consider a starting point $x_0 \in \text{dom } f$ and compute: $f_0 = f(x_0)$ and $g_0 = \nabla f(x_0)$. Set $k = 0$.

Step 2. Using backtracking determine: $t_k = \underset{0 < t \leq 1}{\text{argmin}} f(x_k - t g_k)$.

Step 3. Compute: $z = x_k - t_k g_k$, $g_z = \nabla f(z)$ and $y_k = g_z - g_k$.

Step 4. Compute: $a_k = t_k g_k^T g_k$, $b_k = -t_k y_k^T g_k$ and $\theta_k = a_k / b_k$.

Step 5. Update the variables: $x_{k+1} = x_k - \theta_k t_k g_k$ and compute f_{k+1} and g_{k+1} .

Step 6. Test a criterion for stopping the iterations. If the test is satisfied, then stop;

otherwise consider $k = k + 1$ and go to step 2

The gradient descent (GD) algorithm can be immediately particularized from AGD by skipping steps 3 and 4 and considering $\theta_k = 1$ in step 5 where variables are updated.

Observe that if $a_k > b_k$, then $\theta_k > 1$. In this case $\theta_k t_k > t_k$ and it is possible that $\theta_k t_k \leq 1$ or $\theta_k t_k > 1$, i.e. it is possible that the steplength $\theta_k t_k$ to be greater than 1. On the other hand, if $a_k \leq b_k$, then $\theta_k \leq 1$. In this case $\theta_k t_k \leq t_k \leq 1$, i.e. the steplength $\theta_k t_k$ is reduced. Therefore, if $a_k \neq b_k$, then $\theta_k \neq 1$ and the steplength t_k computed by backtracking will be modified, by its increasing or reducing through factor θ_k , thus avoiding the algorithm to take orthogonal steps along the iterations.

It is worth saying that if $a_k \leq b_k / 2$, then $\Psi_k(0) = a_k - b_k / 2 \leq 0$ and $\theta_k < 1$. For any $\theta \in [0,1]$, $\Psi_k(\theta) \leq 0$. As a consequence from (6) we see that for any $\theta \in (0,1)$, $f(x_k - \theta t_k g_k) < f(x_k)$. In this case, for any $\theta \in [0,1]$, $\theta_k t_k \leq t_k$. However, in our algorithm we selected $\theta_k = \theta_m$ as the point achieving the minimum value of $\Psi_k(\theta)$.

Proposition 2. Suppose that f is a strongly convex function on the level set $S = \{x: f(x) \leq f(x_0)\}$. Then the sequence x_k generated by AGD converges linearly to x^* .

Proof. From (9) we have that $f(x_{k+1}) \leq f(x_k)$, for all k . Since f is bounded below, it follows that

$$\lim_{k \rightarrow \infty} (f(x_k) - f(x_{k+1})) = 0.$$

Since f is strongly convex there are positive constants m and M such that $mI \leq \nabla^2 f(x) \leq MI$ on the level set S . Suppose that $x_k - t g_k \in S$ and $x_k - \theta_m t g_k \in S$, for $0 < t \leq 1$.

We have

$$f(x_k - \theta_m t g_k) = f(x_k - t g_k) - \frac{(a_k - b_k)^2}{2b_k}.$$

But, from strong convexity we have the following quadratic upper bound on $f(x_k - t g_k)$:

$$f(x_k - t g_k) \leq f(x_k) - t \|g_k\|_2^2 + \frac{Mt^2}{2} \|g_k\|_2^2.$$

Observe that for $0 \leq t \leq 1/M$, $-t + Mt^2/2 \leq -t/2$ which follows from the convexity of $-t + Mt^2/2$. Using this result we get:

$$\begin{aligned} f(x_k - t g_k) &\leq f(x_k) - t \|g_k\|_2^2 + \frac{Mt^2}{2} \|g_k\|_2^2 \\ &\leq f(x_k) - \frac{t}{2} \|g_k\|_2^2 \leq f(x_k) - \alpha t \|g_k\|_2^2, \end{aligned}$$

since $\alpha \leq 1/2$.

The backtracking terminates either with $t = 1$ or with a value $t \geq \beta / M$. This provides a lower bound on the decrease in the function f . For $t = 1$ we have:

$$f(x_k - t g_k) \leq f(x_k) - \alpha \|g_k\|_2^2$$

and for $t \geq \beta / M$

$$f(x_k - t g_k) \leq f(x_k) - \frac{\alpha \beta}{M} \|g_k\|_2^2.$$

Therefore, for $0 \leq t \leq 1/M$, we always have

$$f(x_k - tg_k) \leq f(x_k) - \min\left\{\alpha, \frac{\alpha\beta}{M}\right\} \|g_k\|_2^2. \quad (11)$$

On the other hand

$$\frac{(a_k - b_k)^2}{2b_k} \geq \frac{(t\|g_k\|_2^2 - t^2 m \|g_k\|_2^2)^2}{2t^2 M \|g_k\|_2^2} = \frac{(1 - tm)^2}{2M} \|g_k\|_2^2.$$

Now, as above, for $t = 1$

$$\frac{(a_k - b_k)^2}{2b_k} \geq \frac{(1 - m)^2}{2M} \|g_k\|_2^2.$$

For $t \geq \beta / M$

$$\frac{(a_k - b_k)^2}{2b_k} \geq \frac{(1 - \beta m / M)^2}{2M} \|g_k\|_2^2 \geq \frac{(1 - \beta)^2}{2M} \|g_k\|_2^2,$$

since $0 < \beta < 1$ and $m / M \leq 1$.

Therefore,

$$\frac{(a_k - b_k)^2}{2b_k} \geq \min\left\{\frac{(1 - m)^2}{2M}, \frac{(1 - \beta)^2}{2M}\right\} \|g_k\|_2^2. \quad (12)$$

Considering (11) together with (12) we get:

$$f(x_k - \theta_m tg_k) \leq f(x_k) - \min\left\{\alpha, \frac{\alpha\beta}{M}\right\} \|g_k\|_2^2 - \min\left\{\frac{(1 - m)^2}{2M}, \frac{(1 - \beta)^2}{2M}\right\} \|g_k\|_2^2. \quad (13)$$

Therefore

$$f(x_k) - f(x_{k+1}) \geq \left[\min\left\{\alpha, \frac{\alpha\beta}{M}\right\} + \min\left\{\frac{(1 - m)^2}{2M}, \frac{(1 - \beta)^2}{2M}\right\} \right] \|g_k\|_2^2.$$

But, $f(x_k) - f(x_{k+1}) \rightarrow 0$ and as a consequence g_k goes to zero, i.e. x_k converges to x^* . Having in view that $f(x_k)$ is a nonincreasing sequence, it follows that $f(x_k)$ converges to $f(x^*)$.

From (13) we see that

$$f(x_{k+1}) \leq f(x_k) - \left[\min\left\{\alpha, \frac{\alpha\beta}{M}\right\} + \min\left\{\frac{(1 - m)^2}{2M}, \frac{(1 - \beta)^2}{2M}\right\} \right] \|g_k\|_2^2.$$

Combining this with

$$\|g_k\|_2^2 \geq 2m(f(x_k) - f^*)$$

and subtract f^* from both sides of the above inequality we conclude:

$$f(x_{k+1}) - f^* \leq c(f(x_k) - f^*), \quad (14)$$

where

$$c = 1 - \min\left\{2m\alpha, \frac{2m\alpha\beta}{M}\right\} - \min\left\{\frac{(1 - m)^2 m}{M}, \frac{(1 - \beta)^2 m}{M}\right\} < 1. \quad (15)$$

Therefore, $f(x_k)$ converges to f^* at least as fast as a geometric series with an exponent that depends on the backtracking parameters and the condition number bound M / m . Therefore, the convergence is at least linear.

At every iteration k , selecting $\theta_k = \theta_m$ in (5), the AGD algorithm reduces the function values according to (14) where c is given by (15). Since GD algorithm achieves (14) with

$$c = 1 - \min \left\{ 2m\alpha, \frac{2m\alpha\beta}{M} \right\} < 1, \quad (16)$$

it follows that, if $\theta_m \neq 1$, the AGD algorithm is an improvement, i.e. an acceleration, of GD.

4. Numerical Results and Comparisons

In this section we report some numerical results obtained with a Fortran implementation of the above gradient descent algorithms. The full description of these experiments is documented at the web page: <http://www.ici.ro/camo/neculai/ansoft.htm>.

All codes are written in Fortran and compiled with f77 (default compiler settings) on a Workstation 1.8 Ghz.

For the very beginning let us present the following example illustrating the behaviour of these algorithms. In the next series of experiments we present the numerical results obtained with a number of 340 unconstrained optimization test functions.

An illustrative Example. Let us illustrate the behaviour of the accelerated gradient descent algorithm on the following function (trigonometric function):

$$f(x) = \sum_{i=1}^n \left(\left(n - \sum_{j=1}^n \cos x_j \right) + i(1 - \cos x_i) - \sin x_i \right)^2.$$

Considering $x_0 = [0.2, 0.2, \dots, 0.2]$, $\alpha = 0.0001$ and $\beta = 0.8$ in the backtracking procedure, as well as the the following criteria for stopping the iterations

$$\|\nabla f(x_k)\|_{\infty} \leq \varepsilon_g \quad \text{or} \quad t_k |g_k^T d_k| \leq \varepsilon_f |f(x_{k+1})|,$$

with $\varepsilon_g = 10^{-6}$ and $\varepsilon_f = 10^{-20}$, for $n = 100$, the evolution of $|f(x_k) - f^*|$, αa_k and $(a_k - b_k)^2 / (2b_k)$, for different values of backtracking parameter α is illustrated in Figures 1, 2, 3, 4 and 5.

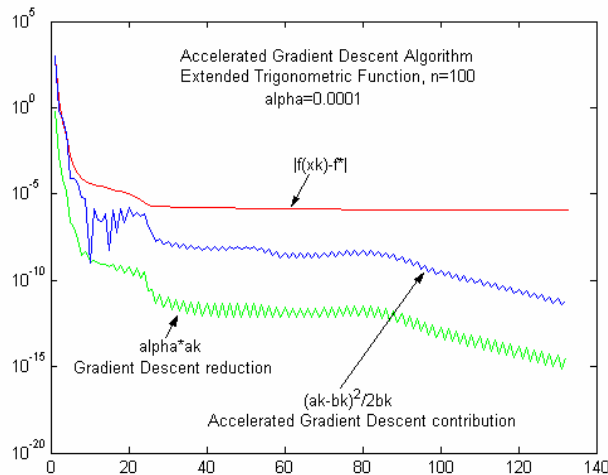


Figure 1. Accelerated Gradient Descent characteristics. $\alpha = 0.0001$.

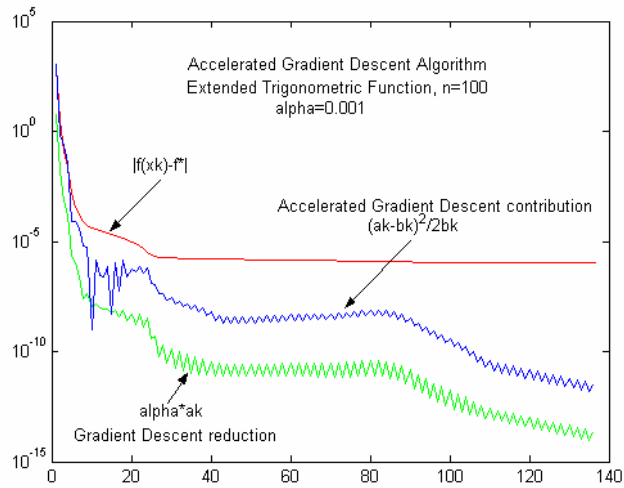


Figure 2. Accelerated Gradient Descent characteristics. $\alpha = 0.001$.

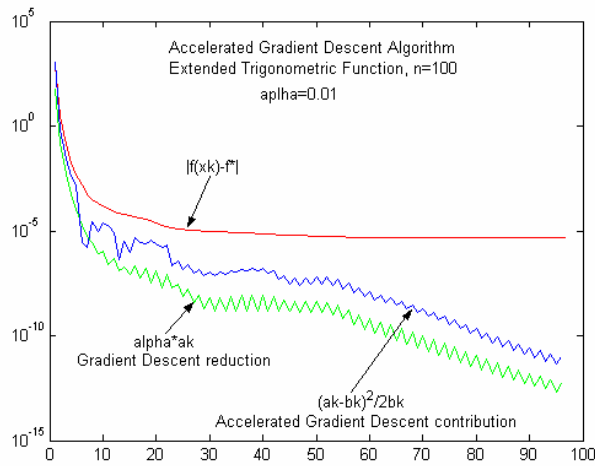


Figure 3. Accelerated Gradient Descent Characteristics. $\alpha = 0.01$.

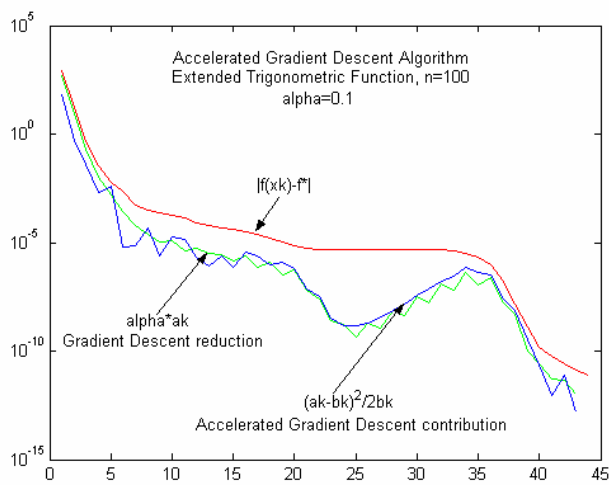


Figure 4. Accelerated Gradient Descent characteristics. $\alpha = 0.1$.

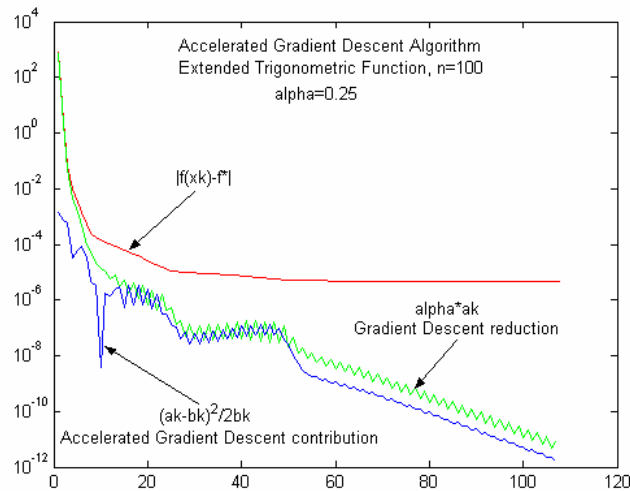


Figure 5. Accelerated Gradient Descent characteristics. $\alpha = 0.25$.

From this example we see that the contribution of the acceleration scheme, given by

$$\frac{(a_k - b_k)^2}{2b_k},$$

is substantial. In comparison with the reduction operated by the gradient descent algorithm, given by αa_k , observe that the contribution of the accelerated scheme is dependent on the value of the backtracking parameter α . For α enough small, for example $\alpha = 0.0001$ as usually is considered in the backtracking procedure, we see that for almost all iterations:

$$\frac{(a_k - b_k)^2}{2b_k} \geq \alpha a_k.$$

Table 1 shows the number of iterations (#iter), number of function evaluations (#fg), cpu time (seconds) (cpu(s)), as well as the average steplength corresponding to GD and AGD algorithms for $n = 100, 200, \dots, 1000$.

Table 1. Algorithm Characteristics. $\alpha = 0.0001$.

	Gradient Descent				Accelerated Gradient Descent			
	#iter	#fg	cpu (s)	\bar{t}_k	#iter	#fg	cpu (s)	$\bar{\theta}_k \bar{t}_k$
100	45	448	0.17	0.48490	131	516	0.22	1.38764
200	69	598	0.49	0.53880	140	543	0.50	0.91710
300	304	1608	2.03	0.69477	137	595	0.82	1.37628
400	64	877	1.49	0.43523	114	536	0.94	1.19269
500	327	1786	3.68	0.68818	148	681	1.42	1.14217
600	73	1056	2.64	0.42819	99	502	1.27	1.34367
700	70	1202	3.51	0.35077	38	292	0.88	0.91164
800	80	1337	4.45	0.37733	25	260	0.82	0.59176
900	95	1505	5.66	0.40425	108	570	2.20	1.15465
1000	62	1052	4.39	0.35699	63	417	1.87	1.13343
TOTAL	1189	11469	28.51		1003	4912	10.94	

Table 2 shows the number of iterations (#iter), number of function evaluations (#fg), cpu time (seconds) (cpu(s)), as well as the average steplength corresponding to GD and AGD algorithms for $n = 1000, 2000, \dots, 10000$.

Table 2. Algorithm Characteristics. $\alpha = 0.0001$.

	Gradient Descent				Accelerated Gradient Descent			
	#iter	#fg	cpu (s)	\bar{t}_k	#iter	#fg	cpu (s)	$\bar{\theta}_k t_k$
1000	62	1052	3.96	0.35047	63	417	1.60	1.11510
2000	96	2033	15.43	0.28254	165	769	5.87	1.16464
3000	119	2411	27.47	0.34527	24	416	4.67	0.42287
4000	69	1924	29.22	0.16058	74	641	9.78	1.10106
5000	125	2770	56.90	0.28841	90	719	13.68	1.31227
6000	91	2846	65.03	0.12231	60	645	14.77	0.94762
7000	156	3537	94.36	0.31129	162	1076	28.89	1.17126
8000	155	4003	122.21	0.26681	67	645	19.72	1.13035
9000	100	3048	104.91	0.16543	55	682	23.40	1.00044
10000	168	4083	156.32	0.27942	142	1051	40.32	1.16849
TOTAL	1141	27707	675.81		902	7061	162.70	

For the next series of experiments, we selected 34 large-scale unconstrained optimization test problems (5 from CUTE library [2]) in extended or generalized form. For each test function we have considered 10 numerical experiments with number of variables $n = 100, 200, \dots, 1000$.

In the following we present the numerical performance of AGD and GD codes, in which we stopped the iterations as soon as one of the above criteria is satisfied.

Table 3 presents the global characteristics, corresponding to these 340 test problems, referring to the total number of iterations, total number of function evaluations and total cpu time for these algorithms.

Table 3. Global characteristics of AGD and GD. 340 problems

Global characteristics	AGD	GD
Total number of iterations	657915	1598990
Total number of function evaluations	17432902	40979772
Total cpu time (seconds)	4547.80	10043.26

Table 4 shows the number of problems, out of 340, for which AGD and GD achieved the minimum number of iterations, minimum number of function evaluations and the minimum cpu time, respectively.

Table 4. Performance of AGD and GD Algorithms. 340 Problems

Performance criterion	# of problems
AGD achieved minimum # of iterations in	320
GD achieved minimum # of iterations in	36
AGD and GD achieved the same # of iterations in	16
AGD achieved minimum # of function evaluations in	322
GD achieved minimum # of function evaluations in	18
AGD achieved minimum cpu time in	325
GD achieved minimum cpu time in	28
AGD and GD achieved the same cpu time in	13

The performance of these algorithms have been evaluated using the profiles of Dolan and Moré [5]. That is, for each algorithm, we plot the fraction of problems for which the algorithm is within a factor of the best number of iterations, the best number of function evaluations, or the best cpu time, respectively. The left side of these Figures gives the percentage of the test problems, out of 340, for which an algorithm performs better; the right side gives the percentage of the test problems that were successfully solved by each of the algorithms. Mainly, the right side represents a measure of an algorithm's robustness.

In Figure 6, 7 and 8 we display the performance profiles of Dolan and Moré, corresponding to AGD and GD, referring to the number of iterations, number of functions evaluations and cpu time, respectively.

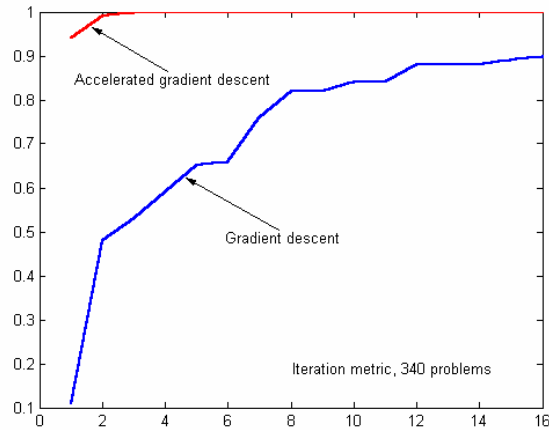


Figure 6. Performance Based on Number of Iterations, 340 Problems.

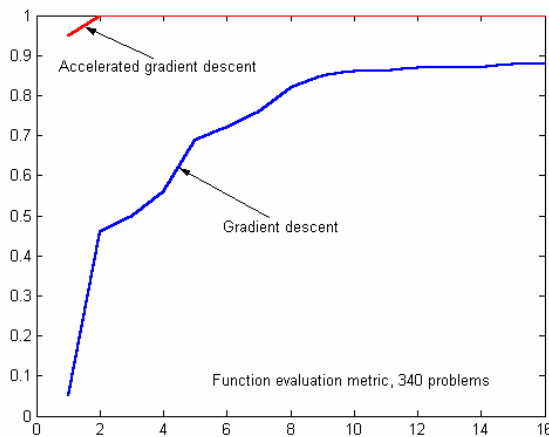


Figure 7. Performance Based on Number of Function Evaluations, 340 Problems

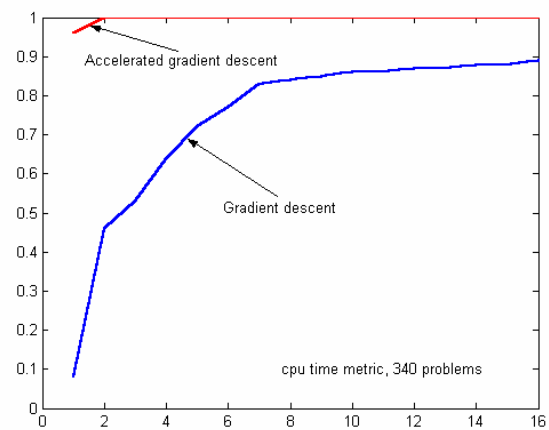


Figure 8. Performance based on cpu time metric, 340 problems

Observe that AGD outperforms GD in the vast majority of problems, and the differences are substantial. From table 3 we see that referring to cpu time AGD is about twice fastest than GD. We explain this difference in behaviour of these algorithms by recalling that as the stationary point is approached, GD method takes small, nearly orthogonal steps. This poor convergence of the GD algorithm at the later iterations can be explained by considering the following expression of function f :

$$f(x_k - t_k g_k) = f(x_k) - t_k \|g_k\|^2 + \frac{1}{2} t_k^2 \gamma_k \|g_k\|^2, \quad (17)$$

where $\gamma_k I \cong \nabla^2 f(z)$ is a scalar approximation of the Hessian at the point z which belongs to the line segment connecting x_k and x_{k+1} . Observe that if x_k is close to a stationary point with zero gradient, and f is continuously differentiable, then $\|g_k\|^2$ will be small. Therefore, $t_k \|g_k\|^2$ in (17) is of a small order of magnitude, its contribution to reduce the function values being almost insignificant. Since the gradient descent method uses only the linear approximation of f to find the search direction, ignoring completely the second order term $(t_k^2 / 2)\gamma_k \|g_k\|^2$, we expect that the direction generated will not be very effective, if the ignored term contributes significantly to the description of f , even for relatively small values of t_k . In AGD this is compensated by modifying the steplength in order to destroy the orthogonality of the successive search directions giving thus the possibility for a substantial progress towards minimum.

5. Conclusions

In this paper we have introduced an acceleration of the gradient descent algorithm by means of a simple multiplicative modification of the steplength given by a backtracking procedure in the classical gradient descent algorithm. The accelerated algorithm is linear convergent with an exponent which is smaller than the exponent corresponding to the gradient descent algorithm.

REFERENCES

1. L. ARMIJO, **Minimization of functions having Lipschitz continuous first partial derivatives**, Pacific Journal of Mathematics, vol. 6, pp. 1-3, 1966.
2. I. BONGARTZ, A.R. CONN, N.I.M. GOULD and P.L. TOINT, **CUTE: constrained and unconstrained testing environments**, ACM Trans. Math. Software, vol. 21, pp. 123-160, 1995.
3. A. CAUCHY, **Méthodes générales pour la résolution des systèmes d'équations simultanées**, C.R. Acad. Sci. Par., vol. 25, pp. 536-538, 1848.
4. J.E. DENNIS and R.B. SCHNABEL, **Numerical Methods for Unconstrained Optimization and Nonlinear Equations**, Englewoods Cliffs, N.J., Prentice-Hall, 1983.
5. E.D. DOLAN, J.J. MORÉ, **Benchmarking optimization software with performance profiles**, Math. Programming, vol. 91, pp. 201-213, 2002.
6. R. FLETCHER, **Practical Methods of Optimization**, John Wiley and Sons, New York, 1987.
7. G.E. FORSYTHE and T.S. MOTZKIN, **Asymptotic properties of the optimum gradient method**, Bull. American Society, vol. 57, pp. 183, 1951.
8. A.A. GOLDSTEIN, **On steepest descent**, SIAM Journal on Control, vol. 3, pp. 147-151, 1965.
9. W.E. HUMPHREY, **A general minimising routine - minfun**, in Recent Advances in Optimisation Techniques, A. Lavi, and T.P. Vogl, (Eds.), John Wiley, 1966.
10. C. LEMARÉCHAL, **A view of line search**, in Optimization and Optimal Control, A. Auslander, W. Oettli and J. Stoer (Eds.) Springer Verlag, pp. 59-78, 1981.
11. J.J. MORÉ and D.J. THUENTE, **On line search algorithms with guaranteed sufficient decrease**, Mathematics and Computer Science Division Preprint MCS-P153-0590, Argonne National Laboratory, Argonne, 1990.
12. F.A. POTRA and Y. SHI, **Efficient line search algorithm for unconstrained optimization**, Journal of Optimization Theory and Applications, vol. 85, pp. 677-704, 1995.
13. M.J.D. POWELL, **Some global convergence properties of a variable-metric algorithm for minimization without exact line searches**, SIAM-AMS Proceedings, Philadelphia, vol. 9, pp. 53-72, 1976.
14. R. SCHINZINGER, **Optimization in electromagnetic system design**, in Recent Advances in Optimisation Techniques, A. Lavi, and T.P. Vogl (Eds.), John Wiley, 1966.
15. SHI ZHEN-JUN, **Convergence of line search methods for unconstrained optimization**, Applied Mathematics and Computation, vol. 157, pp. 393-405, 2004.
16. P. WOLFE, **Convergence conditions for ascent methods**, SIAM Review, vol. 11, pp. 226-235, 1968.