New Accelerated Conjugate Gradient Algorithms as modification of Dai-Yuan's computational scheme for Unconstrained Optimization

Neculai Andrei

Research Institute for Informatics, Center for Advanced Modeling and Optimization, 8-10, Averescu Avenue, Bucharest 1, Romania E-mail: nandrei@ici.ro

Abstract. New accelerated nonlinear conjugate gradient algorithms which are mainly modifications of the Dai and Yuan's for unconstrained optimization are proposed. Using the exact line search, the algorithm reduces to the Dai and Yuan conjugate gradient computational scheme. For inexact line search the algorithm satisfies the sufficient descent condition. Since the step lengths in conjugate gradient algorithms may differ from 1 by two order of magnitude and tend to vary in a very unpredictable manner, the algorithms are equipped with an acceleration scheme able to improve the efficiency of the algorithms. Computational results for a set consisting of 750 unconstrained optimization test problems show that these new conjugate gradient algorithms substantially outperform the Dai-Yuan conjugate gradient algorithm and its hybrid variants, Hestenes-Stiefel, Polak-Ribière-Polyak, CONMIN conjugate gradient algorithms, limited quasi-Newton algorithm LBFGS and compare favourable with CG_DESCENT.

Keywords: Unconstrained optimization; conjugate gradient method; sufficient descent condition; conjugacy condition; Newton direction; numerical comparisons.

AMS 2000 Mathematics Subject Classification: 49M07, 49M10, 90C06, 65K05

1. Introduction

Conjugate gradient methods represent an important class of unconstrained optimization algorithms with strong local and global convergence properties and modest memory requirements. A survey on their definition including 40 conjugate gradient algorithms for unconstrained optimization is given by Andrei [6]. A discussion of development of different versions of nonlinear conjugate gradient methods, with special attention to global convergence properties is presented by Hager and Zhang [21].

In this paper we suggest new nonlinear conjugate gradient algorithms which are mainly modifications of the Dai and Yuan [16] conjugate gradient computational scheme. In these algorithms the direction d_{k+1} is computed as a linear combination between $-g_{k+1}$ and s_k , i.e. $d_{k+1} = -\theta_{k+1}g_{k+1} + \beta_k^N s_k$, where $g_k = \nabla f(x_k)$ and $s_k = x_{k+1} - x_k$. The parameter θ_k is computed in such a way that the direction d_{k+1} is the Newton direction or it satisfies the conjugacy condition. On the other hand, β_k^N is a proper modification of the Dai and Yuan's computational scheme in such a way that the direction d_{k+1} at every iteration satisfies the sufficient descent condition. For the exact line search the proposed algorithms reduce to the Dai and Yuan conjugate gradient computational scheme.

The paper has the following structure. In Section 2 we present the development of the conjugate gradient algorithms with sufficient descent condition as modifications of the Dai-Yuan computational scheme, while in section 3 we prove the global convergence of these algorithms under strong Wolfe line search conditions. In Section 4 we present the accelerated

algorithms, showing their global convergence and in Section 5 we compare the computational performance of the new conjugate gradient schemes against the Dai and Yuan method and its hybrid variants [17], Hestenes and Stiefel [22], Polak-Ribière [27] and Polyak [28], CG_DESCENT by Hager and Zhang [20], CONMIN by Shanno and Phua [30], as well as LBFGS by Liu and Nocedal [23], using 750 unconstrained optimization test problems from the CUTE [12] library along with some other large-scale unconstrained optimization problems presented in [8]. Using the Dolan and Moré performance profiles [19] we prove these new accelerated conjugate gradient algorithms outperform the Dai-Yuan algorithm as well as its hybrid variants, Hestenes-Stiefel, Polak-Ribière-Polyak, CONMIN, LBFGS and compare favourable with CG_DESCENT by Hager and Zhang.

2. Modifications of the Dai-Yuan conjugate gradient algorithm

For solving the unconstrained optimization problem

$$\min\left\{f(x):x\in \mathbb{R}^n\right\},\tag{2.1}$$

where $f: \mathbb{R}^n \to \mathbb{R}$ is continuously differentiable and bounded below we consider a nonlinear conjugate gradient algorithm:

$$x_{k+1} = x_k + \alpha_k d_k, \qquad (2.2)$$

where the stepsize α_k is positive and the directions d_k are computed by the rule:

$$d_{k+1} = -\theta_{k+1}g_{k+1} + \beta_k^N s_k, \quad d_0 = -g_0,$$
(2.3)

where

$$\beta_{k}^{N} = \frac{\left\|g_{k+1}\right\|^{2}}{y_{k}^{T}s_{k}} - \frac{\left\|g_{k+1}\right\|^{2}(s_{k}^{T}g_{k+1})}{(y_{k}^{T}s_{k})^{2}},$$
(2.4)

and θ_{k+1} is a parameter which follows to be determined. Here $y_k = g_{k+1} - g_k$ and $s_k = x_{k+1} - x_k$.

The line search in the conjugate gradient algorithms for α_k computation is often based on the standard Wolfe conditions [31, 32]:

$$f(x_k + \alpha_k d_k) - f(x_k) \le \rho \alpha_k g_k^T d_k, \qquad (2.5)$$

$$g_{k+l}^T d_k \ge \sigma g_k^T d_k \,, \tag{2.6}$$

where d_k is a descent direction and $0 < \rho \le \sigma < 1$.

Observe that if f is a quadratic function and α_k is selected to achieve the exact minimum of f in the direction d_k , then $s_k^T g_{k+1} = 0$ and the formula (2.4) for β_k^N reduces to the Dai and Yuan computational scheme [16]. However, in this paper we refer to general nonlinear functions and inexact line search.

We were led to this computational scheme by modifying the Dai and Yuan algorithm

$$\beta_{k}^{DY} = \frac{g_{k+1}g_{k+1}}{y_{k}^{T}s_{k}},$$

in order to have the sufficient descent condition, as well as some other properties for an efficient conjugate gradient algorithm. Using a standard Wolfe line search, the Dai and Yuan method always generates descent directions and under Lipschitz assumption it is globally convergent. In [13] Dai established a remarkable property relating the descent directions to the sufficient descent condition, showing that if there exist constants γ_1 and γ_2 such that $\gamma_1 \leq ||g_k|| \leq \gamma_2$ for all k, then for any $p \in (0,1)$, there exists a constant c > 0 such that the sufficient descent condition $g_i^T d_i \leq -c||g_i||^2$ holds for at least |pk| indices $i \in [0, k]$, where |j| denotes the largest integer $\leq j$. In our algorithm the parameter β_k is selected in such a manner that the sufficient descent condition is satisfied at every iteration. As we know, despite the strong convergence theory that has been developed for the Dai and

Yuan method, it is susceptible to jamming, that is it begins to take small steps without making significant progress to the minimum. When iterates jam, y_k becomes tiny while $||g_k||$ is bounded away from zero. Therefore, β_k^N is a proper modification of the β_k^{DY} .

Theorem 2.1. If $\theta_{k+1} \ge 1/4$, then the direction $d_{k+1} = -\theta_{k+1}g_{k+1} + \beta_k^N s_k$, $(d_0 = -g_0)$, where β_k^N is given by (2.4) satisfies the sufficient descent condition

$$g_{k+1}^{T}d_{k+1} \leq -\left(\theta_{k+1} - \frac{1}{4}\right) \left\|g_{k+1}\right\|^{2}.$$
(2.7)

Proof. Since $d_0 = -g_0$, we have $g_0^T d_0 = -||g_0||^2$, which satisfy (2.7). Multiplying (2.3) by g_{k+1}^T , we have

$$g_{k+1}^{T}d_{k+1} = -\theta_{k+1} \left\| g_{k+1} \right\|^{2} + \frac{(g_{k+1}^{T}g_{k+1})(g_{k+1}^{T}s_{k})}{y_{k}^{T}s_{k}} - \frac{\left\| g_{k+1} \right\|^{2}(s_{k}^{T}g_{k+1})^{2}}{(y_{k}^{T}s_{k})^{2}}.$$
 (2.8)

Now, using the inequality $u^T v \leq \frac{1}{2} (||u||^2 + ||v||^2)$, where $u, v \in \mathbb{R}^n$, we have:

$$\frac{(g_{k+1}^{T}g_{k+1})(g_{k+1}^{T}s_{k})}{y_{k}^{T}s_{k}} = \frac{\left[(y_{k}^{T}s_{k})g_{k+1}/\sqrt{2}\right]^{T}\left[\sqrt{2}(g_{k+1}^{T}s_{k})g_{k+1}\right]}{(y_{k}^{T}s_{k})^{2}}$$

$$\leq \frac{\frac{1}{2}\left[\frac{1}{2}(y_{k}^{T}s_{k})^{2} \|g_{k+1}\|^{2} + 2(g_{k+1}^{T}s_{k})^{2} \|g_{k+1}\|^{2}\right]}{(y_{k}^{T}s_{k})^{2}}$$

$$= \frac{1}{4}\|g_{k+1}\|^{2} + \frac{(g_{k+1}^{T}s_{k})^{2} \|g_{k+1}\|^{2}}{(y_{k}^{T}s_{k})^{2}}.$$
(2.9)

Using (2.9) in (2.8) we get (2.7). ■

To conclude, the sufficient descent condition from (2.7), the quantity $\theta_{k+1} - 1/4$ is required to be nonnegative. Supposing that $\theta_{k+1} - 1/4 > 0$, then the direction given by (2.3) and (2.4) is a descent direction. Dai and Yuan [16, 17] present conjugate gradient schemes with the property that $g_k^T d_k < 0$ when $y_k^T s_k > 0$. If f is strongly convex or the line search satisfies the Wolfe conditions, then $y_k^T s_k > 0$ and the Dai and Yuan scheme yield descent. In our algorithm observe that, if for all k, $\theta_{k+1} > 1/4$, and the line search satisfies the Wolfe conditions (2.5) and (2.6), then for all k the search direction (2.3) and (2.4) satisfy the sufficient descent condition. It is well known that if the Wolfe line search conditions are satisfied, then $y_k^T s_k > 0$ and the steplength α_k is bounded away from zero [20]. Observe that $y_k^T s_k > 0$ is crucial in (2.4) for β_k^N computation. Note that in (2.7) we bound $g_{k+1}^T d_{k+1}$ by $-(\theta_{k+1} - 1/4) ||g_{k+1}||^2$, while for the computational scheme of Dai and Yuan only the nonnegativity of $g_{k+1}^T d_{k+1}$ is established.

To determine the parameter θ_{k+1} in (2.3) we suggest the following two procedures.

A) When the initial point x_0 is near the solution of (2.1) and the Hessian of function f is a nonsingular matrix we know that the Newton direction is the best line search direction. Therefore, to get a good algorithm for solving (2.1) this is a very good motivation to choose

the parameter θ_k in such a way that for every $k \ge 1$ the direction d_{k+1} given by (2.3) be the Newton direction. Therefore, from the equation

$$-\nabla^2 f(x_{k+1})^{-1} g_{k+1} = -\theta_{k+1} g_{k+1} + \beta_k^N s_k$$
(2.10)

after some algebra we get

$$\theta_{k+1} = \frac{1}{s_k^T \nabla^2 f(x_{k+1}) g_{k+1}} \left[\frac{\|g_{k+1}\|^2}{y_k^T s_k} \left(1 - \frac{s_k^T g_{k+1}}{y_k^T s_k} \right) s_k^T \nabla^2 f(x_{k+1}) s_k + s_k^T g_{k+1} \right].$$
(2.11)

Observe that the choice (2.11) does not imply that d_{k+1} given by (2.3) is the Newton direction. This is only a technical operation to get θ_{k+1} as in (2.11). The salient point in this formula for θ_{k+1} is the presence of the Hessian. For large-scale problems, choices for the update parameter that do not require the evaluation of the Hessian matrix are often preferred in practice to the methods that require the Hessian in each iteration. Therefore, in order to have an algorithm for solving large-scale problems we assume that in (2.10) we use an approximation B_{k+1} of the true Hessian $\nabla^2 f(x_{k+1})$ and let B_{k+1} satisfy the quasi-Newton equation $B_{k+1}s_k = y_k$. This leads us to:

$$\theta_{k+1} = \frac{1}{y_k^T g_{k+1}} \left[\left\| g_{k+1} \right\|^2 - \frac{\left\| g_{k+1} \right\|^2 (s_k^T g_{k+1})}{y_k^T s_k} + s_k^T g_{k+1} \right].$$
(2.12)

Observe that if θ_{k+1} given by (2.12) is greater than or equal to 1/4, then according to Theorem 2.1 the direction (2.3) satisfies the sufficient descent condition (2.7). On the other hand, if in (2.12) $\theta_{k+1} < 1/4$, then we take *ex abrupto* $\theta_{k+1} = 1$ in (2.3).

B) The second procedure is based on the conjugacy condition. Dai and Liao [14] introduced the conjugacy condition $y_k^T d_{k+1} = -ts_k^T g_{k+1}$, where $t \ge 0$ is a scalar. This is indeed very reasonable since in real computation the inexact line search is generally used. However, this condition is very dependent on the nonnegative parameter t, for which we do not know any formula to choose in an optimal manner. Therefore, even if in our developments we use the inexact line search we adopt here a more conservative approach and consider the conjugacy condition $y_k^T d_{k+1} = 0$. This leads us to:

$$\theta_{k+1} = \frac{1}{y_k^T g_{k+1}} \left[\left\| g_{k+1} \right\|^2 - \frac{\left\| g_{k+1} \right\|^2 (s_k^T g_{k+1})}{y_k^T s_k} \right].$$
 (2.13)

As above, if θ_{k+1} given by (2.13) is greater than or equal to 1/4, then according to Theorem 2.1 the direction (2.3) satisfies the sufficient descent condition (2.7). On the other hand, if in (2.13) $\theta_{k+1} < 1/4$, then we take $\theta_{k+1} = 1$ in (2.3).

Observe that since $s_k^T g_{k+1} \rightarrow 0$ along the iterations, θ_k given by (2.12) obtained from the Newton direction paradigm is very similar to (2.13) based on the conjugacy condition. Besides, θ_{k+1} from (2.13) can be written as

$$\theta_{k+1} = \frac{\|g_{k+1}\|^2}{\|g_{k+1}\|^2 - g_k^T g_{k+1}} \left[1 - \frac{(s_k^T g_{k+1})}{y_k^T s_k}\right].$$

Since at every iteration d_k is a descent direction and α_k is computed by the Wolfe line search (2.5) and (2.6), it follows that $g_k^T g_{k+1} \rightarrow 0$. (This is reminiscence from the steepest descent method.) Therefore, along the iterations, $\theta_k \rightarrow 1$.

In [17] Dai and Yuan proved the global convergence of a conjugate gradient algorithm for which $\beta_k = \beta_k^{DY} t_k$, where $t_k \in [-c, 1]$ with $c = (1 - \sigma)/(1 + \sigma)$. Our algorithm is a proper modification of the Dai and Yuan's with the following property. Observe that

$$\beta_{k}^{N} = \frac{\left\| g_{k+1} \right\|^{2}}{y_{k}^{T} s_{k}} \left[1 - \frac{s_{k}^{T} g_{k+1}}{y_{k}^{T} s_{k}} \right] = \beta_{k}^{DY} r_{k}, \qquad (2.14)$$

where

$$r_k = 1 - \frac{s_k^T g_{k+1}}{y_k^T s_k}.$$
 (2.15)

From the second Wolfe condition it follows that $s_k^T g_{k+1} \ge \sigma s_k^T g_k = -\sigma y_k^T s_k + \sigma s_k^T g_{k+1}$, i.e.

$$s_k^T g_{k+1} \ge \frac{-\sigma}{1-\sigma} y_k^T s_k$$

Since by the Wolfe condition $y_k^T s_k > 0$, it follows that $\frac{s_k^T g_{k+1}}{y_k^T s_k} \ge \frac{-\sigma}{1-\sigma}$. Hence $r_k \le \frac{1}{1-\sigma}$. Therefore

Therefore,

$$\beta_k^N \le \beta_k^{DY} \frac{1}{1 - \sigma}.$$
(2.16)

3. Convergence analysis

In this section we analyze the convergence of the algorithm (2.2), (2.3), (2.4) and (2.12) or (2.13) where $d_0 = -g_0$. In the following we consider that $g_k \neq 0$ for all $k \ge 1$, otherwise a stationary point is obtained. Assume that:

- (i) The level set $S = \{x \in \mathbb{R}^n : f(x) \le f(x_0)\}$ is bounded.
- (ii) In a neighborhood N of S, the function f is continuously differentiable and its gradient is Lipschitz continuous, i.e. there exists a constant L > 0 such that $\|\nabla f(x) \nabla f(y)\| \le L \|x y\|$, for all $x, y \in N$.

Under these assumptions on f there exists a constant $\Gamma \ge 0$ such that $\|\nabla f(x)\| \le \Gamma$ for all $x \in S$. In order to prove the global convergence, we assume that the step size α_k in (2.2) is obtained by the strong Wolfe line search, that is,

$$f(x_k + \alpha_k d_k) - f(x_k) \le \rho \alpha_k g_k^T d_k, \qquad (3.1)$$

$$\left|g(x_{k}+\alpha_{k}d_{k})^{T}d_{k}\right| \leq \sigma g_{k}^{T}d_{k}.$$
(3.2)

where ρ and σ are positive constants such that $0 < \rho \le \sigma < 1$.

For any conjugate gradient algorithm with strong Wolfe line search, we have the following results given by lemma 3.1 and lemma 3.2, which were first proved by Zoutendijk [33] and Wolfe [31, 32]. For completeness, we present them here without proofs.

Lemma 3.1. Let α_k be obtained by the strong Wolfe line search (3.1) and (3.2). Suppose that the assumptions (i) and (ii) and d_k is a descent direction. Then

$$\sum_{k=0}^{\infty} -\alpha_k g_k^T d_k < \infty.$$
(3.3)

Lemma 3.2. Let α_k be obtained by the strong Wolfe line search (3.1) and (3.2). Suppose that the assumptions (i) and (ii) and d_k is a descent direction. Then the so-called Zoutendijk condition holds

$$\sum_{k=0}^{\infty} \frac{(g_k^T d_k)^2}{\|d_k\|^2} < \infty.$$
 (3.4)

Based on these results, for conjugate gradient method (2.2) where

$$d_k = -\theta_k g_k + \alpha_{k-1} \beta_{k-1}^N d_{k-1}$$

$$(3.5)$$

and $\theta_k > 1/4$, with strong Wolfe line search, we can prove the following lemma and its corollary which are essential for the convergence of our algorithms. Lemma 3.3 is a variant of the Theorem 2.3 of Dai *et al.* [18].

Lemma 3.3. Suppose that the assumptions (i) and (ii) hold. Consider the conjugate gradient method (2.2) and (3.5) where $\theta_k > 1/4$, with strong Wolfe line search (3.1) and (3.2). Then either

$$\liminf_{k \to \infty} \|g_k\| = 0, \tag{3.6}$$

or

$$\sum_{k=0}^{\infty} \frac{\left\|\boldsymbol{g}_{k}\right\|^{4}}{\left\|\boldsymbol{d}_{k}\right\|^{2}} < \infty.$$
(3.7)

Proof. Since for any $k \ge 0$, $\theta_k > 1/4$, it follows that d_k is a descent direction. From (3.5) since for all $k \ge 0$, $g_k^T d_k < 0$ we have

$$\|d_{k}\|^{2} \ge (\alpha_{k-1}\beta_{k-1}^{N})^{2} \|d_{k-1}\|^{2} - \theta_{k}^{2} \|g_{k}\|^{2}.$$
(3.8)

On the other hand, from (3.5) we get

$$g_{k}^{T}d_{k} - \alpha_{k-1}\beta_{k-1}^{N}g_{k}^{T}d_{k-1} = -\theta_{k}\left\|g_{k}\right\|^{2}$$

Since d_k is a descent direction, it follows that

$$\alpha_{k-1}\beta_{k-1}^{N}g_{k}^{T}d_{k-1} + |g_{k}^{T}d_{k}| = \theta_{k}||g_{k}||^{2}$$

Therefore,

$$\alpha_{k-1} \left| \beta_{k-1}^{N} \right| \left| g_{k}^{T} d_{k-1} \right| + \left| g_{k}^{T} d_{k} \right| \geq \theta_{k} \left\| g_{k} \right\|^{2}.$$

From the strong Wolfe condition we have that

$$\sigma \alpha_{k-1} \left| \beta_{k-1}^{N} \right\| g_{k-1}^{T} d_{k-1} \right| + \left| g_{k}^{T} d_{k} \right| \ge \theta_{k} \left\| g_{k} \right\|^{2}.$$
(3.9)

But for any $a, b, \sigma \ge 0$ the following inequality $(a + \sigma b)^2 \le (1 + \sigma^2)(a^2 + b^2)$ holds. Considering $a = \left|g_k^T d_k\right|$ and $b = \alpha_{k-1} \left|\beta_{k-1}^N \right| g_{k-1}^T d_{k-1} \right|$, then (3.9) yields to

$$(g_k^T d_k)^2 + (\alpha_{k-1} \beta_{k-1}^N)^2 (g_{k-1}^T d_{k-1})^2 \ge c \|g_k\|^4,$$
(3.10)

where $c = \theta_k^2 / (1 + \sigma^2)$ is a positive constant. Therefore, from (3.10) we get

$$\frac{(g_{k}^{T}d_{k})^{2}}{\|d_{k}\|^{2}} + \frac{(g_{k-1}^{T}d_{k-1})^{2}}{\|d_{k-1}\|^{2}} = \frac{1}{\|d_{k}\|^{2}} \left[(g_{k}^{T}d_{k})^{2} + \frac{\|d_{k}\|^{2}}{\|d_{k-1}\|^{2}} (g_{k-1}^{T}d_{k-1})^{2} \right]$$
$$\geq \frac{1}{\|d_{k}\|^{2}} \left[c \|g_{k}\|^{4} + (g_{k-1}^{T}d_{k-1})^{2} \left(\frac{\|d_{k}\|^{2}}{\|d_{k-1}\|^{2}} - (\alpha_{k-1}\beta_{k-1}^{N})^{2} \right) \right].$$

From (3.8) observe that

Therefore,

$$\frac{\left\|d_{k}\right\|^{2}}{\left\|d_{k-1}\right\|^{2}} \ge (\alpha_{k-1}\beta_{k-1}^{N})^{2} - \theta_{k}^{2} \frac{\left\|g_{k}\right\|^{2}}{\left\|d_{k-1}\right\|^{2}}.$$

$$\frac{(g_{k}^{T}d_{k})^{2}}{\left\|d_{k}\right\|^{2}} + \frac{(g_{k-1}^{T}d_{k-1})^{2}}{\left\|d_{k-1}\right\|^{2}} \ge \frac{\left\|g_{k}\right\|^{4}}{\left\|d_{k}\right\|^{2}} \left[c - \theta_{k}^{2} \frac{(g_{k-1}^{T}d_{k-1})^{2}}{\left\|d_{k-1}\right\|^{2}} \frac{1}{\left\|g_{k}\right\|^{2}}\right], \quad (3.11)$$
rom lemma 3.2 we know that

where $\theta_k > 1/4$. From lemma 3.2 we know that

$$\lim_{k \to \infty} \frac{(g_{k-1}^T d_{k-1})^2}{\|d_{k-1}\|^2} = 0$$

Therefore, if (3.6) is not true, then

$$\lim_{k\to\infty} \frac{(g_{k-1}^T d_{k-1})^2}{\|d_{k-1}\|^2} \frac{1}{\|g_k\|^2} = 0.$$

Therefore, from (3.11) we get that

$$\frac{(g_{k}^{T}d_{k})^{2}}{\left\|d_{k}\right\|^{2}} + \frac{(g_{k-1}^{T}d_{k-1})^{2}}{\left\|d_{k-1}\right\|^{2}} \ge c \frac{\left\|g_{k}\right\|^{4}}{\left\|d_{k}\right\|^{2}}$$

holds for all sufficiently large k. Hence, the inequality (3.7) follows from Zoutendijk condition (3.4) in lemma 3.2.

Corollary 3.1. Suppose that the assumptions (i) and (ii) hold and consider any conjugate gradient method (2.2) and (3.5), where d_k is a descent direction, i.e. $\theta_k > 1/4$, and α_k is obtained by the strong Wolfe line search (3.1) and (3.2). If

$$\sum_{k\ge 1} \frac{1}{\|d_k\|^2} = \infty, \qquad (3.12)$$

then

$$\liminf_{k \to \infty} \|g_k\| = 0. \tag{3.13}$$

Proof. Suppose that there is a positive constant γ such that $||g_k|| \ge \gamma$ for all $k \ge 0$. Then, from lemma 3.3 we have

$$\sum_{k\geq 0} \frac{1}{\|d_k\|^2} \leq \frac{1}{\gamma^4} \sum_{k\geq 0} \frac{\|g_k\|^4}{\|d_k\|^2} < \infty.$$

However, this contradicts (3.12) from the corollary 3.1, i.e. the corollary 3.1 is true.

Theorem 3.1. Suppose that the assumptions (i) and (ii) hold and consider the algorithm (2.2), (2.3), (2.4) and (2.12) or (2.13), where d_{k+1} is a descent direction and α_k is obtained by the strong Wolfe line search (3.1) and (3.2). If there exists a constant $\gamma \ge 0$ such $\gamma \le ||\nabla f(x)||$, $1/4 \le \theta_k \le \tau$, where τ is a positive constant and the angle φ_k between g_k and d_k is bounded, i.e. $\cos \varphi_k \le \xi \le 0$ for all k = 0, 1, ..., then the algorithm satisfies $\liminf_{k \to \infty} g_k = 0$.

Proof. Observe that $y_k^T s_k = g_{k+1}^T s_k - g_k^T s_k \ge (\sigma - 1) g_k^T s_k$. But $g_k^T s_k = \|g_k\| \|s_k\| \cos \varphi_k$. Since d_k is a descent direction it follows that $g_k^T s_k \le \|g_k\| \|s_k\| \xi \le 0$ for all k = 0, 1, ..., i.e. $y_k^T s_k \ge -(1 - \sigma) \|g_k\| \|s_k\| \xi$. With these, from (2.16) we have

$$\beta_{k}^{N} \leq \frac{\left\|g_{k+1}\right\|^{2}}{y_{k}^{T} s_{k}} \frac{1}{1-\sigma} \leq \frac{\left\|g_{k+1}\right\|^{2}}{-(1-\sigma)^{2} \xi \left\|g_{k}\right\| \left\|s_{k}\right\|} \leq \frac{\Gamma^{2}}{-(1-\sigma)^{2} \xi \gamma \left\|s_{k}\right\|} = \frac{\eta}{\left\|s_{k}\right\|},$$

where

$$\eta = \frac{\Gamma^2}{-(1-\sigma)^2 \xi \gamma}.$$

Therefore

$$\|d_{k+1}\| \le |\theta_{k+1}| \|g_{k+1}\| + |\beta_k^N| \|s_k\| \le \tau \Gamma + \frac{\eta}{\|s_k\|} \|s_k\| = \tau \Gamma + \eta$$

This relation shows that

$$\sum_{k\geq 1} \frac{1}{\|d_k\|^2} \ge \frac{1}{(\tau\Gamma + \eta)^2} \sum_{k\geq 1} 1 = \infty.$$
(3.14)

Hence, from corollary 3.1 it follows that $\liminf_{k \to \infty} ||g_k|| = 0$.

4. AMDYN and AMDYC Algorithms

Nocedal [25] pointed out that in conjugate gradient methods the step lengths may differ from 1 in a very unpredictable manner. They can be larger or smaller than 1 depending on how the problem is scaled. This is in very sharp contrast to the Newton and quasi-Newton methods, including the limited memory quasi-Newton methods, which accept the unit steplength most of the time along the iterations, and therefore usually they require only few function evaluations per search direction. Numerical comparisons between conjugate gradient methods and the limited memory quasi Newton method by Liu and Nocedal [23] show that the latter is more successful [7]. One explanation of efficiency of this limited memory quasi-Newton method is given by its ability to accept unity step lengths along the iterations. In this section we take advantage of this behavior of conjugate gradient algorithms and consider an acceleration scheme of the above conjugate gradient algorithms. Basically the acceleration of the function values along the iterations (see [5, 9] and [10]). In accelerated algorithm instead of (2.2) the new estimation of the minimum point is computed as

$$x_{k+1} = x_k + \gamma_k \alpha_k d_k,$$

where

$$\gamma_m = -\frac{a_k}{b_k},$$

 $a_k = \alpha_k g_k^T d_k$, $b_k = -\alpha_k (g_k - g_z)^T d_k$, $z = x_k + \alpha_k d_k$ and $g_z = \nabla f(z)$. Hence, if $b_k \neq 0$, then $x_{k+1} = x_k + \gamma_k \alpha_k d_k$, otherwise $x_{k+1} = x_k + \alpha_k d_k$. Therefore, using the definitions of g_k , s_k , y_k and the above acceleration scheme we present the following conjugate gradient algorithms which are accelerated, modified versions of the Dai and Yuan algorithm with Newton direction (AMDYN) or with conjugacy condition (AMDYC).

AMDYN and AMDYC Algorithms

Step 1. Initialization. Select $x_0 \in \mathbb{R}^n$ and the parameters $0 < \rho < \sigma < 1$. Compute $f(x_0)$ and g_0 . Consider $d_0 = -g_0$ and $\alpha_0 = 1/||g_0||$. Set k = 0. Step 2. Test for continuation of iterations. If $||g_k||_{\infty} \le 10^{-6}$, then stop, else set k = k + 1. Step 3. Line search. Compute α_k satisfying the Wolfe line search conditions (2.5) and (2.6). Step 4. Compute: $z = x_k + \alpha_k d_k$, $g_z = \nabla f(z)$ and $y_k = g_k - g_z$. Step 5. Compute: $a_k = \alpha_k g_k^T d_k$, and $b_k = -\alpha_k y_k^T d_k$.

Step 6. Acceleration. If $b_k \neq 0$, then compute $\gamma_k = -a_k / b_k$ and update the variables as $x_{k+1} = x_k + \gamma_k \alpha_k d_k$, otherwise update the variables as $x_{k+1} = x_k + \alpha_k d_k$. Compute f_{k+1} and g_{k+1} . Compute $y_k = g_{k+1} - g_k$ and $s_k = x_{k+1} - x_k$.

Step 7. θ_{k+1} computation. For the algorithm AMDYN, θ_{k+1} is computed as in (2.12). For the algorithm AMDYC, θ_{k+1} is computed as in (2.13). If $\theta_{k+1} < 1/4$, then we set $\theta_{k+1} = 1$.

Step 8. Direction computation. Compute $d = -\theta_{k+1}g_{k+1} + \beta_k^N s_k$, where β_k^N is computed as in (2.4). If

$$g_{k+1}^{T} d \le -10^{-3} \|d\|_{2} \|g_{k+1}\|_{2}, \qquad (4.1)$$

then define $d_{k+1} = d$, otherwise set $d_{k+1} = -g_{k+1}$. Compute the initial guess $\alpha_k = \alpha_{k-1} \|d_{k-1}\| / \|d_k\|$, set k = k+1 and continue with step 2.

It is well known that if f is bounded along the direction d_k then there exists a stepsize α_k satisfying the Wolfe line search conditions (2.5) and (2.6). In our algorithm when the angle between d and $-g_{k+1}$ is not acute enough, then we restart the algorithm with the negative gradient $-g_{k+1}$ [11]. More sophisticated reasons for restarting the algorithms have been proposed in the literature [29], but we are interested in the performance of a conjugate gradient algorithm that uses this restart criterion, associated to a direction satisfying the sufficient descent condition. Under reasonable assumptions, conditions (2.5), (2.6) and (4.1) are sufficient to prove the global convergence of the algorithm.

The initial selection of the step length crucially affects the practical behaviour of the algorithm. At every iteration $k \ge 1$ the starting guess for the step α_k in the line search is computed as $\alpha_{k-1} \|d_{k-1}\|_2 / \|d_k\|_2$. This selection, was considered for the first time by Shanno and Phua in CONMIN [30]. It is also considered in the packages: SCG by Birgin and Martínez [11] and in SCALCG by Andrei [1-4, 7].

For uniformly convex functions, like in [10], we can prove that the sequence generated by AMDYN or AMDYC converges linearly to the solution of the problem (2.1).

Proposition 4.1. Suppose that f is a uniformly convex function on the level set $S = \{x : f(x) \le f(x_0)\}$, and d_k satisfies the sufficient descent condition $g_k^T d_k < -c_1 \|g_k\|^2$, where $c_1 > 0$, and $\|d_k\|^2 \le c_2 \|g_k\|^2$, where $c_2 > 0$. Then the sequence generated by AMDYN or AMDYC converges linearly to x^* , solution to the problem (2.1).

5. Numerical results and comparisons

In this section we present the computational performance of a Fortran implementation of the AMDYN and AMDYC algorithms on a set of 750 unconstrained optimization test problems. We selected 75 large-scale unconstrained optimization problems in extended or generalized form [8]. For each function we have considered ten numerical experiments with the increasing number of variables n = 1000,2000,...,10000. All algorithms implement the Wolfe line search conditions with $\rho = 0.0001$ and $\sigma = 0.9$, and the same stopping criterion $\|g_k\|_{\infty} \leq 10^{-6}$, where $\|\cdot\|_{\infty}$ is the maximum absolute component of a vector. The comparisons of algorithms are given in the following context. Let f_i^{ALG1} and f_i^{ALG2} be the optimal value found by ALG1 and ALG2, for problem i = 1,...,750, respectively. We say that, in the particular problem *i*, the performance of ALG1 was better than the performance of ALG2 if

$$\left| f_i^{ALG1} - f_i^{ALG2} \right| < 10^{-3} \tag{5.1}$$

and the number of iterations, or the number of function-gradient evaluations, or the CPU time of ALG1 was less than the number of iterations, or the number of function-gradient evaluations, or the CPU time corresponding to ALG2, respectively.

All codes are written in double precision Fortran and compiled with f77 (default compiler settings) on an Intel Pentium 4, 1.8GHz workstation. All these codes are authored by Andrei.

In the first set of numerical experiments we compare AMDYN versus AMDYC. In Table 1 we present the number of problems solved by these two algorithms with a minimum number of iterations (#iter), a minimum number of function and its gradient evaluations (#fg) and the minimum cpu time.

	AMDYN	AMDYC	=
# iter	83	105	562
# fg	152	147	451
CPU	143	119	488

Table 1. Performance of AMDYN versus AMDYC. 750 problems.

Both algorithms have similar performances. However, subject to cpu time metric, AMDYN proves to be slightly better. In the following we shall compare AMDYN versus some known conjugate gradient algorithms.

In the second set of numerical experiments we compare AMDYN algorithm with the Dai and Yuan (DY) algorithm. Figure 1 presents the Dolan-Moré performance profile for these algorithms subject to the cpu time metric. We see that AMDYN is top performer, being more successful and more robust than the Dai and Yuan algorithm. When comparing AMDYN with the Dai and Yuan algorithm (Figure 1), subject to the number of iterations, we see that AMDYN was better in 619 problems (i.e. it achieved the minimum number of iterations in 619 problems). DY was better in 27 problems and they achieved the same number of iterations in 60 problems, etc. Out of 750 problems, only for 706 of them does the criterion (5.1) hold.

Dai and Yuan [17] studied the hybrid conjugate gradient algorithms and proposed the following two hybrid methods:

$$\beta_{k}^{hDY} = \max\left\{-\frac{1-\sigma}{1+\sigma}\beta_{k}^{DY}, \min\left\{\beta_{k}^{HS}, \beta_{k}^{DY}\right\}\right\},$$
(5.2)

$$\boldsymbol{\beta}_{k}^{hDY_{z}} = \max\left\{0, \min\left\{\boldsymbol{\beta}_{k}^{HS}, \boldsymbol{\beta}_{k}^{DY}\right\}\right\},\tag{5.3}$$

where $\beta_k^{HS} = y_k^T g_{k+1} / y_k^T s_k$, showing their global convergence when the Lipschitz assumption holds and the standard Wolfe line search is used. The numerical experiments of Dai and Ni [15] proved that the second hybrid method (hDYz) is the better, outperforming the Polak-Ribière [27] and Polyak [28] method. In the third set of numerical experiments we compare the Dolan-Moré performance profile of AMDYN versus Dai-Yuan hybrid conjugate gradient β_k^{hDY} subject to the cpu time metric, as in Figure 2. Observe that the differences are substantial. Again AMDYN is top performer.



Fig. 1. Performance profile of AMDYN versus DY.



Fig. 2. Performance profile of AMDYN versus hDY.

In the fourth set of numerical experiments, in Figure 3, we compare the Dolan-Moré performance profile of AMDYN versus Dai-Yuan hybrid conjugate gradient β_k^{hDYz} subject to the cpu time metric. Again observe that AMDYN is top performer.

In the fifth set of numerical experiments we compare AMDYN versus Hestenes-Stiefel conjugate gradient algorithm $(\beta_k^{HS} = y_k^T g_{k+1} / y_k^T s_k)$. Figure 4 presents the performance profiles of these algorithms. The HS method has the property that the conjugacy condition $y_k^T d_{k+1} = 0$ always holds, independent of the line search. On the other hand, AMDYN algorithm satisfies the Dai-Liao conjugacy condition $y_k^T d_{k+1} = -s_k^T g_{k+1}$ which is a little more relaxed than the pure conjugate condition $y_k^T d_{k+1} = 0$.



Fig. 3. Performance profile of AMDYN versus hDYz.



Fig. 4. Performance profile of AMDYN versus HS.

In the sixth set of numerical comparisons we consider AMDYN versus the Polak-Ribière-Polyak conjugate gradient algorithm ($\beta_k^{PRP} = y_k^T g_{k+1} / g_k^T g_k$). Figure 5 presents the performance profiles of these algorithms subject to cpu time metric. The PRP method, like HS, posses a very important built-in restart feature that addresses directly to jamming. The idea is that PRP (and HS) method automatically adjust the value of the parameter β_k^{PRP} to avoid jamming. In general, the performance of these methods (PRP and HS) is better than the performance of some other conjugate gradient methods (for example DY) [21]. However, from Figure 5 observe that AMDYN is top performer again among these algorithms. AMDYN inherits some convergence properties from the Newton method (see (2.10)).



Fig. 5. Performance profile of AMDYN versus PRP.

In the next set of numerical experiments we compare AMDYN versus CG_DESCENT by Hager and Zhang [20]. Figure 6 presents the Dolan and Moré cpu time performance profile of AMDYN versus CG_DESCENT with Wolfe line search. Presently CG_DESCENT is the practical conjugate gradient algorithm with more reputation. CG_DESCENT is a modification of HS and was devised in order to ensure sufficient descent, independent of the accuracy of the line search. Hager and Zhang [20] proved that the direction d_k in their algorithm satisfies the sufficient descent condition $g_k^T d_k \leq -(7/8) \|g_k\|^2$.



Fig. 6. Performance profile of AMDYN versus CG_DESCENT.

At every iteration, the AMDYN algorithm satisfies the sufficient descent condition (2.7), where $\theta_k \rightarrow 1$. Therefore, at least in the last part of the iterations AMDYN satisfies the sufficient descent condition $g_k^T d_k \leq -(3/4) \|g_k\|^2$. CG_DESCENT has a very advanced line search procedure that utilizes the "approximate Wolfe conditions" which provides a more accurate way to check the usual Wolfe conditions when the iterates are near a local minimum of the function f. On the other hand, AMDYN uses an acceleration scheme which modify the step length given by the classical Wolfe condition (2.5) and (2.6) in order to improve the reduction of the function values along the iterations.

In the following, we compare AMDYN versus COMNIN by Shanno and Phua [30]. Figure 7 presents the performance profiles of these algorithms.



Fig. 7. Performance profile of AMDYN versus CONMIN.

COMNIN by Shanno and Phua [30] is a conjugate gradient algorithm which may be interpreted as a memoryless BFGS quasi-Newton algorithm optimally scaled in the sense of Oren and Spedicato [26]. In CONMIN the scaling is combined with the Powell's restart criterion. The direction d_{k+1} in CONMIN is computed as

$$d_{k+1} = -H_{k+1}g_{k+1} + A_k y_k - B_k s_k, (5.4)$$

where H_{k+1} is the BFGS approximation of the inverse Hessian which at every iteration is initialized with identity matrix and A_k and B_k are specific matrices. The main drawback of this method is that if H_{k+1} contains useful information about the Hessian of the function f, then we are better off using the search direction $d_{k+1} = -H_{k+1}g_{k+1}$ since the addition of the last terms in (5.4) may prevent d_{k+1} from being a descent direction unless the line search is sufficiently accurate. The same is the case for the AMDYN algorithm. The parameter θ_{k+1} in (2.3) given by (2.12) is computed to get as much as possible information from the inverse Hessian by the secant condition. However, the approximation of the inverse Hessian used in AMDYN is scantier that that used in CONMIN. In Figure 7 we have the computational evidence that subject to the cpu time metric, AMDYN is top performer and outperforms COMNIN. Finally we compare AMDYN versus LBFGS (m=3) by Liu and Nocedal [23] as in Figure 8, where m is the number of pairs (s_k, y_k) used. Observe that AMDYN is top performer again.



Fig. 8. Performance profile of AMDYN versus LBFGS (m=3).

One explanation is that the linear algebra in the LBFGS code to update the search direction is more time consuming than the linear algebra in AMDYN. On the other hand the steplength in LBFGS is determined at each iteration by means of the line search routine MCVSRCH, which is a slight modification of the routine CSRCH written by Moré and Thuente [24].

6. Conclusion

We have presented a new conjugate gradient algorithm for solving large-scale unconstrained optimization problems. The parameter β_k is a modification of the Dai and Yuan computational scheme in such a manner that the direction d_k generated by the algorithm satisfies the sufficient descent condition, independent of the line search. Under strong Wolfe line search conditions we proved the global convergence of the algorithm. We present computational evidence that the performance of our algorithms AMDYN and AMDYC was higher than that of the Dai and Yuan conjugate gradient algorithm and its hybrid variants, Hestenes-Stiefel, Polak-Ribière-Polyak, CONMIN, LBFGS (m=3) and compare favourable with CG_DESCENT, for a set consisting of 750 unconstrained optimization problems.

References

- [1] Andrei, N., *Scaled conjugate gradient algorithms for unconstrained optimization*. Computational Optimization and Applications, 38 (2007), pp. 401-416.
- [2] Andrei, N., Scaled memoryless BFGS preconditioned conjugate gradient algorithm for unconstrained optimization. Optimization Methods and Software 22 (2007), pp. 561-571.
- [3] Andrei, N., A scaled BFGS preconditioned conjugate gradient algorithm for unconstrained optimization. Applied Mathematics Letters 20 (2007), pp. 645-650.
- [4] Andrei, N., A scaled nonlinear conjugate gradient algorithm for unconstrained optimization. Optimization, 57 (2008), pp. 549-570.
- [5] Andrei, N., An acceleration of gradient descent algorithm with backtracking for unconstrained optimization. Numerical Algorithms, 42 (2006), pp.63-73.

- [6] Andrei, N., 40 conjugate gradient algorithms for unconstrained optimization. A survey on their definition. ICI Technical Report No. 13/08, March 14, 2008.
- [7] Andrei, N., Performance profiles of conjugate gradient algorithms for unconstrained optimization. Encyclopedia of Optimization, 2nd edition, C.A. Floudas and P.M. Pardalos (Eds.), Springer, New York, vol. P (2009), 2938-2953.
- [8] Andrei, N., *An unconstrained optimization test functions collection*. Advanced Modeling and Optimization. An Electronic International Journal, 10 (2008) 147-161.
- [9] Andrei, N., Accelerated conjugate gradient algorithm with finite difference Hessian / vector product approximation for unconstrained optimization. Journal of Computational and Applied Mathematics, 230 (2009), pp.570-582.
- [10] Andrei, N., Acceleration of conjugate gradient algorithms for unconstrained optimization. Applied Mathematics and Computation, 213 (2009), pp.361-369.
- [11] Birgin, E., Martínez, J.M., A spectral conjugate gradient method for unconstrained optimization, Applied Math. and Optimization, 43 (2001), pp.117-128.
- [12] Bongartz, I., Conn, A.R., Gould, N.I.M., Toint, Ph.L., CUTE: constrained and unconstrained testing environments, ACM Trans. Math. Software, 21 (1995), pp.123-160.
- [13] Dai, Y.H., New properties of a nonlinear conjugate gradient method. Numer. Math., 89 (2001), pp.83-98.
- [14] Dai, Y.H., Liao, L.Z., *New conjugacy con7 ditions and related nonlinear conjugate gradient methods*. Appl. Math. Optim., 43 (2001), pp. 87-101.
- [15] Dai, Y.H. Ni, Q., Testing different conjugate gradient methods for large-scale unconstrained optimization, J. Comput. Math., 21 (2003), pp.311-320.
- [16] Dai, Y.H., Yuan, Y., A nonlinear conjugate gradient method with a strong global convergence property, SIAM J. Optim., 10 (1999), pp.177-182.
- [17] Dai, Y.H. Yuan, Y., An efficient hybrid conjugate gradient method for unconstrained optimization. Annals of Operations Research, 103 (2001), pp.33-47.
- [18] Dai, Y.H., Han, J.Y., Liu, G.H., Sun, D.F., Yin, X., Yuan, Y., Convergence properties of nonlinear conjugate gradient methods. SIAM Journal on Optimization 10 (1999), 348-358.
- [19] Dolan, E.D., Moré, J.J., *Benchmarking optimization software with performance profiles*, Math. Programming, 91 (2002), pp. 201-213.
- [20] Hager, W.W., Zhang, H., A new conjugate gradient method with guaranteed descent and an efficient line search, SIAM Journal on Optimization, 16 (2005) 170-192.
- [21] Hager, W.W., Zhang, H., A survey of nonlinear conjugate gradient methods. Pacific journal of Optimization, 2 (2006), pp.35-58.
- [22] Hestenes, M.R., Stiefel, E.L., *Methods of conjugate gradients for solving linear systems*, J. Research Nat. Bur. Standards, 49 (1952) 409-436.
- [23] Liu, D.C. and Nocedal, J., On the limited memory BFGS method for large scale optimization. Mathematical Programming, 45 (1989), pp.503-528.
- [24] Moré, J.J., Thuente, D.J., *Line search algorithms with guaranteed sufficient decrease*. ACM Transactions on Mathematical Software, 20 (1994) 286-307.
- [25] Nocedal, J., Conjugate gradient methods and nonlinear optimization, in L. Adams and J.L. Nazareth (Eds.) Linear and Nonlinear Conjugate Gradient – Related Methods, SIAM, Philadelphia, 1996, pp.9-23.
- [26] Oren S.S., Spedicato, E., *Optimal conditioning of self-scaling variable metric algorithms*. Math. Programming, 10 (1976), pp.70-90.
- [27] Polak, E., Ribière, G., *Note sur la convergence de méthodes de directions conjuguée*, Revue Francaise Informat. Recherche Opérationnelle, 3e Année 16 (1969), pp.35-43.
- [28] Polyak, B.T., *The conjugate gradient method in extreme problems*, USSR Comp. Math. Math. Phys., 9 (1969), pp.94-112.
- [29] Powell, M.J.D., *Restart procedures for the conjugate gradient method*. Mathematical Programming 12 (1977), pp.241-254.

- [30] Shanno, D.F., Phua, V., Algorithm 500, Minimization of unconstrained multivariate functions, ACM Trans. on Math. Soft., 2, pp.87-94, 1976.
- [31] Wolfe, P., Convergence conditions for ascent methods, SIAM Rev., 11 (1969) pp.226-235.
- [32] Wolfe, P., *Convergence conditions for ascent methods II: some corrections*. SIAM Rev., 13 (1971) pp.185-188.
- [33] Zoutendijk, G., *Nonlinear programming computational methods*. In J. Abadie (Ed.) Integer and Nonlinear Programming, North-Holland, Amsterdam, 1970, pp.37-86.

January 27, 2010