

# Inexact Barzilai–Borwein method for saddle point problems

Yi-Qing Hu and Yu-Hong Dai\*,†

*State Key Laboratory of Scientific and Engineering Computing, Institute of Computational Mathematics and Scientific/Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, P.O. Box 2719, Beijing 100080, People's Republic of China*

## SUMMARY

This paper considers the inexact Barzilai–Borwein (BB) algorithm applied to saddle point problems. To this aim, we study the convergence properties of the inexact BB algorithm for symmetric positive definite linear systems. Suppose that  $g_k$  and  $\tilde{g}_k$  are the exact residual and its approximation of the linear system at the  $k$ th iteration, respectively. We prove the  $R$ -linear convergence of the algorithm if  $\|\tilde{g}_k - g_k\| \leq \eta \|\tilde{g}_k\|$  for some small  $\eta > 0$  and all  $k$ . To adapt the algorithm for solving saddle point problems, we also extend the  $R$ -linear convergence result to the case when the right-hand term  $\|\tilde{g}_k\|$  is replaced by  $\|\tilde{g}_{k-1}\|$ . Although our theoretical analyses cannot provide a good estimate to the parameter  $\eta$ , in practice, we find that  $\eta$  can be as large as the one in the inexact Uzawa algorithm. Further numerical experiments show that the inexact BB algorithm performs well for the tested saddle point problems. Copyright © 2007 John Wiley & Sons, Ltd.

Received 3 October 2005; Revised 25 June 2006; Accepted 30 June 2006

KEY WORDS: saddle point problem; Uzawa algorithm; Barzilai–Borwein method;  $R$ -linear convergence

## 1. INTRODUCTION

We consider the saddle point problem

$$\begin{pmatrix} A & B^T \\ B & -C \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ h \end{pmatrix} \quad (1)$$

\*Correspondence to: Yu-Hong Dai, State Key Laboratory of Scientific and Engineering Computing, Institute of Computational Mathematics and Scientific/Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, P.O. Box 2719, Beijing 100080, People's Republic of China.

†E-mail: dyh@lsec.cc.ac.cn

Contract/grant sponsor: Chinese NSF grants; contract/grant numbers: 10571171, 40233029, 10171104

where  $A \in \mathbb{R}^{n \times n}$  is symmetric positive definite and  $C \in \mathbb{R}^{m \times m}$  is symmetric positive semidefinite. This kind of problem arises frequently from the discretization of elasticity problems, Stokes equations and sometimes linearizations of Navier–Stokes equations. It also has a close relation to nonlinear programming since the problem of minimizing a convex quadratic subject to linear constraints can be converted into the form (1). There have been many methods developed for problem (1), see recent survey paper [1] and book [2]. In this paper, we are interested in a classic algorithm that is due to Uzawa [3]. It can be written as the following algorithm.

*Algorithm 1.1 (Uzawa)*

*Step 1. Initialize  $k=0$  and pick some  $p_0 \in \mathbb{R}^m$ ;*  
*Step 2. Solve  $Au_{k+1} = f - B^T p_k$  for  $u_{k+1} \in \mathbb{R}^n$ ;*  
*Step 3. Calculate  $p_{k+1} = p_k - \alpha(Cp_k - Bu_{k+1} + h)$ ;*  
*Step 4. If not convergent, set  $k = k + 1$  and go to Step 2.*

The elimination of  $u_{k+1}$  in the calculation of  $p_{k+1}$  leads to the iteration

$$p_{k+1} = p_k - \alpha[(BA^{-1}B^T + C)p_k - (BA^{-1}f - h)] \quad (2)$$

Therefore, the Uzawa algorithm is a fixed-parameter first-order Richardson iterative method [4] applied to the linear system

$$(BA^{-1}B^T + C)p = BA^{-1}f - h \quad (3)$$

In the context of optimization, the algorithm can be regarded as a fixed stepsize gradient method for the problem of minimizing a convex quadratic

$$\min \frac{1}{2}p^T \bar{A}p - \bar{b}^T p \quad (4)$$

where  $\bar{A} = BA^{-1}B^T + C$  and  $\bar{b} = (BA^{-1}f - h)$ .

The choice of the parameter  $\alpha$  is important to the efficacy of the Uzawa algorithm. Elman and Golub [5] proposed the following choice:

$$\alpha = \frac{2}{\lambda_1 + \lambda_m} \quad (5)$$

where  $\lambda_1$  and  $\lambda_m$  are the minimal and maximal eigenvalues of the matrix  $\bar{A}$ , respectively. This choice is optimal in the sense that it minimizes the spectral radius of the matrix  $I - \alpha\bar{A}$ . Since the eigenvalues  $\lambda_1$  and  $\lambda_m$  are not known to the users in general, Dai and Yang [6] chose the stepsize as follows:

$$\alpha_k = \frac{\|g_k\|}{\|\bar{A}g_k\|} \quad (6)$$

where  $g_k = \bar{A}p_k - \bar{b}$  and  $\|\cdot\|$  is the two norm. They proved that this sequence of  $\{\alpha_k\}$  tends to the value in (5). In practical computations, however, the gradient method with either (5) or (6) resembles the steepest descent method (see Cauchy [7]), where

$$\alpha_k^{SD} = \frac{g_k^T g_k}{g_k^T \bar{A}g_k} \quad (7)$$

They all become very slow as the condition number of the matrix  $\bar{A}$  deteriorates. Consequently, the use of the Uzawa algorithm is usually with some preconditioning technique. The Uzawa algorithm has received much attention from the numerical linear algebra community, for example, see [8–10].

In 1988, Barzilai and Borwein [11] proposed a different choice for the stepsize in the gradient method. Their basic idea is to regard  $D_k = \alpha_k^{-1} I$  as an approximation of the Hessian matrix  $\bar{A}$  and then to impose some certain quasi-Newton property on the matrix  $D_k$ . More exactly, they minimize  $\|D_k s_{k-1} - y_{k-1}\|_2$  where  $s_{k-1} = x_k - x_{k-1}$  and  $y_{k-1} = g_k - g_{k-1}$ , yielding the following choice of  $\alpha_k$ :

$$\alpha_k^{\text{BB}} = \frac{s_{k-1}^T s_{k-1}}{s_{k-1}^T y_{k-1}} \quad (8)$$

In the quadratic case, the above stepsize is equivalent to

$$\alpha_k = \frac{g_{k-1}^T g_{k-1}}{g_{k-1}^T \bar{A} g_{k-1}} \quad (9)$$

which happens to be the Cauchy stepsize (7) at the previous iteration.

Although the Barzilai–Borwein (BB) stepsize (8) cannot guarantee a descent in the objective function or the gradient norm, the corresponding method is proved to be globally convergent for strict convex quadratics (see Raydan [12]) and the convergence rate is  $R$ -linear (see Dai and Liao [13]). In the two-dimensional quadratic case, Barzilai and Borwein [11] presented a  $R$ -superlinear convergence result for the method. Dai and Fletcher [14] analysed the asymptotic convergence behaviour of the BB method for the higher-dimension case. In practical computations, it was pointed out in [11] that the BB stepsize (8) is far more efficient than the Cauchy stepsize (7). Fletcher [15] presented several linear systems of one million variables, showing that the BB method is comparable with the conjugate gradient method. The BB method has now received many generalizations and applications, for example, see [14, 16–21] and the references therein.

In this paper, we will apply the BB method to solve the saddle point problem (1). Each step of the Uzawa algorithm requires the solution of a symmetric positive definite linear system (see Step 2 of Algorithm 1.1). Elman and Golub [5] showed that this computation can be replaced by an approximate solution produced by an arbitrary iterative method, leading to the inexact Uzawa algorithm. The main purpose of this paper is to establish and analyse inexact BB algorithm for saddle point problems.

The rest of this paper is organized as follows. In the next section, we consider the inexact BB method where the exact gradient  $g_k$  is replaced by its some approximation  $\tilde{g}_k$ . Our study shows that there exists some small constant  $\eta > 0$ , which depends only on the problem dimension and the spectrum of the Hessian matrix, such that the BB method is  $R$ -linearly convergent for symmetric positive definite linear systems if  $\|\tilde{g}_k - g_k\| \leq \eta \|\tilde{g}_k\|$  for all  $k$ . In Section 3, we propose the inexact BB algorithm for the saddle point problem (1). To establish the  $R$ -linear convergence result of this algorithm, we extend the result of Section 2 to the case when the previous gradient norm  $\|\tilde{g}_{k-1}\|$  is used to control the inexactitude  $\|\tilde{g}_k - g_k\|$ . Although the estimate to  $\eta$  in our theoretical analyses can be very small, the numerical experiments in Section 4 show that this parameter  $\eta$  can be reasonably large without harming the convergence of the algorithm in practice. Further numerical results on some saddle point problems demonstrate the usefulness of the inexact BB algorithm. Conclusions and discussions are made in the last section.

## 2. INEXACT BARZILAI–BORWEIN METHOD

In this section, we consider the problem of minimizing a strictly convex quadratic

$$\min f(x) = \frac{1}{2}x^T Ax - b^T x \quad (10)$$

where  $A \in R^{n \times n}$  is symmetric positive definite and  $b \in R^n$ . To solve (10) we study the BB method with the gradient  $g(x) = \nabla f(x) = Ax - b$  computed inexactly and call the method *inexact BB method*. Assuming that  $\tilde{g}_k$  is an approximation to  $g_k$  at the  $k$ th iteration, the inexact BB algorithm for solving (10) can be described as follows.

*Algorithm 2.1 (Inexact BB)*

*Step 1. Initialize  $k = 0$  and pick some  $x_0 \in \mathbb{R}^n$ . Calculate some approximation  $\tilde{g}_0$  of the gradient  $g_0 = \nabla f(x_0)$  and set  $\tilde{\alpha}_0 = \tilde{g}_0^T \tilde{g}_0 / \tilde{g}_0^T A \tilde{g}_0$ ;*

*Step 2. Update  $x_{k+1} = x_k - \tilde{\alpha}_k \tilde{g}_k$  and  $k = k + 1$ ;*

*Step 3. Calculate some approximation  $\tilde{g}_k$  of the gradient  $g_k = \nabla f(x_k)$ ;*

*Step 4. Stop if some termination criterion is satisfied;*

*Step 5. Compute  $s_{k-1} = x_k - x_{k-1}$ ,  $\tilde{y}_{k-1} = \tilde{g}_k - \tilde{g}_{k-1}$  and  $\tilde{\alpha}_k = \frac{s_{k-1}^T s_{k-1}}{s_{k-1}^T \tilde{y}_{k-1}}$ , goto Step 2.*

In the above algorithm, the first stepsize  $\tilde{\alpha}_0$  is calculated by the steepest descent formula (7) with  $g_0$  replaced by an inexact gradient  $\tilde{g}_0$ . This is not expensive if the matrix–vector product  $A\tilde{g}_0$  can be used in computing  $\tilde{g}_1$ , as is the case of this paper.

Denote the error vector  $\zeta_k = \tilde{g}_k - g_k$ . Then, we have the basic relations

$$s_{k-1} = x_k - x_{k-1} = -\tilde{\alpha}_{k-1} \tilde{g}_{k-1} \quad (11)$$

$$g_k = g_{k-1} - \tilde{\alpha}_{k-1} A \tilde{g}_{k-1} \quad (12)$$

$$\tilde{g}_k = (I - \tilde{\alpha}_{k-1} A) \tilde{g}_{k-1} + \zeta_k - \zeta_{k-1} \quad (13)$$

Further, still denoting  $y_{k-1} = g_k - g_{k-1}$ , we have that

$$\tilde{y}_{k-1} = y_{k-1} + \zeta_k - \zeta_{k-1} \quad (14)$$

We are going to analyse Algorithm 2.1 under the condition

$$\|\zeta_k\| \leq \eta \|\tilde{g}_k\| \quad (15)$$

where  $\eta \in (0, 1)$  is some positive constant. It follows from (15) and the definition of  $\zeta_k$  that

$$g_k^T \tilde{g}_k = \tilde{g}_k^T \tilde{g}_k + (g_k - \tilde{g}_k)^T \tilde{g}_k \geq \|\tilde{g}_k\|^2 - \|\zeta_k\| \|\tilde{g}_k\| \geq (1 - \eta) \|\tilde{g}_k\|^2$$

Therefore, we can see that condition (15) ensures the descent property of  $-\tilde{g}_k$  unless  $\tilde{g}_k = 0$ .

Suppose that the eigenvalues of the Hessian matrix  $A$  are

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \quad (16)$$

The following theorem claims that Algorithm 2.1 is well defined if the parameter  $\eta$  satisfies

$$\eta \leq \frac{1}{9} \left( \frac{\lambda_1}{\lambda_n} \right)^2 =: c_1 \quad (17)$$

*Theorem 2.2*

Consider the inexact BB algorithm, namely, Algorithm 2.1 under the assumption (15). If  $\eta$  satisfies (17), the following relations:

$$\frac{2}{\lambda_1 + 2\lambda_n} \leq \tilde{\alpha}_k \leq \frac{2}{\lambda_1} \quad (18)$$

and

$$\|\tilde{g}_{k+1}\| \leq c_2 \|\tilde{g}_k\| \quad (19)$$

where  $c_2 = 9\lambda_n/4\lambda_1 - 1$ , hold for all  $k \geq 0$ . Therefore, the algorithm is well defined.

*Proof*

We establish (18) and (19) by induction. If  $k = 0$ , we have directly by the choice  $\tilde{\alpha}_0 = \tilde{g}_0^T \tilde{g}_0 / \tilde{g}_0^T A \tilde{g}_0$  that  $1/\lambda_n \leq \tilde{\alpha}_0 \leq 1/\lambda_1$ . Thus, (18) holds with  $k = 0$ . Further, fixing  $k = 0$ , we have by (18) that

$$\|I - \tilde{\alpha}_k A\| \leq \max \left\{ 1 - \frac{2\lambda_1}{\lambda_1 + 2\lambda_n}, \frac{2\lambda_n}{\lambda_1} - 1 \right\} = \frac{2\lambda_n}{\lambda_1} - 1 \quad (20)$$

It follows from (13) that

$$\begin{aligned} \|\tilde{g}_{k+1}\| &\leq \|(I - \tilde{\alpha}_k A)\tilde{g}_k\| + \|\zeta_{k+1}\| + \|\zeta_k\| \\ &\leq \|I - \tilde{\alpha}_k A\| \|\tilde{g}_k\| + \eta \|\tilde{g}_{k+1}\| + \eta \|\tilde{g}_k\| \end{aligned} \quad (21)$$

Using (21), (20) and noting that  $\eta \in (0, \frac{1}{9})$ , we can obtain

$$\begin{aligned} \|\tilde{g}_{k+1}\| &\leq (1 - \eta)^{-1} \left[ \frac{2\lambda_n}{\lambda_1} - 1 + \eta \right] \|\tilde{g}_k\| \\ &\leq \left[ \frac{2\lambda_n}{\lambda_1} (1 - \eta)^{-1} - 1 \right] \|\tilde{g}_k\| \\ &\leq \left[ \frac{9\lambda_n}{4\lambda_1} - 1 \right] \|\tilde{g}_k\| \end{aligned} \quad (22)$$

Thus, by the choice of  $c_2$ , we know that (19) is true with  $k = 0$ .

Now we assume that (18) and (19) hold for all  $l \leq k$ , where  $k$  is some integer that satisfies  $k \geq 0$ . Then by (14), (15), the induction assumption and (17), we have the following estimate:

$$\begin{aligned} \left| \frac{s_k^T \tilde{y}_k}{s_k^T s_k} - \frac{s_k^T y_k}{s_k^T s_k} \right| &\leq \left| \frac{s_k^T (\zeta_{k+1} - \zeta_k)}{s_k^T s_k} \right| \leq \frac{\|\zeta_{k+1}\| + \|\zeta_k\|}{\|s_k\|} \\ &\leq \frac{(1 + c_2)\eta \|\tilde{g}_k\|}{\tilde{\alpha}_k \|\tilde{g}_k\|} = \frac{9\lambda_n \eta}{4\lambda_1 \tilde{\alpha}_k} \\ &\leq \frac{9\lambda_n}{4\lambda_1} \cdot \frac{\lambda_1 + 2\lambda_n}{2} \cdot \frac{1}{9} \left( \frac{\lambda_1}{\lambda_n} \right)^2 < \frac{\lambda_1}{2} \end{aligned} \quad (23)$$

On the other hand, we can see from the definition of  $y_k$  that  $\lambda_1 \leq s_k^T y_k / s_k^T s_k \leq \lambda_n$ . It follows from this and (23) that

$$\frac{\lambda_1}{2} \leq \left| \frac{s_k^T \tilde{y}_k}{s_k^T s_k} \right| \leq \lambda_n + \frac{\lambda_1}{2} \quad (24)$$

Thus,  $\tilde{\alpha}_{k+1}$  is well defined and (18) holds with  $k$  replaced by  $k+1$ . Further, it is not difficult to see that the deductions from (20) to (22) are still available and hence (19) is true when the index  $k$  is replaced with  $k+1$ . Therefore by induction, the relations (18) and (19) hold for all  $k \geq 0$ . Algorithm 2.1 is then well defined.  $\square$

In the case that (15) and (17) hold, we can see from Theorem 2.2 that the inexact BB algorithm is well defined provided that the initial stepsize  $\alpha_0$  satisfies the relation (18). Theorem 2.2 also tells us that to guarantee the well-definition of the algorithm, the constant  $\eta$  need be less than the square of the inverse of the condition number  $\lambda_n/\lambda_1$  of the matrix  $A$ . The following theorem extends the  $R$ -linear convergence result in [13] of the BB method.

### Theorem 2.3

Consider the inexact BB algorithm, namely, Algorithm 2.1 under the conditions (15) and (17). There exists some positive constant  $c_3 \leq c_1$ , which depends only on  $\lambda_1$ ,  $\lambda_n$  and the dimension  $n$ , such that if  $\eta \leq c_3$ , the algorithm either gives the solution in finite iterations or converges to the solution  $R$ -linearly.

### Proof

To prove the theorem, we compare the inexact BB algorithm and the (exact) BB method. At first, note that the (exact) BB iterations are uniquely decided by its starting point and initial stepsize. By Lemma 2.4 in [13] and some slight modifications, we can see that there exists some positive integer  $M$ , which depends only on  $\lambda_1$ ,  $\lambda_n$  and the problem dimension  $n$ , such that for any starting point  $z_0 \in \mathbb{R}^n$  and initial stepsize  $\beta_0$  satisfying  $c_L \leq \beta_0 \leq c_U$  ( $c_L$  and  $c_U$  are some fixed positive constants), the  $M$ th point  $z_M$  generated by the (exact) BB method satisfies

$$\|\nabla f(z_M)\| \leq \frac{1}{2} \|\nabla f(z_0)\| \quad (25)$$

Let us fix  $c_L = 2/(\lambda_1 + 2\lambda_n)$  and  $c_U = 2/\lambda_1$  and take some  $M$  satisfying the above statement. For any point  $x_k$  and stepsize  $\tilde{\alpha}_k$  generated by Algorithm 2.1, we consider the (exact) BB iterations  $\{z_{k+l}; l \geq 0\}$  with  $z_k = x_k$  and  $\beta_k = \tilde{\alpha}_k$ . The gradient of  $f$  at  $z_{k+l}$  is denoted by  $h_{k+l}$ . It is easy to see that  $h_k = g_k$  and for  $l \geq 0$ ,

$$h_{k+l+1} = (I - \beta_{k+l} A) h_{k+l} \quad \text{and} \quad \beta_{k+l+1} = \frac{h_{k+l}^T h_{k+l}}{h_{k+l}^T A h_{k+l}} \quad (26)$$

Further, we take the smallest integer  $m' \leq M$  such that  $\|h_{k+m'}\| \leq \frac{1}{2} \|h_k\|$ . In this case, we have that

$$\|h_{k+l}\| > \frac{1}{2} \|h_k\| \quad \text{for all } l = 0, 1, \dots, m' - 1 \quad (27)$$

We now consider the quantities  $\phi_{k+l} = \|\tilde{g}_{k+l} - h_{k+l}\|$  and  $\psi_{k+l} = |\tilde{\alpha}_{k+l} - \beta_{k+l}|$ . Using (13), (26), (15) and (19), we get that for  $l \geq 0$

$$\begin{aligned} \phi_{k+l+1} &\leq \|(I - \beta_{k+l}A)h_{k+l} - (I - \tilde{\alpha}_{k+l})A\tilde{g}_{k+l}\| + \|\zeta_{k+l+1}\| + \|\zeta_{k+l}\| \\ &\leq \|I - \beta_{k+l}A\|\phi_{k+l} + \|A\|\|\tilde{g}_{k+l}\|\psi_{k+l} + \eta(1 + c_2)\|\tilde{g}_{k+l}\| \\ &\leq \left(\frac{\lambda_n}{\lambda_1} - 1\right)\phi_{k+l} + \lambda_n\|\tilde{g}_{k+l}\|\psi_{k+l} + \eta(1 + c_2)\|\tilde{g}_{k+l}\| \end{aligned} \quad (28)$$

On the other hand, we have by direct calculations that

$$\psi_{k+l+1} = \left| \frac{\tau_{k+l} + \|h_{k+l}\|^2\theta_{k+l}}{h_{k+l}^T A h_{k+l} [\tilde{g}_{k+l}^T A \tilde{g}_{k+l} + \theta_{k+l}]} \right| \quad (29)$$

where

$$\tau_{k+l} = \|h_{k+l}\|^2 \tilde{g}_{k+l}^T A \tilde{g}_{k+l} - h_{k+l}^T A h_{k+l} \|\tilde{g}_{k+l}\|^2 \quad (30)$$

and

$$\theta_{k+l} = \tilde{\alpha}_{k+l}^{-1} \tilde{g}_{k+l}^T (\zeta_{k+l} - \zeta_{k+l+1}) \quad (31)$$

For  $\tau_{k+l}$ , we have the estimate

$$\begin{aligned} |\tau_{k+l}| &\leq |\tilde{g}_{k+l}^T A \tilde{g}_{k+l} [\|h_{k+l}\|^2 - \|\tilde{g}_{k+l}\|^2] + \|\tilde{g}_{k+l}\|^2 [\tilde{g}_{k+l}^T A \tilde{g}_{k+l} - h_{k+l}^T A h_{k+l}]]| \\ &\leq [\tilde{g}_{k+l}^T A \tilde{g}_{k+l} (\|h_{k+l}\| + \|\tilde{g}_{k+l}\|) + \|\tilde{g}_{k+l}\|^2 \|A(h_{k+l} + \tilde{g}_{k+l})\|] \phi_{k+l} \\ &\leq 2\lambda_n \|\tilde{g}_{k+l}\|^2 (\|h_{k+l}\| + \|\tilde{g}_{k+l}\|) \phi_{k+l} \end{aligned} \quad (32)$$

For  $\theta_{k+l}$ , we have the estimate

$$|\theta_{k+l}| \leq \frac{(1 + c_2)(\lambda_1 + 2\lambda_n)}{2} \eta \|\tilde{g}_{k+l}\|^2 \leq \frac{\lambda_1}{2} \|\tilde{g}_{k+l}\|^2 \quad (33)$$

Denote the constant  $c_4 = (1 + c_2)(\lambda_1 + 2\lambda_n)/\lambda_1$ . Using (30) and (33), we can get from (29)

$$\begin{aligned} \psi_{k+l+1} &\leq \frac{2\lambda_n \|\tilde{g}_{k+l}\|^2 (\|h_{k+l}\| + \|\tilde{g}_{k+l}\|) \phi_{k+l} + (c_4 \lambda_1 / 2) \eta \|h_{k+l}\|^2 \|\tilde{g}_{k+l}\|^2}{h_{k+l}^T A h_{k+l} [\tilde{g}_{k+l}^T A \tilde{g}_{k+l} - (\lambda_1 / 2) \|\tilde{g}_{k+l}\|^2]} \\ &\leq \frac{2\lambda_n (\|h_{k+l}\| + \|\tilde{g}_{k+l}\|) \phi_{k+l} + (c_4 \lambda_1 / 2) \eta \|h_{k+l}\|^2}{(\lambda_1 / 2) \|h_{k+l}\|^2} \\ &\leq \frac{4\lambda_n (\|\tilde{g}_{k+l}\| + \|h_{k+l}\|) \phi_{k+l}}{\lambda_1 \|h_{k+l}\|^2} + c_4 \eta \end{aligned} \quad (34)$$

It follows from (19) that  $\|\tilde{g}_{k+l}\| \leq c_2^l \|\tilde{g}_k\|$ . By (26), we can deduce that  $\|h_{k+l}\| \leq c_5^l \|h_k\|$ , where  $c_5 = \max\{1, \lambda_n/\lambda_1 - 1\}$  is constant. Since  $h_k = g_k$ , we have from (15) and  $\eta \leq c_1 < \frac{1}{9}$  that  $\|\tilde{g}_k\| \leq$

$\frac{9}{8}\|h_k\|$ . Using these analyses, the relation (27), (34) and  $m' \leq M$ , we can obtain the estimate

$$\psi_{k+l+1}\|\tilde{g}_{k+l+1}\| \leq c_6\phi_{k+l} + c_4\eta\|\tilde{g}_{k+l+1}\| \quad \text{for } l = 0, 1, \dots, m' - 2 \quad (35)$$

where  $c_6$  is the constant given by

$$c_6 = \frac{9\lambda_n(9c_2^{M-2} + 8c_5^{M-2})c_2^{M-1}}{4\lambda_1} \quad (36)$$

Since the first stepsize  $\beta_k$  is chosen to be  $\tilde{\alpha}_k$ , we have that  $\psi_k = 0$ . In addition, (15) implies that  $\phi_k \leq \eta\|\tilde{g}_k\|$ . Thus, we have from (28)

$$\phi_{k+1} \leq \left(\frac{\lambda_n}{\lambda_1} - 1\right)\phi_k + \eta(1 + c_2)\|\tilde{g}_k\| \leq \left(c_2 + \frac{\lambda_n}{\lambda_1}\right)\eta\|\tilde{g}_k\| \quad (37)$$

For  $l \in [1, m' - 1]$ , we have by (28) and (35) (with  $l$  replaced by  $l - 1$ ) and  $\|\tilde{g}_{k+l}\| \leq c_2^M\|\tilde{g}_k\|$  that

$$\phi_{k+l+1} \leq \left(\frac{\lambda_n}{\lambda_1} - 1\right)\phi_{k+l} + c_6\lambda_n\phi_{k+l-1} + (1 + c_2 + c_4\lambda_n)\eta\|\tilde{g}_{k+l}\| \quad (38)$$

Note that all the constants  $c_1, c_2, c_4, c_5, c_6$  and the integer  $M$  are dependent only on  $\lambda_1, \lambda_n$  and possibly the dimension  $n$ . From  $\phi_k \leq \eta\|\tilde{g}_k\|$ , (37), (38) and  $\|\tilde{g}_{k+l}\| \leq c_2^M\|\tilde{g}_k\|$ , we know that there exists some constant  $c_7$ , which depends only on  $\lambda_1, \lambda_n$  and the dimension  $n$ , such that

$$\phi_{k+m'} \leq c_7\eta\|\tilde{g}_k\| \quad (39)$$

Taking  $c_3 = \min\{1/(3 + 6c_7), c_1\}$ , we obtain from (39),  $\|h_{k+m'}\| \leq \frac{1}{2}\|h_k\|$ ,  $\|h_k\| \leq (1 + \eta)\|\tilde{g}_k\|$  and  $\eta \leq c_3$  that

$$\|\tilde{g}_{k+m'}\| \leq \|h_{k+m'}\| + \phi_{k+m'} \leq \left[\frac{(1 + \eta)}{2} + c_7\eta\right]\|\tilde{g}_k\| \leq \frac{2}{3}\|\tilde{g}_k\| \quad (40)$$

To complete the proof, we define a subsequence  $\{k_i\}$  with  $k_1 = 2$  for Algorithm 2.1. If  $k_i$  has been decided, we choose  $k_{i+1} = k_i + m_i$ , where  $m_i \in [1, m]$  is so chosen

$$\|\tilde{g}_{k_{i+1}}\| \leq \frac{2}{3}\|\tilde{g}_{k_i}\| \quad (41)$$

By the analysis in the previous paragraph, we know that this is possible. It then follows that

$$\|\tilde{g}_{k_i}\| \leq \left(\frac{2}{3}\right)^{i-1}\|\tilde{g}_{k_1}\| \quad (42)$$

with  $k_i = k_1 + \sum_{j=1}^{i-1} m_j \leq k_1 + M(i - 1)$ . Consequently, we have that

$$\limsup_{i \rightarrow \infty} \|\tilde{g}_{k_i}\|^{1/k_i} \leq \left(\frac{2}{3}\right)^{1/M} \quad (43)$$

From the above relation and (19), we know that  $\|\tilde{g}_k\|$  and hence  $\|g_k\|$  converges to zero  $R$ -linearly.  $\square$

The importance of Theorems 2.2 and 2.3 is in that, to guarantee the well-definition of Algorithm 2.1 and inherent the  $R$ -linear convergence of the (exact) BB method, the calculation error of the gradient can be less than some constant proportion of the gradient norm. The



constant depends only on the dimension  $n$  and the minimal and maximal eigenvalues of the matrix  $A$ . However, the current estimate to the constant  $\eta$  in the proof of Theorem 2.3 may be very small, since the integer  $M$  can be very large. In practice, the choice of the value  $\eta$  is optimistic, as will be seen in our numerical experiments of Section 4.

As one application of Theorem 2.3, we can show that the (exact) BB method is locally  $R$ -linearly convergent for twice continuously differentiable functions. Suppose that  $f(x)$  is the function to be minimized and  $x^*$  is a point at which  $\nabla f(x^*)$  and its Hessian  $H^*$  is positive definite. Then at some neighbourhood of  $x^*$ , the (exact) BB method for the minimization of  $f(x)$  can be regarded as the inexact BB method for minimizing the following quadratic:

$$q(x) = f(x^*) + \frac{1}{2}(x - x^*)^T H^* (x - x^*) \quad (44)$$

In the case that  $f$  is twice continuously differentiable, it is not difficult to establish the relation

$$\|\nabla f(x) - \nabla q(x)\| = o(\|\nabla f(x)\|) \quad (45)$$

Hence the condition (15) must be satisfied when  $x_k$  tends to  $x^*$  and hence  $R$ -linear convergence can be established. This remark weakens the assumption that the objection function  $f$  is two times Lipschitz continuously differentiable for the CBB method in [22], in which case the following relation holds:

$$\|\nabla f(x) - \nabla q(x)\| = O(\|\nabla f(x)\|^2) \quad (46)$$

### 3. INEXACT BB METHOD FOR SADDLE POINT PROBLEMS

The Uzawa algorithm for the saddle point problem (1) requires the solution of a linear system at each iteration (see Step 2 of Algorithm 1.1). In practice, it is usually expensive to solve the subproblem exactly. Elman and Golub [5] proposed to replace Step 2 of Algorithm 1.1 by

$$Au_{k+1} = f - B^T p_k + \delta_k \quad (47)$$

where the vector  $\delta_k$  is the residual of the approximation solution  $u_{k+1}$  to the system  $Av = f - B^T p_k$ . They suggested that a natural choice for the magnitude of  $\delta_k$  is

$$\|\delta_k\| \leq \tau \|Cp_{k-1} - Bu_k + h\| \quad (48)$$

This is because the quantity  $Cp_{k-1} - Bu_k + h$  is the residual of the second block row of (1) for the approximation solution pair  $(u_k, p_{k-1})$  and this quantity has already been calculated for the update of  $p_k$  in the previous step.

If the subproblem at Step 1 of Algorithm 1.1 is exactly solved, we obtain the exact residual to system (3):

$$g_k = \bar{A}p_k - \bar{b} = Cp_k - BA^{-1}(f - B^T p_k) + h \quad (49)$$

When the subproblem is solved inexactly by (47), the inexact gradient  $\tilde{g}_k$  can be written as

$$\tilde{g}_k = Cp_k - Bu_{k+1} + h = g_k - BA^{-1}\delta_k \quad (50)$$

From (48) and the first equality of (50) with  $k$  replaced by  $k - 1$ , the error vector  $\delta_k$  is required to satisfy

$$\|\delta_k\| \leq \tau \|\tilde{g}_{k-1}\| \quad (51)$$

Combining the BB method and the inexact idea of Golub and Elman, we give an inexact BB algorithm for saddle point problem (1).

*Algorithm 3.1 (Inexact BB algorithm for saddle point problems)*

*Step 1. Initialize  $k=0$  and  $p_0 \in \mathbb{R}^m$ . Choose some big constant  $\rho > 1$  and some initial stepsize  $\tilde{\alpha}_0 \in [\rho^{-1}, \rho]$ ;*  
*Step 2. Compute  $u_{k+1}$  such that  $Au_{k+1} = f - B^T p_k + \delta_k$  with  $\delta_k$  satisfying (48);*  
*Step 3. Compute  $\tilde{g}_k = Cp_k - Bu_{k+1} + h$ ,  $p_{k+1} = p_k - \tilde{\alpha}_k \tilde{g}_k$  and set  $k = k + 1$ ;*  
*Step 4. Stop if some termination criterion is satisfied;*  
*Step 5. Calculate  $s_{k-1} = p_k - p_{k-1}$  and  $\tilde{y}_{k-1} = \tilde{g}_k - \tilde{g}_{k-1}$ . Compute the next stepsize  $\tilde{\alpha}_k$  by  $[\max\{\rho^{-1}, \min\{\frac{s_{k-1}^T \tilde{y}_{k-1}}{s_{k-1}^T s_{k-1}}, \rho\}\}]^{-1}$  and goto Step 2.*

The introduction of the constant  $\rho > 1$  ensures that the stepsize  $\tilde{\alpha}_k$  and hence Algorithm 3.1 is well defined. The previous section considers the inexact BB algorithm under the assumption (15). For Algorithm 3.1, however, we have by (50) and (51) that

$$\|\zeta_k\| = \|\tilde{g}_k - g_k\| \leq \bar{\eta} \|\tilde{g}_{k-1}\| \quad (52)$$

where  $\bar{\eta} = \tau \|BA^{-1}\|$ . Since it is possible that  $\|g_k\|$  can be arbitrarily smaller than  $\|g_{k-1}\|$  in the (exact) BB method and hence it is likely that  $\|\tilde{g}_k\|$  is far smaller than  $\|\tilde{g}_{k-1}\|$ , (52) does not imply (15) for any small constant  $\tau$ . Therefore, we cannot establish the  $R$ -linear convergence of Algorithm 3.1 directly from Theorem 2.3. At the same time, we can see that unlike (15), condition (52) cannot ensure the descent property of  $-\tilde{g}_k$  at every iteration.

The above difficulties can be circumvented by noting that if  $\|\tilde{g}_k\|$  is significantly less than  $\|\tilde{g}_{k-1}\|$ , then a good approximation  $\tilde{g}_k$  of the gradient  $g_k$  has been obtained. On the other hand, it follows from (13),  $\tilde{\alpha}_k \in [\rho^{-1}, \rho]$  and (52) that

$$\|\tilde{g}_k\| \leq \lambda_n \rho \|\tilde{g}_{k-1}\| + \bar{\eta} \|\tilde{g}_{k-2}\| \quad (53)$$

Here and below we assume that  $\rho$  is a very large constant such that

$$\rho \geq \max \left\{ \frac{2}{\lambda_1}, \lambda_n + \frac{\lambda_1}{2} \right\} \quad (54)$$

The above relation and  $\lambda_1 \leq \lambda_n$  imply that  $\lambda_n \rho \geq 2$ . From (52), we have

$$\|g_k\| \leq \|\tilde{g}_k\| + \bar{\eta} \|\tilde{g}_{k-1}\| \quad (55)$$

Relations (53) and (55) hint that, to establish the  $R$ -linear convergence of Algorithm 3.2, we need to consider a subsequence of  $k_i$  such that the approximation gradients  $\tilde{g}_{k_i-1}$  and  $\tilde{g}_{k_i}$  have some properties simultaneously.

*Theorem 3.2*

Consider Algorithm 3.1 for saddle point problem (1) under conditions (15) and (17). Assume that  $\rho$  is a big constant that satisfies (54). Then there exists some positive constant  $\tau_1$ , which depends only on  $\lambda_1$ ,  $\lambda_n$  and the dimension  $n$ , such that if (51) holds for all  $k$  and  $\tau \leq \tau_1$ , the algorithm either gives the solution in finite iterations or converges to the solution  $R$ -linearly.

*Proof*

Similarly to the proof of Theorem 2.3, we find some constant  $c_3 \in (0, 1)$  and integer  $m$  which depend only on  $\lambda_1$ ,  $\lambda_n$ , the dimension  $n$  and the parameter  $\rho$  such that, if  $\|\xi_j\| \leq \eta \|\tilde{g}_j\|$  for all  $j$  and if  $\eta \leq c_3$ , then for any index  $k$ , there exists some integer  $m' \leq M$  satisfying  $\|\tilde{g}_{k+m'}\| \leq \frac{2}{3} \|\tilde{g}_k\|$ . Denote  $M_1 = \lceil \log(2\lambda_n \rho) / \log 1.5 \rceil M$ . Then if  $\|\xi_j\| \leq \eta \|\tilde{g}_j\|$  for all  $j$  and if  $\eta \leq c_3$ , then for any index  $k$ , there exists some integer  $m' \leq M_1$  such that

$$\|\tilde{g}_{k+m'}\| \leq \left(\frac{2}{3}\right)^{M_1/M} \|\tilde{g}_k\| \leq \frac{1}{2\lambda_n \rho} \|\tilde{g}_k\| \quad (56)$$

Now, we denote the constants

$$c_8 = \frac{1}{2\lambda_n \rho (\lambda_n \rho + 1)^{M_1-1}}, \quad \tau_1 = c_3 c_8 \|BA^{-1}\|^{-1} \quad (57)$$

and define the subsequence  $\{k_i\}$  in the following way. Pick the least index  $k_1 \geq 1$  such that  $\|\tilde{g}_{k_1}\| \geq \frac{2}{3} \|\tilde{g}_{k_1-1}\|$ . If this is not possible, we have  $\|\tilde{g}_k\| \leq \frac{2}{3} \|\tilde{g}_{k-1}\|$  for all  $k \geq 1$  and  $\{\|\tilde{g}_k\|\}$  is a  $Q$ -linear convergence sequence. Assume that for some  $i \geq 1$ ,  $k_i$  has been chosen with the property

$$\|\tilde{g}_{k_i}\| \geq \frac{2}{3} \|\tilde{g}_{k_i-1}\| \quad (58)$$

By this, (55) and  $\bar{\eta} \leq \frac{2}{3}$  and using the induction principle, it is not difficult to show that

$$\|\tilde{g}_{k_i+l}\| \leq (\lambda_n \rho + 1)^l \|\tilde{g}_{k_i}\| \quad \text{for all } l \geq 1 \quad (59)$$

Thus, if  $\|\tilde{g}_{k_i+l}\| < c_8 \|\tilde{g}_{k_i+l-1}\|$  for some  $l \in [1, M_1]$ , we have by this, (59) and the choice of  $c_8$  that  $\|\tilde{g}_{k_i+l}\| \leq 1/2\lambda_n \rho \|\tilde{g}_{k_i}\|$ . If this is not the case, we have  $\|\tilde{g}_{k_i+l}\| \geq c_8 \|\tilde{g}_{k_i+l-1}\|$  for  $l = 1, \dots, M_1$ . It follows from this, (52) and (57) that  $\|\xi_{k_i+l}\| \leq c_3 \|\tilde{g}_{k_i+l}\|$  for  $l = 1, \dots, M_1$ . By the statement in the first paragraph of the proof and the choice of  $\bar{\eta}$ , we must have that  $\|\tilde{g}_{k_i+M_1}\| \leq 1/2\lambda_n \rho \|\tilde{g}_{k_i}\|$ . Therefore, there always exists some integer  $m' \in [1, M_1]$  such that

$$\|\tilde{g}_{k_i+m'}\| \leq \frac{1}{2\lambda_n \rho} \|\tilde{g}_{k_i}\| \quad (60)$$

We then take  $k_{i+1}$  to be the least integer that is not less than  $k_i + m' + 1$  and satisfies  $\|\tilde{g}_{k_{i+1}}\| \geq \frac{2}{3} \|\tilde{g}_{k_{i+1}-1}\|$ . If  $k_{i+1} = +\infty$ , we know that  $\{\|\tilde{g}_k\|\}$  is  $R$ -linearly convergent. Therefore, we assume that the above-defined  $\{k_i\}$  is an infinite sequence.

If  $k_{i+1} = k_i + m' + 1$ , we have by (53), (60), (59), (57),  $m' \leq M_1$  and the definition of  $\bar{\eta}$  that

$$\|\tilde{g}_{k_{i+1}}\| \leq \lambda_n \rho \|\tilde{g}_{k_i+m'}\| + \bar{\eta} \|\tilde{g}_{k_i+m'-1}\|$$

$$\begin{aligned}
&\leq \lambda_n \rho \left( \frac{1}{2\lambda_n \rho} \|\tilde{g}_{k_i}\| \right) + \bar{\eta}(\lambda_n \rho + 1)^{m'-1} \|\tilde{g}_{k_i}\| \\
&\leq \left[ \frac{1}{2} + \frac{c_3}{2\lambda_n \rho} \right] \|\tilde{g}_{k_i}\| \leq \frac{3}{4} \|\tilde{g}_{k_i}\|
\end{aligned} \tag{61}$$

The facts that  $c_3 \in (0, 1)$  and  $\lambda_n \rho \geq 2$  are also used for the last inequality.

If  $k_{i+1} > k_i + m' + 1$ , we denote  $j_0 = k_{i+1} - (k_i + m' + 1)$ . The choice of  $k_{i+1}$  implies that

$$\|\tilde{g}_{k_i+m'+j}\| \leq \left(\frac{2}{3}\right)^j \|\tilde{g}_{k_i+m'}\| \quad \text{for } j = 1, \dots, j_0 \tag{62}$$

It follows by (53), (62), (57) and the fact that  $3\bar{\eta} \leq \lambda_n \rho$  that

$$\begin{aligned}
\|\tilde{g}_{k_{i+1}}\| &\leq \lambda_n \rho \|\tilde{g}_{k_i+m'+j_0}\| + \bar{\eta} \|\tilde{g}_{k_i+m'+j_0-1}\| \\
&\leq \left(\frac{2}{3}\right)^{j_0} \lambda_n \rho \|\tilde{g}_{k_i+m'}\| + \bar{\eta} \left(\frac{2}{3}\right)^{j_0-1} \|\tilde{g}_{k_i+m'}\| \\
&\leq \left(\frac{2}{3}\right)^{j_0} \left( \frac{1}{2} \|\tilde{g}_{k_i}\| \right) + \bar{\eta} \left(\frac{2}{3}\right)^{j_0-1} \left( \frac{1}{2\lambda_n \rho} \|\tilde{g}_{k_i}\| \right) \\
&\leq \left(\frac{2}{3}\right)^{j_0} \left[ \frac{1}{2} + \frac{3\bar{\eta}}{4\lambda_n \rho} \right] \|\tilde{g}_{k_i}\| \leq \left(\frac{2}{3}\right)^{j_0} \left[ \frac{3}{4} \|\tilde{g}_{k_i}\| \right]
\end{aligned} \tag{63}$$

Combining the two possible cases of  $k_{i+1}$  and noting that  $m' \leq M_1$ , we always have the following relation:

$$\|\tilde{g}_{k_{i+1}}\| \leq \frac{3}{4} \left(\frac{2}{3}\right)^{\max\{0, k_{i+1}-k_i-M_1-1\}} \|\tilde{g}_{k_i}\| \tag{64}$$

Denote the constant  $c_9 = \left(\frac{3}{4}\right)^{1/(M_1+1)}$ . Since  $M_1 \geq 1$ , we have that  $\frac{2}{3} < c_9 < 1$ . With the choice of  $c_9$ , we can show by (64) that  $\|\tilde{g}_{k_{i+1}}\| \leq c_9^{k_{i+1}-k_i} \|\tilde{g}_{k_i}\|$ . The recursion of this relation leads to

$$\|\tilde{g}_{k_{i+1}}\| \leq c_9^{k_{i+1}-k_1} \|\tilde{g}_{k_1}\| \tag{65}$$

In addition, by (59), the definition of  $k_{i+1}$  and  $m' \leq M_1$ , it is not difficult to see that

$$\max\{\|\tilde{g}_{k_i+1}\|, \dots, \|\tilde{g}_{k_{i+1}-1}\|\} \leq (\lambda_n \rho + 1)^{M_1-1} \|\tilde{g}_{k_i}\| \tag{66}$$

Therefore, we know by (65) and (66) that

$$\limsup_{k \rightarrow \infty} \|\tilde{g}_k\|^{1/k} \leq c_9 < 1 \tag{67}$$

which implies that  $\{\|\tilde{g}_k\|\}$  and hence by (55),  $\{\|g_k\|\}$  are  $R$ -linearly convergent.  $\square$

Again, the proof to Theorem 3.2 provides a pessimistic estimate to the largest admissible value to  $\tau_1$ . Nevertheless, the numerical experiments in the next section show that  $\tau_1$  and hence  $\bar{\eta}$  can be much larger.

## 4. NUMERICAL EXPERIMENTS

Our numerical experiments in this section are divided into two parts. In the first part, we test Algorithm 2.1 with the inexactitude case (52) for random symmetric and positive definite linear systems. Specifically, we observe how the value of  $\bar{\eta}$  in (52) influences the performance of the algorithm and how its choice depends on the problem dimension  $n$  and the smallest eigenvalue  $\lambda_1$  and the condition number  $\kappa$  of the matrix  $A$ .

*Problem 4.1*

Consider the symmetric positive definite linear system  $Ax = b$ , where  $x \in \mathbb{R}^{n \times n}$ . The coefficient matrix  $A$  is formed by

$$A = PDP^T \quad \text{where } P = (1 - 2\omega_1\omega_1^T)(1 - 2\omega_2\omega_2^T)(1 - 2\omega_3\omega_3^T) \quad (68)$$

and  $\omega_1, \omega_2$  and  $\omega_3$  are unit vectors generated by the uniform distribution in  $\mathbb{R}^n$  and  $D$  is a diagonal matrix. Given the smallest eigenvalue  $\lambda_1$  and the condition number  $\kappa$ , the  $i$ th diagonal entry  $D_{i,i}$  of the matrix  $D$  is set to

$$D_{i,i} = \exp\left(\log \lambda_1 + \frac{i-1}{n-1} \log \kappa\right) \quad (69)$$

To generate the right-hand term  $b$ , we randomly generate a solution  $x^* \in \mathbb{R}^n$  with  $x_i^* \in [-1, 1]$ . Then, we set  $b = Ax^*$ . The starting point is  $x_0 = 0$ .

Given  $t_{\min}$ ,  $t_{\max}$  and some positive integer  $n_{\bar{\eta}}$ , we tested the inexact BB method with the inexactitude case (52) using the following values for  $\bar{\eta}$ :

$$\bar{\eta}_j = \exp\left(\log t_{\min} + \sqrt{(j-1)/(n_{\bar{\eta}}-1)} \log\left(\frac{t_{\max}}{t_{\min}}\right)\right), \quad j = 1, \dots, n_{\bar{\eta}}$$

Since  $\|\tilde{g}_{-1}\|$  is not available for the first iteration, we assumed that  $\|\tilde{g}_{-1}\| = 0$ , which means that the exact gradient  $g_0$  is used. The initial stepsize is set to be the Cauchy stepsize  $\alpha_0 = g_0^T g_0 / g_0^T A g_0$  as in Step 1 of Algorithm 2.1. For  $k \geq 1$ , to generate an approximation  $\tilde{g}_k$  of the exact gradient  $g_k$ , we first generate a random vector  $v_k$  and then set

$$\tilde{g}_k = g_k + \bar{\eta} \frac{\|\tilde{g}_{k-1}\|}{\|v_k\|} v_k$$

The above choice of  $\tilde{g}_k$  is such that the equality in relation (52) holds. Although condition (52) cannot guarantee the descent property of  $-\tilde{g}_k$ , Theorem 3.2 tells us that the inexact BB method still converges and the convergence rate is  $R$ -linear.

In our tests, we fix  $t_{\min} = 10^{-3}$ ,  $t_{\max} = 0.5$  and  $n_{\bar{\eta}} = 100$  and observe the influence of  $\kappa$ ,  $n$  and  $\lambda_1$ . For each value of  $\kappa$ ,  $n$  and  $\lambda_1$ , we do 100 tests and use the stopping condition

$$\|g_k\| \leq 10^{-6} \|g_0\| \quad (70)$$

where the exact gradient is used for the purpose of comparison. All tests for this part were done with MATLAB 6.5.0.

At first, to observe the influence of  $\kappa$ , we fix  $n = 100$  and  $\lambda_1 = 1$  and use the following three values for  $\kappa$ :  $10^2$ ,  $10^3$  and  $10^4$ . For  $j = 1, \dots, n_{\bar{\eta}}$ , we denote by  $\text{Iter}(j)$  the average iteration numbers

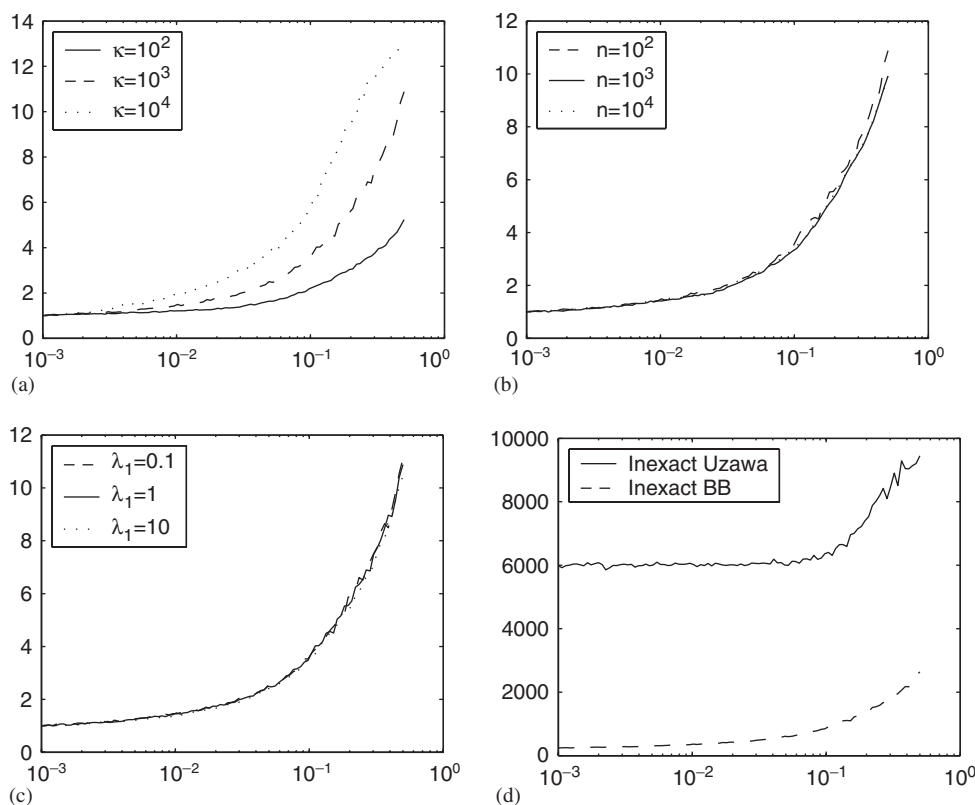


Figure 1. Note: the  $x$ -axis in the above figures corresponds to the value of  $\bar{\eta} = \eta_j$ . The  $y$ -axis in (a)–(c) corresponds to the relative value of  $\text{Iter}(j)$  with respect to  $\text{Iter}(1)$ ; the  $y$ -axis in (d) stands for the absolute value of  $\text{Iter}(j)$ .

of the 100 tests required for the inexactitude rule (52) with  $\bar{\eta} = \bar{\eta}_j$ . For the three different values of  $\kappa$ , Figure 1(a) plots the corresponding curves  $\text{Iter}(j)/\text{Iter}(1)$  vs  $\bar{\eta}_j$ . Secondly, we fix  $\kappa = 10^3$  and  $\lambda_1 = 1$  and vary the problem dimension  $n$  to be  $10^2$ ,  $10^3$  and  $10^4$ , respectively. Thirdly, we fix  $\kappa = 10^3$  and  $n = 10^2$  and vary the minimal eigenvalue  $\lambda_1$  to be 0.1, 1 and 10, respectively. See Figures 1(b) and (c) for the corresponding curves. Considering the log scale of the horizontal axis, we can see from the three figures that basically, as the inexactness parameter  $\bar{\eta}$  increases, the required average iteration number linearly increases and hence the  $R$ -linear factor of the inexact BB algorithm linearly decreases. The decrement of the  $R$ -linear factor is strongly affected by the condition number of the problem, but has little relation with the problem dimension or the minimal eigenvalue of the coefficient matrix  $A$ . In the numerical sense, this feature resembles that of the inexact Uzawa algorithm, whose  $Q$ -linear factor mainly depends on  $\kappa$  and  $\bar{\eta}$  (see Theorem 2.2 in [5]).

A comparison was also made between the inexact BB algorithm and the inexact Uzawa algorithm. In this case, we fix  $n = 100$ ,  $\kappa = 10^3$  and  $\lambda_1 = 1$ . Therefore, the only difference between the two algorithms is the choice of the stepsize  $\alpha_k$ . The inexact Uzawa algorithm uses (5), which means

that  $\alpha_k \equiv \frac{2}{101}$ , whereas the inexact BB algorithm decides the stepsize according to the information at the most recent two points except the initial stepsize is chosen to be the Cauchy stepsize. In Figure 1(d), we plot the curves of the required average iteration numbers  $\text{Iter}(j)$  vs  $\bar{\eta}_j$ . From the figure, we can see that the inexact BB algorithm is far more efficient than the inexact Uzawa algorithm. In the case when  $\bar{\eta} = \bar{\eta}_1 = 10^{-3}$ , to reach the stopping condition (70), the inexact BB algorithm and the inexact Uzawa algorithm require 242.6 and 5983.5 iterations on the average. From Figure 1(d), we can also see that the influence of the inexactness parameter  $\bar{\eta}$  is similar to the performance of the two algorithms.

In the second part of our numerical experiments, we test Algorithm 3.1 on saddle point problems arising from the finite-element discretization of Stokes equations.

#### Problem 4.2

This problem is related to the Stokes equations:

$$\begin{aligned} -\Delta u + \nabla p &= f & \text{in } \Omega = (0, 1) \times (0, 1) \\ -\text{div } u &= h & \text{in } \Omega \\ u &= 0 & \text{on } \partial\Omega \\ \int_{\Omega} p &= 0 \end{aligned} \tag{71}$$

We discretize (71) in the same way as that in [5]. More exactly, the discretization takes uniform triangular meshes on  $\Omega$ , use continuous piecewise linear velocities on a mesh of width  $d$  and use continuous piecewise linear pressures on a mesh of width  $2d$  (this discretization is called as  $P_1(h)P_1(2h)$  in [5]). We are then led to the saddle point problem (1) with  $n = 2(\frac{1}{d} - 1)^2$  and  $m = (\frac{2}{d} - 1)^2$ . For this problem,  $f$  is randomly generated with its elements in  $[-1, 1]$  and  $h$  is set to zero.

In this experiment, we vary  $\tau$  to see the influence of the inexactitude of the subproblem on the performance of Algorithm 3.1. For each value of  $\tau$ , we generate 20 random experiments and observe the average performance of Algorithm 3.1. Note that the purpose of this work mainly focuses on the inexact BB method used for computing  $p_k$ . For ease in coding, we simply use the (exact) BB method with no preconditioner to solve the subproblem for  $u_{k+1}$  (see Step 2 of Algorithm 3.1). It is certain that other methods possibly with some preconditioner can be used to solve the subproblem, for example, conjugate gradient methods or other gradient methods. Considering the special structure of  $A$ , some other effective iterative or direct methods can be also found. When the (exact) BB method is used for the subproblem, at least one step is computed. The maximal iteration number, INMAX say, for the subproblem is set to 100. If this number is exceeded, the inner solver exits with the point having the minimal residuals. Here, we note that some other values were also used for INMAX, for example, 250 and 400. However, the difference between the numerical results is not significant.

The tests in Problems 4.2 and 4.3 were done with MATLAB 7.1.0. We list the numerical results of Problem 4.2 in Table I, where ‘# OUT’ means the average number of outer iterations, ‘# IN’ stands for the average number of inner iterations per outer iteration and ‘time’ is the required average central processing unit (CPU) time in second. Listed in the last column is the number of outer iterations required by the inexact Uzawa algorithm with the same value of  $\tau$  for the problem, which can be found in [5] (only the case  $d = \frac{1}{32}$  is available).

Table I. Results of Algorithm 3.1 for Problem 4.2.

$d$	$\tau$	# OUT	# IN	time	Uzawa
1/16	1/64	89.0	78.9	0.995	—
1/16	1/16	99.4	68.4	0.993	—
1/16	1/4	105.3	56.5	0.911	—
1/16	1	146.4	43.7	1.113	—
1/32	1/64	194.5	99.1	9.832	(500)
1/32	1/16	185.3	98.3	9.245	426
1/32	1/4	185.2	95.4	8.915	427
1/32	1	194.7	86.6	8.839	431

Table I indicates that as  $\tau$  increases, the average number of outer iterations increases and the cost for inner iteration decreases. In this experiment, the suggested value for  $\tau$  is  $\frac{1}{4}$ . From the table, we also see that the inexact BB algorithm requires less than one half of the outer iterations by the inexact Uzawa algorithm. This means that if the same solver is used for the subproblem, the inexact BB algorithm, which does not need to estimate the minimal and maximal eigenvalues of  $A$ , is even faster.

#### Problem 4.3

This problem comes from a different discretization of the Stokes equations with non-zero diverse of velocity, namely,  $h \neq 0$ . Specifically, we discrete the velocity with continuous piecewise linear finite elements on uniform triangular mesh of width  $d$  as before. However, for the pressure function we use the elements suggested in [23, 24]. More exactly, a piecewise constant function on the adjacent blocks on uniform square mesh of width  $d$  is used:

$$\phi_{ij} = \begin{cases} 1, & x \in (x_{i-1}, x_i), \quad y \in (y_{i-1}, y_i) \\ -1, & x \in (x_i, x_{i+1}), \quad y \in (y_{i-1}, y_i) \\ -1, & x \in (x_{i-1}, x_i), \quad y \in (y_i, y_{i+1}) \\ 1, & x \in (x_i, x_{i+1}), \quad y \in (y_i, y_{i+1}) \end{cases}$$

The discretization leads to the saddle point problem (1) with  $n = 2(1/d - 1)^2$  and  $m = (1/d - 1)^2$ . To guarantee the positive definiteness of the matrix  $C + BA^{-1}B^T$ , we choose  $C = dI_m$  with  $I_m$  being the identity matrix in  $\mathbb{R}^{m \times m}$ . It is easy to see that  $C$  tends to zero as  $d$  tends to zero. For this problem, both  $f$  and  $g$  are randomly generated with their elements in  $[-1, 1]$ .

For this problem, we vary the value of  $\tau$  to be  $\frac{1}{16}$ ,  $\frac{1}{8}$ ,  $\frac{1}{4}$  and  $\frac{1}{2}$ . The maximum for the inner iterations, INMAX, is set to 10. The results of Algorithm 3.1 for Problem 4.3 are reported in Table II. From the table, we see that the choice of  $\tau = \frac{1}{8}$  is preferred, but the performance of Algorithm 3.1 is not sensitive to the parameter  $\tau$ . However, we observed that the performance of Algorithm 3.1 heavily depends on the choice of INMAX. In the case of INMAX = 10, we see from the column # IN of Table II that the inner solver often reaches INMAX iterations. We found that this is also the case if INMAX = 50. Take  $d = \frac{1}{32}$  and  $\tau = \frac{1}{8}$  as an example. If we set INMAX = 50, the average number of inner iterations is 47.8, whereas the algorithm still needs



Table II. Results of Algorithm 3.1 for Problem 4.3.

$d$	$\tau$	# OUT	# IN	Time
1/8	1/16	36.4	9.7	0.103
1/8	1/8	34.4	9.6	0.097
1/8	1/4	35.2	9.3	0.102
1/8	1/2	37.8	9.2	0.111
1/16	1/16	41.9	9.8	0.494
1/16	1/8	42.9	9.7	0.503
1/16	1/4	42.2	9.5	0.490
1/16	1/2	44.0	9.4	0.516
1/32	1/16	76.7	9.9	4.953
1/32	1/8	70.6	9.8	4.442
1/32	1/4	72.5	9.7	4.588
1/32	1/2	78.3	9.7	5.075
1/64	1/16	153.8	9.9	57.856
1/64	1/8	144.3	9.9	52.254
1/64	1/4	157.2	9.9	59.388
1/64	1/2	150.9	9.8	55.596

68.2 outer iterations on the average. The total time is 8.016, about double the one by choosing  $\text{INMAX} = 10$ .

The numerical experiments in Problems 4.2 and 4.3 suggest that the inexact BB algorithm is an efficient alternative to the inexact Uzawa algorithm. On the other hand, due to its non-monotonic feature, the BB algorithm might not be a good option for the inner solver. Further numerical experiments are still required to understand the behaviour of the BB algorithm.

## 5. CONCLUSIONS AND DISCUSSIONS

In this paper, we have analysed the inexact BB method with the inexactitude rules (15) and (52). The analysis with rule (15) could help us in understanding the (exact) BB method for unconstrained optimization, since the latter can be regarded as an inexact BB method for some quadratic function if  $x_k$  tends to  $x^*$ . Consequently, we are able to prove that the (exact) BB method is locally  $R$ -linearly convergent for twice continuously differentiable functions, a result stronger than the one in [22] for the cyclic BB method. Another interesting point with the BB method is that, in the previous analysis of the BB method, the non-monotonicity is introduced by the choice of stepsize  $\alpha_k$ . Since the inexactitude rules (15) allows the possibility that  $-\tilde{g}_k$  is an uphill search direction, Theorem 3.2 tells us that it would be also fine to introduce some suitable non-monotonicity in choosing search directions without affecting the  $R$ -linear convergence.

To adapt the inexact BB algorithm for solving saddle point problems, we also analysed the rule (52) carefully and provide  $R$ -linear convergence result in the case. These analyses are based on those with the inexact rule (15). However, our theoretical analyses cannot provide a good estimate for either the parameter  $\eta$  in (15) or the  $\bar{\eta}$  in (52), although our numerical experiments show that the latter one can be as large as the one in the inexact Uzawa algorithm. It still remains under study how to estimate the parameters  $\eta$  and  $\bar{\eta}$  theoretically. The good solution of this problem is related to the question how to establish theoretical evidence showing that the (exact) BB method is

faster than the Uzawa algorithm or the steepest descent method in the any dimension case. Some evidence in low dimensions has been established in [11, 14].

To solve the saddle point problems (1), References [9, 23, 25] introduce some preconditioner for the inexact Uzawa method:

$$\begin{aligned} u_{k+1} &= Q^{-1}(f - B^T p) \\ p_{k+1} &= p_k + \bar{Q}^{-1}(Cp_k - Bu_{k+1} + h) \end{aligned} \quad (72)$$

where  $Q$  and  $\bar{Q}$  are some approximation to  $A$  and  $\bar{A} = BA^{-1}B^T + C$  as before. Some extensions are also made to the case when  $A$ ,  $B$  and  $C$  are nonlinear operators. Since the BB method does not need to estimate any eigenvalue of the coefficient matrix and is far better than the Uzawa algorithm or the steepest descent method, it might be worthwhile to study the above issues with the inexact BB methods. As the referees commented, the use of preconditioning is indispensable in building fast and practical methods. Here, we should note that the preconditioning technique was used to the BB method first by Molina and Raydan [26]. Another important work related to this topic is Hernández-Ramos [27], which proposed preconditioned gradient (Uzawa) and conjugate gradient algorithms for saddle point problems. On the other hand, there have been many contenders of the BB stepsize (8), see [14, 19], for example. Our future work is then to establish an efficient inexact and preconditioning BB-like method for saddle point problems.

#### ACKNOWLEDGEMENTS

The authors are very grateful to Professor Marcos Raydan in Universidad Central de Venezuela and two anonymous referees for their useful comments and suggestions that improved the quality of this paper greatly.

#### REFERENCES

1. Benzi M, Golub GH, Liesen J. Numerical solution of saddle point problems. *Acta Numerica* 2005; **14**:1–137.
2. Elman HC, Silvester DJ, Wathen AJ. *Finite Elements and Fast Iterative Solvers: With Applications in Incompressible Fluid Dynamics*. Oxford University Press: Oxford, 2005.
3. Arrow K, Hurwicz L, Uzawa H. *Studies in Nonlinear Programming*. Stanford University Press: Stanford, CA, 1958.
4. Varga PS. *Matrix Iterative Analysis*. Prentice-Hall: Englewood Cliffs, NJ, 1962.
5. Elman HC, Golub GH. Inexact and preconditioned Uzawa algorithms for saddle point problems. *SIAM Journal on Numerical Analysis* 1994; **31**(6):1645–1661.
6. Dai YH, Yang XQ. A new gradient method with an optimal stepsize property. *Computational Optimization and Applications* 2006; **33**(1):73–88.
7. Cauchy A. Méthode générale pour la résolution des systèmes d'équations simultanées. *Comptes Rendus de l'Académie des Sciences, Paris* 1847; **25**:536–538.
8. Cao ZH. Fast Uzawa algorithm for generalized saddle point problems. *Applied Numerical Mathematics* 2003; **46**(2):157–171.
9. Cheng XL. On the nonlinear inexact Uzawa algorithm for saddle-point problems. *SIAM Journal on Numerical Analysis* 2000; **37**(6):1930–1934.
10. Hu QY, Zou J. Two new variants of nonlinear inexact Uzawa algorithms for saddle-point problems. *Numerische Mathematik* 2002; **93**(2):333–359.
11. Barzilai J, Borwein JM. Two-point step size gradient methods. *IMA Journal of Numerical Analysis* 1988; **8**:141–148.
12. Raydan M. On the Barzilai and Borwein choice of steplength for the gradient method. *IMA Journal of Numerical Analysis* 1993; **13**(3):321–326.

13. Dai YH, Liao LZ.  $R$ -linear convergence of the Barzilai and Borwein gradient method. *IMA Journal of Numerical Analysis* 2002; **22**(1):1–10.
14. Dai YH, Fletcher R. On the asymptotic behaviour of some new gradient methods. *Mathematical Programming, Series A* 2005; **103**(3):541–559.
15. Fletcher R. On the Barzilai–Borwein method. In *Optimization and Control with Applications*, Qi L, Teo K, Yang XQ (eds). Springer: Berlin, 2005; 235–256.
16. Raydan M. The Barzilai and Borwein gradient method for the large scale unconstrained minimization problem. *SIAM Journal on Optimization* 1997; **7**(1):26–33.
17. Glunt W, Hayden TL, Raydan M. Molecular conformations from distance matrices. *Journal of Computational Chemistry* 1993; **14**(1):114–120.
18. Birgin EG, Chambouleyron I, Martínez JM. Estimation of the optical constants and the thickness of thin films using unconstrained optimization. *Journal of Computational Physics* 1999; **151**(2):862–880.
19. Friedlander A, Martínez JM, Molina B, Raydan M. Gradient method with retards and generalizations. *SIAM Journal on Numerical Analysis* 1999; **36**(1):275–289.
20. Dai YH, Yuan JY, Yuan YX. Modified two-point stepsize gradient methods for unconstrained optimization. *Computational Optimization and Applications* 2002; **22**(1):103–109.
21. Liu WB, Dai YH. Minimization algorithms based on supervisor and searcher cooperation. *Journal of Optimization Theory and Applications* 2001; **111**(2):359–379.
22. Dai YH, Hager W, Schittkowski K, Zhang H. The cyclic Barzilai–Borwein method for unconstrained optimization. *IMA Journal of Numerical Analysis* 2006; **26**(3):604–627.
23. Bramble JH, Pasciak JE, Vassilev AT. Analysis of the inexact Uzawa algorithm for saddle point problems. *SIAM Journal on Numerical Analysis* 1997; **34**(3):1072–1092.
24. Bramble JH, Pasciak JE. A preconditioning technique for indefinite systems resulting from mixed approximations of elliptic problems. *Mathematics of Computation* 1988; **50**:1–18.
25. Bansch E, Morin P, Nochetto RH. An adaptive Uzawa FEM for the stokes problem: convergence without the Inf-Sup condition. *SIAM Journal on Numerical Analysis* 2002; **40**(4):1207–1229.
26. Molina B, Raydan M. Preconditioned Barzilai–Borwein method for the numerical solution of partial differential equations. *Numerical Algorithms* 1996; **13**(1):45–60.
27. Hernández-Ramos LM. Alternating oblique projections for coupled linear systems. *Numerical Algorithms* 2005; **38**(4):285–303.