

## ITERATIVE METHODS FOR THE FORWARD-BACKWARD HEAT EQUATION <sup>\*1)</sup>

Xiao-liang Cheng

(Department of Mathematics, Zhejiang University, Hangzhou 310028, China)

Jie Sun

(Department of Mathematics, Zhejiang University, Hangzhou 310028, China)

(Department of Mathematics, Zhejiang University of Finance and Economics,  
Hangzhou 310012, China)

### Abstract

In this paper we propose the finite difference method for the forward-backward heat equation. We use a coarse-mesh second-order central difference scheme at the middle line mesh points and derive the error estimate. Then we discuss the iterative method based on the domain decomposition for our scheme and derive the bounds for the rates of convergence. Finally we present some numerical experiments to support our analysis.

*Mathematics subject classification:* 65N22, 65M06, 35K05, 65N55.

*Key words:* Forward-backward heat equation, Finite difference method, Iterative method, Coarse mesh.

### 1. Introduction

In this paper, we consider the following boundary value problem of a forward-backward parabolic equation:

$$\begin{cases} a(x)u_t - u_{xx} = f(x, t), & (x, t) \in \Omega = (-1, 1) \times (0, 1), \\ u(x, 0) = 0, & 0 \leq x \leq 1, \\ u(x, 1) = 0, & -1 \leq x \leq 0, \\ u(1, t) = 0, u(-1, t) = 0, & 0 < t < 1, \end{cases} \quad (1.1)$$

where  $a(x) > 0$  for  $x > 0$ ,  $a(x) < 0$  for  $x < 0$  and  $a(0) = 0$ . For example,  $a(x) = x$  or  $a(x) = x^m$  with  $m$  the odd integer. The problem (1.1) arises in a variety of applications such as randomly accelerated particle problem and fluid flow near a boundary where separation occurs, see [1, 2] for the details. So far there are several numerical approach to this problem, for example, the finite difference method[1], least square method[5] and Galerkin finite element method[3, 4, 7, 8].

The purpose of this paper is to present a finite difference scheme to equation (1.1). Unlike the standard way in [1], we use a coarse mesh second-order central difference scheme at the mesh points lie on the middle line  $x = 0$ ,  $0 < t < 1$ . We prove the error estimates  $O(\tau + h^2 + H^3)$  with time mesh size  $\tau$  and space mesh size  $h$  and coarse mesh size  $H$ . Then we discuss the iterative method based on the domain decomposition method for our scheme and obtain bounds of the convergent rate with  $1 - H$ , which is better than that  $1 - h$  in [1]. In the last section we present some numerical results to support our analysis.

### 2. The Difference Scheme

We first specify the grids. Let  $h = 1/M$  and  $x_i = ih$  for  $i = 0, \pm 1, \pm 2, \dots, \pm M$ . Let  $\tau = 1/N$  and  $t_j = j\tau$  for  $j = 0, 1, \dots, N$ .

---

\* Received February 3, 2004; final revised July 19, 2004.

<sup>1)</sup> The project was supported by National Science Foundation(Grant No. 10471129).

We use the backward and forward difference scheme on domain  $x > 0$  and  $x < 0$  respectively and second order central difference scheme on the line  $x = 0$  with coarse mesh  $H = m_0 h$  for some given positive integer  $m_0$ . Denote  $z_i^j$  ( $-M + 1 \leq i \leq M - 1, 1 \leq j \leq N - 1$ ) is the approximation solution for the exact solution at point  $(ih, j\tau)$ . Then

$$\begin{cases} a_i \frac{z_i^{j+1} - z_i^j}{\tau} - \frac{z_{i-1}^{j+1} - 2z_i^{j+1} + z_{i+1}^{j+1}}{h^2} = f_i^{j+1}, & 1 \leq i \leq M - 1, 0 \leq j \leq N - 2, \\ z_i^0 = 0, & 1 \leq i \leq M - 1, \\ z_M^j = 0, & 1 \leq j \leq N - 1; \end{cases} \quad (2.1)$$

$$\begin{cases} a_i \frac{z_i^{j+1} - z_i^j}{\tau} - \frac{z_{i-1}^j - 2z_i^j + z_{i+1}^j}{h^2} = f_i^j, & -M + 1 \leq i \leq -1, 1 \leq j \leq N - 1, \\ z_i^N = 0, & -M + 1 \leq i \leq -1, \\ z_{-M}^j = 0, & 1 \leq j \leq N - 1; \end{cases} \quad (2.2)$$

and

$$-\frac{z_{m_0}^j - 2z_0^j + z_{-m_0}^j}{H^2} = f_0^j, \quad 1 \leq j \leq N - 1. \quad (2.3)$$

Here  $a_i = a(ih)$  and  $f_i^j = f(ih, j\tau)$ . For  $m_0 = 1$ , it is the same method proposed in [1].

It is convenient to introduce the set of all mesh points by  $\mathcal{N}$ ,

$$\mathcal{N} = \{(ih, j\tau) \mid -M + 1 \leq i \leq M - 1, 1 \leq j \leq N - 1\}.$$

Further we split  $\mathcal{N} = \mathcal{N}_w \cup \mathcal{N}_v \cup \mathcal{N}_\psi$  into three disjoint subsets as follows,

$$\mathcal{N}_w = \{(ih, j\tau) \mid 1 \leq i \leq M - 1, 1 \leq j \leq N - 1\},$$

$$\mathcal{N}_v = \{(ih, j\tau) \mid -M + 1 \leq i \leq -1, 1 \leq j \leq N - 1\},$$

$$\mathcal{N}_\psi = \{(0, j\tau) \mid 1 \leq j \leq N - 1\}.$$

We write the linear system (2.1)-(2.3) in the matrix form  $\mathcal{P}Z = F$  with

$$\mathcal{P} = \begin{pmatrix} A_{vv} & 0 & A_{v\psi} \\ 0 & A_{ww} & A_{w\psi} \\ A_{v\psi} & A_{w\psi} & A_{\psi\psi} \end{pmatrix} \quad (2.4)$$

The vector  $Z = (Z_v, Z_w, Z_\psi)^T$  with

$$Z_v = (z_{-M+1}^1, z_{-M+1}^2, \dots, z_{-1}^1, \dots, z_{-1}^{N-1})^T,$$

$$Z_w = (z_1^1, z_1^2, \dots, z_{M-1}^1, \dots, z_{M-1}^{N-1})^T,$$

and  $Z_\psi = (z_0^1, z_0^2, \dots, z_0^{N-1})^T$ . And  $F$  the vector defined on the mesh points  $\mathcal{N}$  of the function  $f(x, t)$ .

Let  $u$  be the exact solution of problem (1.1). Denote the error

$$E_i^j = u(ih, j\tau) - z_i^j, \quad (ih, j\tau) \in \mathcal{N}.$$

We use the maximum norm

$$\|E\|_{\mathcal{N}} = \max_{(ih, j\tau) \in \mathcal{N}} |E_i^j|.$$

Now we will prove the following error estimates.

**Theorem 2.1.** *Suppose that  $\frac{1}{2}|\partial^2 u / \partial t^2|$  and  $\frac{1}{12}|\partial^4 u / \partial x^4|$  are bounded by constant  $C_0$  on  $\bar{\Omega}$ , the closure of  $\Omega$ . Then*

$$\|E\|_{\mathcal{N}} \leq \frac{1}{2}C_0(\tau + h^2 + H^3). \quad (2.5)$$

*Proof.* It is easy to see that from Taylor's series, we have the truncation error in (2.1)-(2.3)

$$R_i^j = (\mathcal{P}E)_i^j = \begin{cases} K_i^j(\tau + h^2), & (ih, j\tau) \in \mathcal{N}_v \cup \mathcal{N}_w, \\ K_0^j H^2, & (0, j\tau) \in \mathcal{N}_\psi, \end{cases} \quad (2.6)$$

where  $|K_i^j| \leq C_0$  and  $|K_0^j| \leq C_0$ . As in [1], it is not difficult to see that the matrix  $\mathcal{P}$  has positive diagonal entries, nonpositive offdiagonal entries, and is irreducibly diagonally dominant. Hence  $\mathcal{P}$  is an M matrix[10]. That is all the entries of the inverse of matrix  $\mathcal{P}$  are nonnegative. It is equivalent to the maximum principle in the finite difference method and it is useful to derive the error estimates.

Borrow the trick of [9], we construct

$$\alpha_i^j = 1 - \frac{i^2}{M^2}, \quad \beta_i^j = H(1 - \frac{|i|}{M}) \quad (2.7)$$

for  $-M + 1 \leq i \leq M - 1, 1 \leq j \leq N - 1$ . Obviously  $|\alpha_i^j| \leq 1$  and  $|\beta_i^j| \leq H$ . Then we can verify that

$$(\mathcal{P}\alpha)_i^j \geq 2, \quad (2.8)$$

and

$$(\mathcal{P}\beta)_i^j \geq 0, \quad i \neq 0 \quad (2.9)$$

$$(\mathcal{P}\beta)_0^j = -\frac{H(1 - \frac{m_0}{M}) - 2H + H(1 - \frac{m_0}{M})}{H^2} = 2. \quad (2.10)$$

Putting

$$\xi_i^j = \frac{1}{2}C_0(\tau + h^2)\alpha_i^j + \frac{1}{2}C_0H^2\beta_i^j, \quad (2.11)$$

then we can prove

$$(\mathcal{P}(\xi \pm E))_i^j \geq 0.$$

Thus by the M matrix property of  $\mathcal{P}$ , we have

$$|E_i^j| \leq \xi_i^j \leq \frac{1}{2}C_0(\tau + h^2) + \frac{1}{2}C_0H^3, \quad -M + 1 \leq i \leq M - 1, 1 \leq j \leq N - 1.$$

It completes the proof.

This result is not surprising for the presence of the  $H^3$  term. In [9], they first derive the similar result for the numerical approximation of heat equation by domain decomposition method—an explicit forward difference formula on the coarse mesh at the interface mesh points, and implicit backward difference formula on the fine mesh in sub-domains. So they weaken the condition  $\tau \leq \frac{1}{2}h^2$  to  $\tau \leq \frac{1}{2}H^2$ . Here we use the coarse mesh to improve the convergent rate of the iterative algorithm discussed in next section.

### 3. Iterative Method

In this section, we discuss the iterative method based on the domain decomposition method to our scheme (2.1)-(2.3). To express the idea, we use the matrix form (2.4). We first give an initial guessing value at grid points on the line  $x = 0$  with  $Z_\psi^0 = (z_0^{1,0}, z_0^{2,0}, \dots, z_0^{N-1,0})^T$ , then we solve the problem separately in domain  $x > 0$  and  $x < 0$  by

$$A_{vv}Z_v^1 = F_v - A_{v\psi}Z_\psi^0, \quad (3.1)$$

$$A_{ww}Z_w^1 = F_w - A_{w\psi}Z_\psi^0, \quad (3.2)$$

and update  $Z_\psi^1$  by

$$A_{\psi\psi}Z_\psi^1 = F_\psi - A_{v\psi}Z_v^1 - A_{w\psi}Z_w^1. \quad (3.3)$$

Repeating above procedure we have the following iterative algorithm. For initial guess value  $\phi^{j,0}$ , ( $1 \leq j \leq N - 1$ ), let  $k = 1, 2, \dots$ , we solve the two subsystems as follows:

$$\left\{ \begin{aligned} a_i \frac{z_i^{j+1,k} - z_i^{j,k}}{\tau} - \frac{z_{i-1}^{j+1,k} - 2z_i^{j+1,k} + z_{i+1}^{j+1,k}}{h^2} &= f_i^{j+1}, \quad 1 \leq i \leq M - 1, \quad 0 \leq j \leq N - 2, \\ z_0^{j,k} &= \phi^{j,k-1}, \quad 1 \leq j \leq N - 1, \\ z_i^{0,k} &= 0, \quad 1 \leq i \leq M - 1, \\ z_M^{j,k} &= 0, \quad 1 \leq j \leq N - 1; \end{aligned} \right. \tag{3.4}$$

$$\left\{ \begin{aligned} a_i \frac{z_i^{j+1,k} - z_i^{j,k}}{\tau} - \frac{z_{i-1}^{j,k} - 2z_i^{j,k} + z_{i+1}^{j,k}}{h^2} &= f_i^j, \quad -M + 1 \leq i \leq -1, \quad 1 \leq j \leq N - 1, \\ z_0^{j,k} &= \phi^{j,k-1}, \quad 1 \leq j \leq N - 1, \\ z_i^{N,k} &= 0, \quad -M + 1 \leq i \leq -1, \\ z_{-M}^{j,k} &= 0, \quad 1 \leq j \leq N - 1; \end{aligned} \right. \tag{3.5}$$

and

$$\phi^{j,k} = \frac{1}{2} \left( z_{m_0}^{j,k} + z_{-m_0}^{j,k} + H^2 f_0^j \right), \quad 1 \leq j \leq N - 1. \tag{3.6}$$

The linear system (3.4) is easy to solve. We only need to solve a tridiagonal linear system for  $(z_1^{1,k}, z_2^{1,k}, \dots, z_{M-1}^{1,k})$ , then for  $(z_1^{2,k}, z_2^{2,k}, \dots, z_{M-1}^{2,k})$  and so on. Thus we need solve  $N - 1$  tridiagonal linear systems for (3.4). Equation (3.5) is the same as (3.4) but from  $z^{N-1,k}$  to  $z^{1,k}$ .

**Remark 3.1.** We can also use the relaxation method [1], replacing (3.3) by

$$\left\{ \begin{aligned} A_{\psi\psi} \tilde{Z}_\psi^1 &= F_\psi - A_{v\psi} Z_v^1 - A_{w\psi} Z_w^1, \\ Z_\psi^1 &= \omega Z_\psi^0 + (1 - \omega) \tilde{Z}_\psi^1, \quad 0 < \omega \leq 1. \end{aligned} \right. \tag{3.7}$$

In [11], they proposed an algorithm which is similar to above algorithm by the non-overlap domain decomposition method. They derived some convergent rate in some different norm.

We next derive the convergent rate for the iterative algorithm (3.4)-(3.6).

**Theorem 3.1.** Let  $\phi^{j,k}$  ( $1 \leq j \leq N - 1, k = 0, 1, 2, \dots$ ) be the solutions of equations (3.4)-(3.6),  $z_0^j$  ( $1 \leq j \leq N - 1$ ) be the solution of equations (2.1)-(2.3). Then

$$\max_{1 \leq j \leq N-1} (|z_0^j - \phi^{j,k+1}|) \leq (1 - H) \max_{1 \leq j \leq N-1} (|z_0^j - \phi^{j,k}|). \tag{3.8}$$

So  $\phi^{j,k}$  converge to  $z_0^j$  with the rate  $1 - H$  as  $k \rightarrow \infty$ .

*Proof.* Let the iteration error

$$\varepsilon_i^{j,k} = z_i^j - z_i^{j,k}, \quad -M \leq i \leq M, \quad i \neq 0, \quad 1 \leq j \leq N - 1.$$

Then the error  $\varepsilon_i^{j,k}$  satisfies:

$$\left\{ \begin{aligned} a_i \frac{\varepsilon_i^{j+1,k} - \varepsilon_i^{j,k}}{\tau} - \frac{\varepsilon_{i-1}^{j+1,k} - 2\varepsilon_i^{j+1,k} + \varepsilon_{i+1}^{j+1,k}}{h^2} &= 0, \quad 1 \leq i \leq M - 1, \quad 0 \leq j \leq N - 2, \\ \varepsilon_0^{j,k} &= z_0^j - \phi^{j,k-1}, \quad 1 \leq j \leq N - 1 \\ \varepsilon_i^{0,k} &= 0, \quad 1 \leq i \leq M - 1, \\ \varepsilon_M^{j,k} &= 0, \quad 1 \leq j \leq N - 1; \end{aligned} \right. \tag{3.9}$$

$$\left\{ \begin{array}{l} a_i \frac{\varepsilon_i^{j+1,k} - \varepsilon_i^{j,k}}{\tau} - \frac{\varepsilon_{i-1}^{j,k} - 2\varepsilon_i^{j,k} + \varepsilon_{i+1}^{j,k}}{h^2} = 0, \quad -M+1 \leq i \leq -1, \quad 1 \leq j \leq N-1, \\ \varepsilon_0^{j,k} = z_0^j - \phi^{j,k-1}, \quad 1 \leq j \leq N-1 \\ \varepsilon_i^{N,k} = 0, \quad -M+1 \leq i \leq -1, \\ \varepsilon_{-M}^{j,k} = 0, \quad 1 \leq j \leq N-1; \end{array} \right. \quad (3.10)$$

and

$$z_0^j - \phi^{j,k} = \frac{1}{2}(\varepsilon_{m_0}^{j,k} + \varepsilon_{-m_0}^{j,k}), \quad 1 \leq j \leq N-1. \quad (3.11)$$

Denote

$$S_k = \max_{1 \leq j \leq N-1} (|z_0^j - \phi^{j,k}|).$$

For the error linear system (3.9), we again write (3.9) in the matrix form  $\mathcal{G}\Upsilon^k = R$ . Here

$$\Upsilon^k = (\varepsilon_1^{1,k}, \varepsilon_1^{2,k}, \dots, \varepsilon_{M-1}^{1,k}, \dots, \varepsilon_{M-1}^{N-1,k})^T$$

and

$$R = (\frac{z_0^1 - \phi^{1,k-1}}{h^2}, \dots, \frac{z_0^{N-1} - \phi^{N-1,k-1}}{h^2}, 0, 0, \dots, 0)^T.$$

Similar to the proof of theorem 2.1, the matrix  $\mathcal{G}$  also has positive diagonal entries, non-positive off-diagonal entries, and is irreducibly diagonally dominant. So  $\mathcal{G}$  is M matrix. We construct

$$\eta_i^{j,k} = S_{k-1}(1 - |i|/M), \quad 1 \leq i \leq M-1, \quad 1 \leq j \leq N-1.$$

Thus we can obtain

$$(\mathcal{G}(\eta \pm \varepsilon))_i^{j,k} \geq 0, \quad 1 \leq i \leq M-1, \quad 1 \leq j \leq N-1.$$

So we have the inequality

$$|\varepsilon_i^{j,k}| \leq \eta_i^{j,k}, \quad 1 \leq i \leq M-1, \quad 1 \leq j \leq N-1. \quad (3.12)$$

As the same way we obtain

$$|\varepsilon_i^{j,k}| \leq \eta_i^{j,k}, \quad -M+1 \leq i \leq -1, \quad 1 \leq j \leq N-1. \quad (3.13)$$

Finally from the relation (3.11), we get for all  $1 \leq j \leq N-1$

$$\begin{aligned} |z_0^j - \phi^{j,k}| &\leq \frac{1}{2}(|\varepsilon_{m_0}^{j,k}| + |\varepsilon_{-m_0}^{j,k}|) \\ &\leq \frac{1}{2}(S_{k-1}(1 - |m_0|/M) + S_{k-1}(1 - |-m_0|/M)) \\ &\leq (1 - H)S_{k-1}. \end{aligned} \quad (3.14)$$

It is the result (3.8) and it completes the proof.

**Remark 3.2.** From theorem 2.1, we would expect to choose  $\tau \approx h^2 \approx H^3$ . If such choices are made,  $H^3 = h^2$  for example, then we will see the convergent rate of iterative method is  $1 - h^{2/3}$ . It is better than  $1 - h$  in [1].

**Remark 3.3.** We can also use different mesh size in  $N_v$  and  $N_w$  parts as long as  $H = m_1 h_1 = m_2 h_2$  for some integers  $m_1$  and  $m_2$ . We refer to the paper of [12] for parabolic problems.

### 4. Numerical Experiments

In this section we present some numerical results for the convergent rate of the iterative method (3.4)-(3.6) with  $a(x) = x$ . From the equation (3.9)-(3.11), we test the limit of

$$r = \lim_{k \rightarrow \infty} \frac{\max_{1 \leq j \leq N-1} (|z_0^j - \phi^{j,k+1}|)}{\max_{1 \leq j \leq N-1} (|z_0^j - \phi^{j,k}|)}$$

for the random initial error  $z_j^0 - \phi^{j,0}$ . See Table 1 and Table 2. We only report the case for fixed time mesh size  $N = 10$ . The numerical result is the same as the case  $m_0 = 1$  in [1] and is confirmed by our analysis in our paper  $r \approx 1 - H$ .

Table 1: The convergent rate of iterative method for N=10

$r$	$m_0 = 1$	$m_0 = 2$	$m_0 = 3$	$m_0 = 4$	$m_0 = 5$
$M = 10$	0.8940	0.7886	0.6841	0.5810	0.4797
$M = 20$	0.9470	0.8941	0.8413	0.7887	0.7363
$M = 30$	0.9647	0.9294	0.8941	0.8589	0.8237
$M = 40$	0.9735	0.9470	0.9205	0.8941	0.8677
$M = 50$	0.9788	0.9576	0.9364	0.9153	0.8941
$M = 60$	0.9823	0.9645	0.9470	0.9294	0.9117

Table 2: The convergent rate of iterative method for H=1/10

	$m_0 = 1, M = 10$	$m_0 = 5, M = 50$	$m_0 = 15, M = 150$	$m_0 = 20, M = 200$
$r$	0.8940	0.8941	0.8941	0.8941

## References

- [1] V. Vanaja and R.B. Kellogg, Iterative methods for a forward-backward heat equation, *SIAM J. Numer. Anal.*, **27** (1990), 622-635.
- [2] A.K. Aziz, D.A. French, S. Jensen and R.B. Kellogg, Origins, analysis, numerical analysis, and numerical approximation of a forward-backward parabolic problem, *Math. Model. Numer. Anal.*, **33** (1999), 895-922.
- [3] H. Lu and J. Maubach, A finite element method and variable transformations for a forward-backward heat equation, *Numer. Math.*, **81** (1998), 249-272.
- [4] A.K. Aziz and J.L. Liu, A Galerkin method for the forward-backward heat equation, *Math. Comput.*, **56** (1991), 35-44.
- [5] A.K. Aziz and J.L. Liu, A weighted least squares method for the forward-backward heat equation, *SIAM J. Numer. Anal.*, **28** (1991), 156-167.
- [6] W.-Z. Dai and Raja Nassari, An unconditionally stable hybrid FE-FD scheme for solving a 3-D heat transport equation in a cylindrical thin film with sub-microscale thickness, *J. Comput. Math.*, **21** (2003), 555-568.
- [7] D.A. French, Discontinuous Galerkin finite element methods for a forward-backward heat equation, *Appl. Numer. Math.*, **28** (1998), 37-44.
- [8] D.A. French, Continuous Galerkin finite element methods for a forward-backward heat equation, *Numer. Methods Partial Differ. Equat.*, **15** (1999), 257-265.
- [9] C.N. Dawson, Q. Du and T.F. Dupont, A finite difference domain decomposition algorithm for numerical solution of the heat equation, *Math. Comput.*, **57** (1991), 63-71.
- [10] R.S. Varga, *Matrix iterative analysis*, Prentice-Hall, Englewood Cliffs, NJ, 1962.
- [11] H. Han and D. Yin, A non-overlap domain decomposition method for the forward-backward heat equation, *J. Comput. Appl. Math.*, **159** (2003), 35-44.
- [12] Q. Du, M. Mu and Z.N. Wu, Efficient parallel algorithms for parabolic problems, *SIAM J. Numer. Anal.*, **39** (2001), 1469-1487.