# On iterative solving of linear algebraic equations in $p$-$h$-versions of finite element methods

V.P. Il'in   and   S.S. Radionov

The goal of this paper is investigation of the efficiency of iterative preconditioned cojugate gradient methods for solving linear systems of equations which arise in $p$-$h$-version of finite element methods. The results of numerical experiments are presented for the model boundary value problem with different values of $h$, $p$ and iterative parameters. Some conclusions on the comparative costs of algorithms are made on the base of analysis of matrix structure of linear systems.

## 1.   Introduction

We consider an efficiency of modern iterative algorithms in application for sparse linear systems of equations in $p$-$h$-versions of finite element methods [1]. The aim is to compare the numerical costs for different values of $p$, $h$ and some variants of iterative methods for different input data. For simplicity the Poisson equation

$$-\frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} = f(x, y) \tag{1.1}$$

is solved in a square computational domain with the Dirichlet or Neumann conditions on the different parts of boundary.

In Section 2 the structures of local and global stiffness matrices are described for the square elements. Section 3 includes the formulas of symmetric successive over relaxation method in the efficient Eisenstat implementation with conjugate gradient acceleration, see [2, 3, 4]. The discussion on the estimates of computer costs is presented in comparison with the direct Cholesky factorization algorithm. In the last section we present and discuss the results of numerical experiments to analyse separately the influence of preconditioning, iterative parameters, neccessary exactness and characteristics of boundary value problems.

## 2.   The matrices of $p$-$h$-finite element methods

For standard square element in a local coordinates

$$R = \{-1 \le x \le 1, \quad -1 \le y \le 1\}$$

there are used three types of basic functions with the total number $N_t = (p+1)^2$:

4 nodal functions

$$N_{k,l} = \frac{1+kx}{2} \cdot \frac{1+ly}{2}, \quad k,l = \pm 1; \tag{2.1}$$

$4(p-1)$ side functions

$$S_{k,l}^{(m)} = Q_m(x \cdot |l| + y \cdot |k|) \frac{1+kx}{2} \cdot \frac{1+ly}{2}, \quad |k| + |l| = 1,$$

$$Q_m(z) = \sqrt{\frac{2m+1}{2}} \int_{-1}^{z} P_m(z')dz',$$

$$P_m(z) = \frac{1}{2^m \cdot m!} \frac{d^m}{dz^m}(z^2-1)^m,$$

$$m = 1, 2, \ldots, p-1,$$

$(p-1)^2$ internal functions

$$I_{k,l} = Q_k(x) \cdot Q_l(y), \quad k,l = 1, \ldots, p-1. \tag{2.3}$$

We unify the basic functions into set $\varphi = \{\varphi_q(x,y)\}$ with the ordering

$$\varphi = \{N, S, I\},$$

where $N$, $S$ and $I$ mean the subsets of nodal, side and internal functions. In more detail

$$N = \{N_{1,1}, N_{-1,1}, N_{1,-1}, N_{-1,-1}\},$$

$$S = \{S_{1,0}^{(1)}, \ldots, S_{1,0}^{(p-1)}, S_{-1,0}^{(1)}, \ldots, S_{-1,0}^{(p-1)},$$

$$S_{0,1}^{(1)}, \ldots, S_{0,1}^{(p-1)}, S_{0,-1}^{(1)}, \ldots, S_{0,-1}^{(p-1)}\},$$

$$I = \{I_{1,1}, \ldots, I_{1,p-1}, I_{2,1}, \ldots, I_{2,p-1}, \ldots, I_{p-1,1}, \ldots, I_{p-1,p-1}\}.$$

It is natural to connect the different basic functions with geometric object of element. We will consider the nodal functions corresponded to

the mesh points, the groups $S_{1,0}^{(k)}$, $S_{-1,0}^{(k)}$, $S_{0,1}^{(k)}$, $S_{0,-1}^{(k)}$, $k = 1, \ldots, p - 1$, of side function corresponded to the right, left, upper and low element sides, and the groups $I_{k,q}$, $q = 1, \ldots, p - 1$; $I_{q,l}, q = 1, \ldots, p - 1$ corresponded to $k$-th vertical or $l$-th horizontal line inside the element.

Local stiffnes matrix is defined as

$$a = \{a_{q,q'}\}, \quad q, q' = 1, 2, \ldots, (p+1)^2.$$

Here $a_{q,q'}$ is scalar product

$$a_{q,q'} = (\varphi_q, \varphi_{q'}) = \int\limits_{R} \left( \frac{\partial \varphi_q}{\partial x} \cdot \frac{\partial \varphi_{q'}}{\partial x} + \frac{\partial \varphi_q}{\partial y} \cdot \frac{\partial \varphi_{q'}}{\partial y} \right) dx dy. \qquad (2.4)$$

The block structure of matrix $A$ is of the following type:

$$A = \begin{bmatrix} NN & NS & NI \\ SN & SS & SI \\ IN & IS & II \end{bmatrix}.$$

Here $NN$, $SS$ and $II$ denote the square submatrices of the orders 4, $4(p-1)$ and $(p-1)^2$ which correspond to scalar products of nodal, side and internal basic functions. The blocks $NS = (SN)^t$, $SI = (IS)^t$ and $NI = (IN)^t = 0$ are rectangular submatrices.

In particular, the matrix

$$NN = \begin{bmatrix} \frac{2}{3} & -\frac{1}{6} & -\frac{1}{6} & -\frac{1}{3} \\ -\frac{1}{6} & \frac{2}{3} & -\frac{1}{3} & -\frac{1}{6} \\ -\frac{1}{6} & -\frac{1}{3} & \frac{2}{3} & -\frac{1}{6} \\ -\frac{1}{3} & -\frac{1}{6} & -\frac{1}{6} & \frac{2}{3} \end{bmatrix}$$

has 16 nonzero entries. The matrix $SS$ is block diagonal one with the block of order 2:

$$SS = \begin{bmatrix} S_1 & S_2 & 0 & 0 \\ S_2 & S_1 & 0 & 0 \\ 0 & 0 & S_1 & S_2 \\ 0 & 0 & S_2 & S_1 \end{bmatrix}.$$

One diagonal block corresponds to the vertical mesh intervals and the other one to horisontal. The subblocks $S_1$ correspond to connections between the basic functions from the common sides, and $S_2$ – from the opposite.

These subblocks have the same tridiagonal structures (between each pair of nonzero diagonals there is one zero diagonal):

$$
S_k = \begin{bmatrix}
+ & 0 & - & \ddots & \ddots & & \\
0 & + & \ddots & - & \ddots & & \\
- & \ddots & + & \ddots & - & & \\
\ddots & - & \ddots & + & 0 & \\
\ddots & \ddots & - & 0 & + &
\end{bmatrix}.
$$

So, the matrix $SS$ has $8(3p - 7)$ nonzero elements. Of course, the above consideration is valid for $p > 3$. If $p = 2$, the matries $S_k$ are scalar values in fact, and for $p = 3$, $S_k$ are diagonal matries.

The matrix (of block order $p - 1$)

$$
II = \begin{bmatrix}
B & 0 & -C & 0 & . & . & . \\
0 & B & 0 & \ddots & \ddots & 0 & . \\
-C & 0 & \ddots & \ddots & \ddots & . & . \\
0 & \ddots & \ddots & \ddots & \ddots & . & 0 \\
. & \ddots & \ddots & \ddots & \ddots & 0 & -C \\
. & 0 & \ddots & \ddots & 0 & B & 0 \\
. & . & . & 0 & -C & 0 & B
\end{bmatrix}
$$

includes diagonal positive definite submatrix $C$, diagonal square blocks of order $p - 1$ with three nonzero diagonals (as $S_k$)

$$
B = \begin{bmatrix}
+ & 0 & - & & & \\
0 & \ddots & \ddots & \ddots & 0 & \\
- & \ddots & \ddots & \ddots & \ddots & \\
& \ddots & \ddots & \ddots & \ddots & - \\
& 0 & \ddots & \ddots & \ddots & 0 \\
& & & - & 0 & +
\end{bmatrix}
$$

and the total nonzero entries 1 for $p = 2$ and $(p - 1)(5p - 13)$ for $p \geq 3$. It is useful to remark that if side and internal functions at each element side or "inner line" will be numbered in "red-black" ordering (firstly – all odd functions and secondly – even), then $S_1$, $S_2$ and $B$ become block diagonal

matrices of block order 2, and each subblocks will be "usual" tridiagonal matrix. Nondiagonal block $NS$ has the form

$$NS = [NS_1, NS_2, NS_3, NS_4]$$

$$NS_1 = \begin{bmatrix} - & - & 0 & \cdots & 0 \\ + & + & 0 & \cdots & 0 \\ - & + & 0 & \cdots & 0 \\ + & - & 0 & \cdots & 0 \end{bmatrix}, \quad NS_2 = \begin{bmatrix} + & + & 0 & \cdots & 0 \\ - & - & 0 & \cdots & 0 \\ + & - & 0 & \cdots & 0 \\ - & + & 0 & \cdots & 0 \end{bmatrix},$$

$$NS_3 = \begin{bmatrix} - & - & 0 & \cdots & 0 \\ - & + & 0 & \cdots & 0 \\ + & + & 0 & \cdots & 0 \\ + & - & 0 & \cdots & 0 \end{bmatrix}, \quad NS_4 = \begin{bmatrix} + & + & 0 & \cdots & 0 \\ + & - & 0 & \cdots & 0 \\ - & - & 0 & \cdots & 0 \\ - & + & 0 & \cdots & 0 \end{bmatrix}.$$

Here the numbers of total and nonzero entries are $16(p-1)$ and $(32)$. The blocks $NS_1$, $NS_2$, $NS_3$, $NS_4$ have $p-1$ columns and are corresponded to the sets of basic functions $S_{1,0}^m$, $S_{-1,0}^m$, $S_{0,1}^m$, $S_{0,-1}^m$. Above "-" and "+" denote the signs of matrix nonzero entries.

The matrix $SI$ is rectangular one and has $4 \times (p-1)$ block structure

$$SI = \begin{bmatrix} S_{1,0}^{(1)} & S_{1,0}^{(2)} & \cdots & S_{1,0}^{(p-1)} \\ S_{-1,0}^{(1)} & S_{-1,0}^{(2)} & \cdots & S_{-1,0}^{(p-1)} \\ S_{0,1}^{(1)} & S_{0,1}^{(2)} & \cdots & S_{0,1}^{(p-1)} \\ S_{0,-1}^{(1)} & S_{0,-1}^{(2)} & \cdots & S_{0,-1}^{(p-1)} \end{bmatrix}$$

with the square blocks of the equal same order $p-1$. The blocks $S_{1,0}^{(m)}$ and $S_{-1,0}^{(m)}$ are diagonal matrices for $m = 1, 2$ and zero – for others $m$. The blocks $S_{0,1}^{(m)}$ and $S_{0,1}^{(m)}$ have nonzero entries in the two first columns at the row with number $m$ only. So, the matrix $SI$ has the total $4(p-1)^3$ entries and $8(p-1)$ nonzero elements. In more detail, the entries of matrix block $SI$ are defined by the scalar products of basic functions:

$$\left(S_{1,0}^{(m)}, I_{k,l}\right) = \begin{cases} -1/\sqrt{6}, & \text{for } l = 1, \ m = k, \\ -1/3\sqrt{10}, & \text{for } l = 2, \ m = k, \\ 0, & \text{otherwise}, \end{cases}$$

$$\left(S_{-1,0}^{(m)}, I_{k,l}\right) = \begin{cases} -1/\sqrt{6}, & \text{for } l = 1, \ m = k, \\ -1/3\sqrt{10}, & \text{for } l = 2, \ m = k. \\ 0, & \text{otherwise}, \end{cases}$$

$$\left(S_{0,1}^{(m)}, I_{k,l}\right) = \begin{cases} -1/\sqrt{6}, & \text{for } l = 1, \ m = k, \\ -1/3\sqrt{10}, & \text{for } l = 2, \ m = k, \\ 0, & \text{otherwise}, \end{cases}$$

$$\left(S_{0,-1}^{(m)}, I_{k,l}\right) = \begin{cases} -1/\sqrt{6}, & \text{for } l = 1, \ m = k, \\ -1/\sqrt{3}, & \text{for } l = 2, \ m = k, \\ 0, & \text{otherwise}. \end{cases}$$

In the block $NS$ there are only the following nonzero entries:

$$\left(N_{1,-1}, S_{0,-1}^{(m)}\right) = \left(N_{1,1}, S_{1,0}^{(m)}\right)$$
$$= \left(N_{1,1}, S_{-1,0}^{(m)}\right) = \left(N_{1,1}, S_{0,1}^{(m)}\right) = \begin{cases} -1/2\sqrt{6}, & m = 1, \\ -1/6\sqrt{10}, & m = 2, \end{cases}$$

$$\left(N_{1,1}, S_{1,0}^{(m)}\right) = \left(N_{1,1}, S_{0,1}^{(m)}\right)$$
$$= \left(N_{1,1}, S_{0,-1}^{(m)}\right) = \left(N_{-1,1}, S_{1,0}^{(m)}\right) = \begin{cases} 1/2\sqrt{6}, & m = 1, \\ 1/6\sqrt{10}, & m = 2, \end{cases}$$

$$\left(N_{1,1}, S_{1,0}^{(m)}\right) = \left(N_{-1,-1}, S_{-1,0}^{(m)}\right)$$
$$= \left(N_{1,1}, S_{0,1}^{(m)}\right) = \left(N_{-1,-1}, S_{0,-1}^{(m)}\right) = \begin{cases} -1/2\sqrt{6}, & m = 1, \\ 1/6\sqrt{10}, & m = 2, \end{cases}$$

$$\left(N_{-1,-1}, S_{0,1}^{(n)}\right) = \left(N_{-1,1}, S_{0,-1}^{(m)}\right)$$
$$= \left(N_{1,1}, S_{1,0}^{(m)}\right) = \left(N_{-1,-1}, S_{1,0}^{(m)}\right) = \begin{cases} 1/2\sqrt{6}, & m = 1, \\ -1/6\sqrt{10}, & m = 2. \end{cases}$$

Nonzero entries of diagonal block $SS$ are constructed by the scalar products of basic functions corresponded to the common sides:

$$\left(S_{k,l}^{(m)}, S_{k,l}^{(m')}\right) = \begin{cases} \dfrac{2}{3} + \dfrac{1}{2(2m-1)(2m+3)}, & m' = m, \\[4mm] \dfrac{-1}{2\sqrt{(2m+1)(2m+5)(2m+3)}}, & m' = m+2, \end{cases}$$

and corresponded to the opposite sides only:

$$\left(S_{k,l}^{(m)}, S_{-k,-l}^{(m')}\right) = \begin{cases} \dfrac{1}{3} - \dfrac{1}{2(2m-1)(2m+3)}, & m' = m, \\[3mm] \dfrac{1}{2\sqrt{(2m+1)(2m+5)(2m+3)}}, & m' = m+2. \end{cases}$$

And at last, nonzero entries of block $II$ are:

$$(I_{k,l}, I_{k',l'}) = \dfrac{2}{(2m-1)(2m+3)}$$

$$+ \begin{cases} \dfrac{2}{(2n-1)(2n+3)}, & \text{for } k = k' = m, \ l = l' = n, \\[3mm] 0, & \begin{aligned}&\text{for } k = k' = m, \ l' \neq l, \ l' \neq l+2,\\ &\text{or } l = l' = m, \ k' \neq k, \ k' \neq k+2,\end{aligned} \\[3mm] -\dfrac{1}{\sqrt{(2n+1)(2n+5)(2n+3)}}, & \begin{aligned}&\text{for } k = k' \pm 2 = m, \ l = l' = n\\ &\text{or } l = l' + 2 = m, \ k = k' = n.\end{aligned} \end{cases}$$

For describing the global stiffness matrix $A$ we have to define the ordering of all basic functions. Let the computational domain be a square divided into $n \times n$ elements by meshlines $x_i = x_0 + ih$, $y_j = y_0 + jh$; $i, j = 0, 1, \ldots, n$, $h$ is some meshsize. So, there are $n^2$ elements, $(n+1)^2$ nodes, $n(n+1)$ vertical sides (segments between nodes) and $n(n+1)$ horizontal sides. The total number of unknowns and the order of the global system of equations

$$Au = f, \quad A = D - L - U, \tag{2.5}$$

with square matrix $A$ equals to

$$N = (n+1)^2 + 2(p-1)n(n+1) + (p-1)^2 n^2 = (np+1)^2. \tag{2.6}$$

Here $D = \operatorname{diag}(A)$ and $L, U$ are strictly low and upper triangular parts of matrix $A$.

In (2.5) $u = \{u_k\}$ means the global vector which entries are numbered in the following way:

node entries in the row-by-row ordering –

$$\kappa = i + (j-1)(n+1); \quad i, j = 1, 2, \ldots, n+1,$$

vertical side entries in the similar ordering –

$$\kappa = (n+1)^2 + 1, \ldots, (n+1)^2 + (p-1)(n+1)n,$$

horizontal side entries –

$$\kappa = (n+1)[n+1+(p-1)n]+1,\ldots,(n+1)^2+2(p-1)(n+1)n,$$

internal entries in the element-by-element ordering –

$$\kappa = (n+1)[n+1+2(p-1)n]+1,\ldots,N.$$

Such matrix $A$ is constructed without taking into account the boundary conditions and is called unconstrained matrix. For obtaining the resulting constrained matrix it is necessary to modify the rows of the matrix $A$ and right side vector $f$ in (2.5), corresponded to the boundary nodes and sides.

If some basic function is connected with the boundary side or node with the Neumann condition, then corresponding row of matrix $A$ doesn't change. In the case of the Dirichlet condition the row becomes with the unit at the main diagonal and zero other entries.

The structure of the global stiffness matrix for the described ordering can be presented as a similar local one:

$$A = \begin{bmatrix} GNN & GNS & 0 \\ GSN & GSS & GSI \\ 0 & GIS & GII \end{bmatrix}.$$

Here blocks are

$GNN$ – square submatrix of order $(n+1)^2$,

$GSS$  – square block of order $2(p-1)n(n+1)$,

$GII$  – square submatrix of order $(p-1)^2n^2$,

$GNS = (GSN)'$ – rectangular matrix with $(n+1)^2$ rows and
$\qquad\qquad\quad 2(p-1)n(n+1)$ columns,

$GSI = (GIS)'$ – rectangular matrix of the form
$\qquad\qquad\quad [2(p-1)n(n+1)] \cdot [(p-1)^2n^2].$

Let us consider the structure of each block now (in an unconstrained form, i.e., without considering of boundary condition). Matrix $GNN$ is nine-diagonal with the semi-bandwidth $d = n+2$. In other words, this matrix is block tridiagonal of block order $n+1$ and each block is tridiagonal one (of order $n+1$). Nonzero entries of $GNN$ equal to 16/3 at the main diagonal and 2/3 at the others. The total number of nonzero entries of $GNN$ equals to

$$Q_{NN} = (3n+1)^2.$$

Submatrix $GSS$ is block-diagonal of block order 2 (first block is corresponded to the vertical sides and the second – to horizontal ones. Each diagonal block $GSS_k$, $k = 1, 2$, is block-diagonal matrix of block order $n$ with the equal blocks $GSS_0$ at the main diagonal ($GSS_k = \text{diag}\{GSS_0\}$), $GSS_0$ is block-tridiagonal matrix of block order $n + 1$ with the subblocks $S_{11} = \text{diag}\{S_1\}$, $S_{22} = \text{diag}\{S_2\}$, which are block-diagonal matries. Here $S_1$, $S_2$ are the same as in $SS$, see above. In other words, $GSS$ is nine-diagonal with semi-bandwidth $d = (p - 1)n$:

$$GSS_0 = \begin{array}{|c|c|c|c|}
\hline
S_{11} & S_{22} & & \\
\hline
S_{22} & S_{11} & S_{22} & \\
\hline
& S_{22} & S_{11} & S_{22} \\
\hline
& & S_{22} & S_{11} \\
\hline
\end{array}, \quad
GSS = \begin{array}{|c|c|}
\hline
GSS_1 & 0 \\
\hline
0 & GSS_2 \\
\hline
\end{array}.$$

So, the total number of nonzero entries of $GSS$ equals to

$$Q_{SS} = 2n(3n + 1)(3p - 7).$$

This equality is valid for $q > 3$. For $p = 2$ and $p = 3$ the values $(3p - 7)$ have to be replaced by 1 and 2 correspondingly. The biggest diagonal block $GII$ is block-diagonal matrix of block order $n^2$ and the order of each block is $(p - 1)^2$ (it is exactly the local stiffness matrix $II$, see above). So the total number of nonzero entries of $GII$ equals to

$$Q_{II} = n^2(p - 1)(5p - 13)$$

for $p \geq 3$ and $Q_{II} = n^2$ for $p = 2$.

Offdiagonal block $GNS$ has $(n + 1)^2$ rows, $2(p - 1)n(n + 1)$ columns and in general 24 nonzero entries in the each row. Maximum distance from nonzero entry to main diagonal equals to $2n(n + 1)(p - 1)$ approximately.

In block form $GNS$ can be divided by vertical line into two equal subblocks, the left one is corresponded to the "vertical" side functions and right – to the "horizontal" functions.

The left subblock is presented as two-block-diagonal matrix with $n + 1$ block rows and $n$ block columns. Corresponding sub-subblocks are identical rectangular matrices with $n + 1$ rows, $(n + 1)(p - 1)$ columns and 6 nonzero diagonals for $p \geq 3$ (3 – for $p = 2$).

The right subblock is block-tridiagonal matrix with $n + 1$ block rows and columns. Each sub-subblock is matrix with $n + 1$ rows, $n(p - 1)$ columns and 4 nonzero diagonals for $p \geq 3$ (2 – for $p = 2$).

$$GNS = \quad \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots$$

Block $GSI = (GIS)'$ has $2(p-1)n(n+1)$ rows, $(p-1)^2 n^2$ columns and in general 4 nonzero entries in the each row. The biggest distance from nonzero entry to main diagonal equals to $(p-1)^2 n^2$ approximately.

In block form $GSI$ is divided by vertical line into two equal subblocks. The left one corresponds to vertical sides and presents the block-diagonal matrix (of block order $n$) with identical sub-subblocks, consisted of $((n+1)(p-1)$ rows, $n(p-1)^2$ columns and 4 nonzero diagonals. The right subblock indcludes $n+1$ block rows, $n$ block columns and two block diagonals with identical sub-subblocks consisted of $n(p-1)$ columns and nonzero diagonals.

The total number of nonzero entries of $GSI$ equals to

$$Q_{SI} = 8n^2,$$

and of matrix $A$ –

$$Q_A = Q_{NN} + Q_{SS} + Q_{II} + 2Q_{NS} + 2Q_{SI} \approx n^2(12 + 5p^2).$$

The last approximate equality is valid for $p > 3$, and

$$Q_A \approx 57n^2 \text{ for } p = 3, \quad Q_A \approx 48n^2 \text{ for } p = 2.$$

## 3.  An iterative algorithms

We apply the preconditioned conjugate gradient methods in the form

$$
\begin{aligned}
&r^\circ = f - Au^\circ, \quad p^\circ = B^{-1}r^\circ, \\
&u^{k+1} = u^k + \alpha_k p^k, \quad \alpha_k = (r^k, B^{-1}r^k)/(Ap^k, p^k), \\
&r^{k+1} = r^k - \alpha_k Ap^k, \quad p^{k+1} = B^{-1}r^{k+1} + \beta_k p^k, \\
&\beta_k = (r^{k+1}, B^{-1}r^{k+1})/(r^k, B^{-1}r^k).
\end{aligned}
\tag{3.1}
$$

Here $B$ is symmetric preconditioning matrix

$$B = (G - L)G^{-1}(G - U) \qquad (3.2)$$

with some diagonal matrix $G = 1/\omega D$, $0 < \omega < 2$, corresponded to the well-known Symmetric Successive Over Relaxation (SSOR) method.

Now we describe an efficient implementation of iterative process (3.1), based on some modification of Eisenstat's idea [4]. Let us introduce the vectors

$$\bar{u}^k = U_b u^k, \quad \bar{p}^k = U_b p^k, \quad \bar{r}^k = L_b^{-1} r^k, \quad \bar{f} = L_b^{-1} f,$$

$$B = L_b U_b, \quad U_b = L_b' = G^{1/2} - G^{-1/2} U \qquad (3.3)$$

and matrix

$$\bar{A} = L_b^{-1} A U_b = (E - I)^{-1} + (E - \bar{U})^{-1}$$
$$+ (\bar{E} - \bar{L})^{-1}(2E - \bar{D})(E - \bar{U})^{-1},$$
$$\bar{L} = G^{-1/2} L G^{-1/2}, \quad \bar{U} = G^{-1/2} U G^{-1/2}, \qquad (3.4)$$
$$\bar{D} = G^{-1/2} D G^{-1/2}.$$

Then from 3.1 for new vectors there is equivalent iterative process

$$\bar{r}^\circ = \bar{f} - \bar{A}\bar{u}^\circ, \quad \bar{p}^\circ = \bar{r}^\circ,$$
$$\bar{u}^{k+1} = \bar{u}^k + \alpha_k \bar{p}^k, \quad \bar{r}^{k+1} = \bar{r}^k - \alpha_k \bar{A}\bar{p}^k, \qquad (3.5)$$
$$\bar{p}^{k+1} = \bar{r}^{k+1} + \beta_k \bar{p}^k$$

with identity preconditioning matrix $\bar{K} = L_b^{-1} K U_b^{-1} = E$ and the same parameters $\alpha_k$, $\beta_k$.

It is essential, that multiplication the vector by the matrix $\bar{A}$ is cheap enough operation:

$$\bar{A}\bar{p}^k = (E - \bar{L})^{-1}\bar{t}^k + \bar{q}^k,$$
$$\bar{q}^k = (E - \bar{U})^{-1}\bar{p}^k, \quad \bar{t}^k = \bar{p}^k + (2E - \bar{D})\bar{q}^k. \qquad (3.6)$$

For not very sparse matrix $A$ (more than 3 nonzero entries in each row in average) the implementation of formulas (3.9) is twice cheaper approximately, in compare with (3.1). Of course, it is supposed that the matrices $\bar{L}$, $\bar{U}$ and $\bar{D}$ are computed before iterations and resulting vector $u^n = U_b^{-1}\bar{u}^n$ is computed after the iterations.

It is evidently that matrix $A$ is positive definite, but only for $p = 1$ it is of the Stiltjes type (non-positive offdiagonal elements).

An estimates of the convergence velocity of SSOR-iterations for such kind of equations is unknown and our goal is only to make experimental investigations.

# 4.   The results of numerical experiments

The iterations in experiments were performed until the validity of inequality

$$\| r^k \|_2 / \| r^\circ \|_2 \le \varepsilon \tag{4.1}$$

for given small enough $\varepsilon$.

Our main criteria of efficiency of algorithms is the number of iterations. The total number of arithmetic operations for solution of system is evaluated approximately by

$$Q = n_\epsilon \cdot (n_c + n_a) \tag{4.2}$$

where $n_a$ is the cost of multiplication $\bar{A} \times \bar{p}$, $n_c = 10N$ is the rest cost of conjugate gradient algorithm per one iteration and $n_\epsilon$ is the number of iterations.

As followes from (3.6), the values $n_a$ can be estimated by the approximate equality

$$n_a = 2(Q_A + N),$$

where $N$, $Q_A$ are the order and the total number of nonzero entries of matrix $A$. By means of evaluation of $Q_A$, $N$ for different $p$ we can write

$$n_a = \begin{cases} 2[(3n + 1)^2 + (n + 1)^2] \approx 20n^2, & p = 1, \\ 104n^2, & p = 2, \\ 132n^2, & p = 3, \\ 2[n^2(12 + 5p^2) + (1 + pn)^2], & p > 3. \end{cases}$$

Calculations were made in double precisions on the IBM PC AT with coprocessor and frequency 12 Mg, in the environment MS-DOS 4.0 with available core $\sim$550 Kbyte which permitted to compute variants for $p \le 17$ only.

In the following tables there are numerical results for the model problem in the unit square $\Omega = (0, 1) \times (0, 1)$ with $f(x, y) = 0$ in (1.1) and the Neumann boundary condition $\partial u / \partial n = g(x, y)$ except one angle point of the computational square with the Dirichlet condition. The boundary conditions are defined from the choosen exact solution

$$u_e = Re\left(\frac{1}{a^2 + z^2} + \frac{1}{a^2 - z^2}\right), \quad a = 1.05, \quad z = x + \sqrt{-1}y. \tag{4.3}$$

In each cell of Table 1 there are the numbers of iterations (above) for SSOR-preconditioner, $\omega = 1.5$, $\epsilon = 10^{-5}$, and numbers of unknowns below for different $n$ and $p$.

**Table 1.** Numbers of unknowns and iterations for SSORCG, $\epsilon = 10^{-5}$, $\omega = 1.5$

| n \ p | 1 | 2 | 3 | 4 | 6 | 8 | 9 | 12 | 16 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 3 | 8 | 13 | 18 | 25 | 33 | 35 | 49 | 62 |
|   | 4 | 9 | 16 | 25 | 49 | 81 | 100 | 169 | 289 |
| 2 | 7 | 15 | 19 | 25 | 36 | 44 | 46 | 59 | 74 |
|   | 9 | 25 | 49 | 81 | 169 | 289 | 361 | 625 | 1089 |
| 4 | 9 | 20 | 21 | 30 | 38 | 46 | 48 | 62 | 76 |
|   | 25 | 81 | 169 | 289 | 625 | 1089 | 1369 | 2401 | 4225 |
| 8 | 14 | 22 | 23 | 33 | 42 | 51 | 52 | – | – |
|   | 81 | 289 | 625 | 1089 | 2401 | 4225 | 5329 | – | – |
| 16 | 20 | 26 | 27 | 36 | – | – | – | – | – |
|   | 289 | 1089 | 2401 | 4225 | – | – | – | – | – |
| 31 | 33 | 36 | – | – | – | – | – | – | – |
|   | 1024 | 3269 | – | – | – | – | – | – | – |

**Table 2.** Time costs of iterative solution (SSORCG, $\epsilon = 10^{-5}$, $\omega = 1.5$) and computing of matrices

| n \ p | 1 | 2 | 3 | 4 | 6 | 8 | 9 | 12 | 16 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 2.3 | 2.9 | 3.8 | 4.7 | 7.8 | 12.9 | 16.8 | 36.3 | 87.3 |
|   | 0.4 | 0.6 | 0.9 | 1.3 | 2.4 | 4.4 | 5.4 | 12.0 | 24.2 |
| 2 | 3.3 | 4.6 | 6.2 | 7.8 | 12.9 | 20.8 | 25.8 | 50.6 | 115.1 |
|   | 0.6 | 1.1 | 2.3 | 4.2 | 10.0 | 19.3 | 24.3 | 51.1 | 107.2 |
| 4 | 5.6 | 8.5 | 12.3 | 15.6 | 25.6 | 40.0 | 49.7 | 94.0 | 205.0 |
|   | 0.8 | 3.8 | 7.8 | 16.3 | 38.3 | 74.9 | 95.9 | 206.2 | 427.9 |
| 8 | 11.3 | 22.3 | 30.6 | 40.4 | 66.1 | 105.0 | 129.7 | – | – |
|   | 2.6 | 13.9 | 30.7 | 66.2 | 161.9 | 321.9 | 404.6 | – | – |
| 16 | 36.4 | 117.4 | 151.9 | 193.0 | – | – | – | – | – |
|   | 10.7 | 61.4 | 137.0 | 280.5 | – | – | – | – | – |
| 31 | 255.5 | 1172.6 | – | – | – | – | – | – | – |
|   | 59.2 | 304.7 | – | – | – | – | – | – | – |

Table 2 includes the time costs of iterative solution and computing of the global stiffness matrices (in seconds) for the same experiments.

For analysing of the influence of values $\varepsilon$ in (4.1) on the numbers of iterations and the accuracy of iterative solution $u^k$ we show in Table 3 the values of $k_\epsilon$ and the errors

$$\delta = \frac{\| u_{ex} \| - \| u_{fe} \|}{\| u_{ex} \|} \qquad (4.4)$$

for the similar experiments (SSORCG, $\omega = 1.5$) with $\epsilon = 10^{-3}$, $10^{-4}$, $10^{-5}$ and $p = 2$ only ($u_{ex}$ and $u_{fe}$ are exact and finite element solution, $\| u \| = (Au, u))^{1/2}$).

**Table 3.** The numbers of iterations and errors of iterative solutions ($p = 2$, $\omega = 1.5$)

| $\epsilon \setminus n$ | 1 | 2 | 4 | 8 | 16 | 31 |
|---|---|---|---|---|---|---|
| $10^{-3}$ | 7 <br> $3.9 \cdot 10^{-1}$ | 11 <br> $2.4 \cdot 10^{-1}$ | 13 <br> $1.1 \cdot 10^{-1}$ | 15 <br> $4.4 \cdot 10^{-2}$ | 19 <br> $1.1 \cdot 10^{-2}$ | 29 <br> $8.9 \cdot 10^{-4}$ |
| $10^{-4}$ | 8 <br> $3.9 \cdot 10^{-1}$ | 13 <br> $2.4 \cdot 10^{-1}$ | 17 <br> $1.1 \cdot 10^{-1}$ | 18 <br> $4.4 \cdot 10^{-2}$ | 22 <br> $1.1 \cdot 10^{-2}$ | 33 <br> $8.9 \cdot 10^{-4}$ |
| $10^{-5}$ | 8 <br> $3.9 \cdot 10^{-1}$ | 15 <br> $2.4 \cdot 10^{-1}$ | 20 <br> $1.1 \cdot 10^{-1}$ | 22 <br> $4.4 \cdot 10^{-2}$ | 26 <br> $1.1 \cdot 10^{-2}$ | 36 <br> $8.9 \cdot 10^{-4}$ |

Let us make one more observation – what values of $p$ and $n$ are optimal for necessary accuracy $\Delta$. In Table 4 there are for different $p$ the values of the time $t_o$ of SSORCG iterations (in seconds, $\omega = 1.5$), numbers of elements $n_o$ (in one direction) and numbers of corresponding iterations under condition

$$n_o = \min\{n : \delta \le \Delta\}.$$

In every cell the upper value corresponds to $\Delta = 10^{-3}$ and low to $\Delta = 10^{-4}$.

**Table 4.** The values of iteration time costs and errors for different $p$ and $n_o(p, \Delta)$

| $p \setminus n_0$ | 3 | 4 | 6 | 8 | 12 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|
| $t_0$ | 117 <br> 259 | 83 <br> 181 | 60 <br> 125 | 43 <br> 120 | 51 <br> 113 | 17 <br> 77 | 20 <br> 89 | 24 <br> 107 |
| $n_0$ | 15 <br> 21 | 9 <br> 13 | 5 <br> 7 | 3 <br> 5 | 2 <br> 3 | 1 <br> 2 | 1 <br> 2 | 1 <br> 2 |
| $\delta$ | $8 \cdot 10^{-4}$ <br> $8 \cdot 10^{-5}$ | $8 \cdot 10^{-4}$ <br> $6 \cdot 10^{-5}$ | $4 \cdot 10^{-4}$ <br> $7 \cdot 10^{-5}$ | $2 \cdot 10^{-4}$ <br> $5 \cdot 10^{-5}$ | $6 \cdot 10^{-5}$ <br> $1 \cdot 10^{-5}$ | $3 \cdot 10^{-4}$ <br> $3 \cdot 10^{-5}$ | $2 \cdot 10^{-4}$ <br> $1 \cdot 10^{-5}$ | <br> $4 \cdot 10^{-6}$ |

For understanding the sources of efficiency of the preconditioned conjugate gradient method it is usefull to obtain separately the influence of preconditioning and conjugate gradient acceleration. In Table 5 we give the numbers of iterations ($\epsilon = 10^{-5}$ as in Table 1) for SSOR-steepest descent (SSORSD) method ($\omega = 1.5$, $\beta_k \equiv 0$ in (3.1) or (3.8)) and for "explicit" conjugate gradient iterative process (in (3.1) $B \equiv E$ is identity matrix).

It is important to remark that for large enough $n$, $p$ the numbers of arithmetical operations per one iteration are approximately the same for each method – SSORCG, SSORSD and explicit CG.

**Table 5.** The numbers of iterations for SSORSD (above) and explicit CG (below)

| n \ p | 2 | 3 | 4 | 6 | 8 | 9 | 12 | 16 |
|-------|------|------|------|-----|-----|-----|-----|-----|
| 1     | 41   | 39   | 99   | 177 | 275 | 275 | 489 | 727 |
|       | 7    | 13   | 20   | 37  | 53  | 60  | 88  | 120 |
| 2     | 49   | 63   | 107  | 183 | 263 | 265 | 445 | 655 |
|       | 20   | 28   | 41   | 58  | 78  | 79  | 111 | 133 |
| 4     | 117  | 139  | 165  | 219 | 291 | 313 | 479 | 699 |
|       | 33   | 37   | 57   | 74  | 91  | 88  | 125 | 152 |
| 8     | 349  | 393  | 427  | 491 | 555 | 583 | –   | –   |
|       | 45   | 49   | 64   | 84  | 102 | 101 | –   | –   |
| 16    | 1151 | 1259 | 1351 | –   | –   | –   | –   | –   |
|       | 72   | 75   | 90   | –   | –   | –   | –   | –   |

We didn't make the search of the optimal iterative parameter. But some experiments are presented in Table 6 for SSORCG, $\epsilon = 10^{-5}$. In each cell the upper number corresponds to $\omega = 1.1$, the middle – to $\omega = 1.2$ and the lower – to $\omega = 1.3$.

**Table 6.** Numbers of iterations for different $\omega$ (1.1, 1.2, 1.3), SSORCG, $\epsilon = 10^{-5}$

| n \ p | 2  | 3  | 4  | 8  | 12 | 16 |
|-------|----|----|----|----|----|----|
| 1     | –  | 10 | 15 | 26 | 40 | 53 |
|       | 7  | 10 | 13 | 27 | 40 | 53 |
|       | 7  | –  | –  | –  | –  | –  |
| 2     | –  | 15 | 21 | 35 | 49 | 61 |
|       | 13 | 16 | 21 | 36 | 49 | 63 |
|       | 13 | –  | –  | –  | –  | –  |
| 4     | –  | 17 | 25 | 39 | 53 | 64 |
|       | 17 | 17 | 25 | 39 | 53 | 65 |
|       | 17 | –  | –  | –  | –  | –  |
| 8     | –  | 19 | 28 | 44 | –  | –  |
|       | 19 | 19 | 27 | 45 | –  | –  |
|       | 19 | –  | –  | –  | –  | –  |
| 16    | –  | 27 | 32 | –  | –  | –  |
|       | 26 | 27 | 33 | –  | –  | –  |
|       | 25 | –  | –  | –  | –  | –  |

In some sense, the considered discrete boundary value problem (the Dirichlet condition in one point only) is the worst. Intuitively, for the other boundary value conditions the number of iterations has to be decrease. An additional experiments confirm such conjecture. For example,

we give in Table 7 the numbers of iterations for the Dirichlet boundary value problem with the conditions, corresponded to exact solution of the Laplace equation $u = \sin \pi x \cdot \sh \pi y / \sh \pi$. In each square of Table 7 the numbers are appropriated to $\epsilon = 10^{-5}, 10^{-6}, 10^{-7}$, zero initial data $u^{\circ}$ and $\omega = 1, 5$.

**Table 7.** Numbers of iteration for the Dirichlet boundary value problem ($\epsilon = 10^{-5}, 10^{-6}, 10^{-7}$)

| n \ p | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 2 | 1, 1 ,1 | 6, 6 ,6 | 9, 11, 12 | 11, 16, 19 |
| 4 | 5, 6 ,7 | 11, 15, 19 | 13, 19, 25 | 17, 25, 32 |
| 8 | 6, 8 ,10 | 11, 18, 25 | 12, 19, 26 | 16, 27, 38 |
| 16 | 8, 11, 14 | 13, 20, 27 | | |

There are similar pictures for another Dirichlet boundary value problems with different solutions. The analysis of preceding experimental data permits to make some preliminary conclusions.

**a.** The application of conjugate gradient acceleration has very high efficiency (in compare with steepest descent) for every values of $n$ and $p$. The speed up by SSOR-preconditioner for large $n$, $p$ is approximately two, for small $n$ and $p$ the preconditioning effect decreases.

**b.** The sensivity of SSORCG method to values of iterative parameter $\omega$ is small enough and differences in time costs are approximately 10–20 percent (an optimal values $\omega_{opt} \approx 1.1 \div 1.3$).

**c.** The increasing of "stop-iteration" criteria $\epsilon$ in (4.1) from $10^{-5}$ to $10^{-3}$ doesn't change the resulting accuracy practically, see Table 3.

**d.** For given accuracy the optimal pairs $(n, p)$ correspond to minimal $n = 1$ or 2 (in the sense of minimal iteration time cost).

Of course, there are a lot of open questions on the application and optimization of iterative methods for solving high order finite element equations:

- The choosing of iterative parameters and criteria $\epsilon$ for necessary accuracy $\Delta$ of numerical solution or given values $n$ and $p$;

- The role of "implicitness" of incomplete factorization methods for efficiency of iterative process;

- Adapting of iterative algorithms for different properties of original continuous problem (coefficients of equations, boundary conditions, domain shape and so on);

- The influence of uknowns ordering on the constructing of optimal iterative processes.

These and other aspects present the wide field of activity in the topic of consideration to find the main goal: the search of optimal method for solving a class of problems with given accuracy.

# References

[1] B. Szabo, I. Babuska, Finite Element Analysis, John Wiley & Sons, 1991.

[2] L.A. Hageman, D.M. Young, Applied Iterative Methods, Academic Press, 1981.

[3] V.P. Il'in, Iterative Incomplete Factorization Methods, World Scientific Publishing Co., 1992.

[4] S.C. Eisenstat, Efficient implementation of a class of preconditioned conjugate gradient methods, SIAM J. Sci. Stat. Comput., Vol. 2, 1–4, 1981.