

MIXED FINITE ELEMENT METHOD FOR THE STEADY STOKES EQUATIONS

WILL FREY*

Abstract. NOT FINISHED A mixed finite element method is derived for the steady state Stokes flow equation. Derivation of the weak form solution is provided and the corresponding system is written in a system of linear equations. An example of SIAM L^AT_EX macros is presented. Various aspects of composing manuscripts for SIAM's journal series are illustrated with actual examples from accepted manuscripts. SIAM's stylistic standards are adhered to throughout, and illustrated.

Key words. steady stokes flow, finite element method, mixed finite elements

1. Introduction. The Stokes equation (or *Stokes flow*) is derived by making certain simplifying assumptions and applying them to the non-linear Navier-Stokes equations. The flow of an incompressible viscous fluid in an n -dimensional domain (with $n = 2$ or 3)[?]:

$$(1.1) \quad \nabla^2 u + \nabla p = -f \quad \text{in } \Omega,$$

$$(1.2) \quad \nabla u = 0 \quad \text{in } \Omega,$$

$$(1.3) \quad u = u_0 \quad \text{on } \Omega.$$

This paper presents a sample file for the use of SIAM's L^AT_EX macro package. It illustrates the features of the macro package, using actual examples culled from various papers published in SIAM's journals. It is to be expected that this sample will provide examples of how to use the macros to generate standard elements of journal papers, e.g., theorems, definitions, or figures. This paper also serves as an example of SIAM's stylistic preferences for the formatting of such elements as bibliographic references, displayed equations, and equation arrays, among others. Some special circumstances are not dealt with in this sample file; for such information one should see the included documentation file.

Note: This paper is not to be read in any form for content. The conglomeration of equations, lemmas, and other text elements were put together solely for typographic illustrative purposes and don't make any sense as lemmas, equations, etc.

1.1. The Stokes Equation. The Stokes Equation, otherwise known as *Stokes Flow*, is the linearized form of the Navier-Stokes Equation.

characterizing SNS-matrices [?], [?]. There has also been interest in strong forms of sign-nonsingularity [?]. In this paper we give a new generalization of SNS-matrices and investigate some of their basic properties.

Let $S = [s_{ij}]$ be a $(0, 1, -1)$ -matrix of order n and let $C = [c_{ij}]$ be a real matrix of order n . The pair (S, C) is called a *matrix pair of order n* . Throughout, $X = [x_{ij}]$ denotes a matrix of order n whose entries are algebraically independent indeterminates over the real field. Let $S \circ X$ denote the Hadamard product (entrywise product) of

*Thanks to Dave Wells for providing me with mesh parsing software and teaching me how to use gmsh. Thanks also for the endless Borggaard for the "rules of thumb" and other guidance. Thanks also to Jeff Borggaard for his helpful "rules of thumb."

S and X . We say that the pair (S, C) is a *sign-nonsingular matrix pair of order n* , abbreviated *SNS-matrix pair of order n* , provided that the matrix

$$A = S \circ X + C$$

is nonsingular for all positive real values of the x_{ij} . If $C = O$ then the pair (S, O) is a SNS-matrix pair if and only if S is a SNS-matrix. If $S = O$ then the pair (O, C) is a SNS-matrix pair if and only if C is nonsingular. Thus SNS-matrix pairs include both nonsingular matrices and sign-nonsingular matrices as special cases.

The pairs (S, C) with

$$S = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

and

$$S = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 2 & 0 \\ 3 & 0 & 0 \end{bmatrix}$$

are examples of SNS-matrix pairs.

1.2. A remuneration list. In this paper we consider the evaluation of integrals of the following forms:

$$(1.4) \quad \int_a^b \left(\sum_i E_i B_{i,k,x}(t) \right) \left(\sum_j F_j B_{j,l,y}(t) \right) dt,$$

$$(1.5) \quad \int_a^b f(t) \left(\sum_i E_i B_{i,k,x}(t) \right) dt,$$

where $B_{i,k,x}$ is the i th B-spline of order k defined over the knots $x_i, x_{i+1}, \dots, x_{i+k}$. We will consider B-splines normalized so that their integral is one. The splines may be of different orders and defined on different knot sequences x and y . Often the limits of integration will be the entire real line, $-\infty$ to $+\infty$. Note that (??) is a special case of (??) where $f(t)$ is a spline.

There are five different methods for calculating (??) that will be considered:

1. Use Gauss quadrature on each interval.
2. Convert the integral to a linear combination of integrals of products of B-splines and provide a recurrence for integrating the product of a pair of B-splines.
3. Convert the sums of B-splines to piecewise Bézier format and integrate segment by segment using the properties of the Bernstein polynomials.
4. Express the product of a pair of B-splines as a linear combination of B-splines. Use this to reformulate the integrand as a linear combination of B-splines, and integrate term by term.
5. Integrate by parts.

Of these five, only methods 1 and 5 are suitable for calculating (??). The first four methods will be touched on and the last will be discussed at length.

1.3. Some displayed equations and $\{eqnarray\}$ s. By introducing the product topology on $R^{m \times m} \times R^{n \times n}$ with the induced inner product

$$(1.6) \quad \langle (A_1, B_1), (A_2, B_2) \rangle := \langle A_1, A_2 \rangle + \langle B_1, B_2 \rangle,$$

we calculate the Fréchet derivative of F as follows:

$$(1.7) \quad \begin{aligned} F'(U, V)(H, K) &= \langle R(U, V), H\Sigma V^T + U\Sigma K^T - P(H\Sigma V^T + U\Sigma K^T) \rangle \\ &= \langle R(U, V), H\Sigma V^T + U\Sigma K^T \rangle \\ &= \langle R(U, V)V\Sigma^T, H \rangle + \langle \Sigma^T U^T R(U, V), K^T \rangle. \end{aligned}$$

In the middle line of (??) we have used the fact that the range of R is always perpendicular to the range of P . The gradient ∇F of F , therefore, may be interpreted as the pair of matrices:

$$(1.8) \quad \nabla F(U, V) = (R(U, V)V\Sigma^T, R(U, V)^T U\Sigma) \in R^{m \times m} \times R^{n \times n}.$$

Because of the product topology, we know

$$(1.9) \quad \mathcal{T}_{(U, V)}(\mathcal{O}(m) \times \mathcal{O}(n)) = \mathcal{T}_U \mathcal{O}(m) \times \mathcal{T}_V \mathcal{O}(n),$$

where $\mathcal{T}_{(U, V)}(\mathcal{O}(m) \times \mathcal{O}(n))$ stands for the tangent space to the manifold $\mathcal{O}(m) \times \mathcal{O}(n)$ at $(U, V) \in \mathcal{O}(m) \times \mathcal{O}(n)$ and so on. The projection of $\nabla F(U, V)$ onto $\mathcal{T}_{(U, V)}(\mathcal{O}(m) \times \mathcal{O}(n))$, therefore, is the product of the projection of the first component of $\nabla F(U, V)$ onto $\mathcal{T}_U \mathcal{O}(m)$ and the projection of the second component of $\nabla F(U, V)$ onto $\mathcal{T}_V \mathcal{O}(n)$. In particular, we claim that the projection $g(U, V)$ of the gradient $\nabla F(U, V)$ onto $\mathcal{T}_{(U, V)}(\mathcal{O}(m) \times \mathcal{O}(n))$ is given by the pair of matrices:

$$(1.10) \quad g(U, V) = \left(\frac{R(U, V)V\Sigma^T U^T - U\Sigma V^T R(U, V)^T}{2} U, \frac{R(U, V)^T U\Sigma V^T - V\Sigma^T U^T R(U, V)}{2} V \right).$$

Thus, the vector field

$$(1.11) \quad \frac{d(U, V)}{dt} = -g(U, V)$$

defines a steepest descent flow on the manifold $\mathcal{O}(m) \times \mathcal{O}(n)$ for the objective function $F(U, V)$.

2. Main results. Let (S, C) be a matrix pair of order n . The determinant

$$\det(S \circ X + C)$$

is a polynomial in the indeterminates of X of degree at most n over the real field. We call this polynomial the *indicator polynomial* of the matrix pair (S, C) because of the following proposition.

THEOREM 2.1. *The matrix pair (S, C) is a SNS-matrix pair if and only if all the nonzero coefficients in its indicator polynomial have the same sign and there is at least one nonzero coefficient.*

Proof. Assume that (S, C) is a SNS-matrix pair. Clearly the indicator polynomial has a nonzero coefficient. Consider a monomial

$$(2.1) \quad b_{i_1, \dots, i_k; j_1, \dots, j_k} x_{i_1 j_1} \cdots x_{i_k j_k}$$

occurring in the indicator polynomial with a nonzero coefficient. By taking the x_{ij} that occur in (??) large and all others small, we see that any monomial that occurs in the indicator polynomial with a nonzero coefficient can be made to dominate all others. Hence all the nonzero coefficients have the same sign. The converse is immediate. \square

For SNS-matrix pairs (S, C) with $C = O$ the indicator polynomial is a homogeneous polynomial of degree n . In this case Theorem ?? is a standard fact about SNS-matrices.

LEMMA 2.2 (Stability). *Given $T > 0$, suppose that $\|\epsilon(t)\|_{1,2} \leq h^{q-2}$ for $0 \leq t \leq T$ and $q \geq 6$. Then there exists a positive number B that depends on T and the exact solution ψ only such that for all $0 \leq t \leq T$,*

$$(2.2) \quad \frac{d}{dt} \|\epsilon(t)\|_{1,2} \leq B(h^{q-3/2} + \|\epsilon(t)\|_{1,2}).$$

The function $B(T)$ can be chosen to be nondecreasing in time.

THEOREM 2.3. *The maximum number of nonzero entries in a SNS-matrix S of order n equals*

$$\frac{n^2 + 3n - 2}{2}$$

with equality if and only if there exist permutation matrices such that $P|S|Q = T_n$ where

$$(2.3) \quad T_n = \begin{bmatrix} 1 & 1 & \cdots & 1 & 1 & 1 \\ 1 & 1 & \cdots & 1 & 1 & 1 \\ 0 & 1 & \cdots & 1 & 1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 1 & 1 \\ 0 & 0 & \cdots & 0 & 1 & 1 \end{bmatrix}.$$

We note for later use that each submatrix of T_n of order $n - 1$ has all 1s on its main diagonal.

We now obtain a bound on the number of nonzero entries of S in a SNS-matrix pair (S, C) in terms of the degree of the indicator polynomial. We denote the strictly upper triangular $(0,1)$ -matrix of order m with all 1s above the main diagonal by U_m . The all 1s matrix of size m by p is denoted by $J_{m,p}$.

PROPOSITION 2.4 (Convolution theorem). *Let*

$$a * u(t) = \int_0^t a(t - \tau)u(\tau)d\tau, \quad t \in (0, \infty).$$

Then

$$\widehat{a * u}(s) = \widehat{a}(s)\widehat{u}(s).$$

LEMMA 2.5. For $s_0 > 0$, if

$$\int_0^\infty e^{-2s_0 t} v^{(1)}(t) v(t) dt \leq 0,$$

then

$$\int_0^\infty e^{-2s_0 t} v^2(t) dt \leq \frac{1}{2s_0} v^2(0).$$

Proof. Applying integration by parts, we obtain

$$\begin{aligned} \int_0^\infty e^{-2s_0 t} [v^2(t) - v^2(0)] dt &= \lim_{t \rightarrow \infty} \left(-\frac{1}{2s_0} e^{-2s_0 t} v^2(t) \right) + \frac{1}{s_0} \int_0^\infty e^{-2s_0 t} v^{(1)}(t) v(t) dt \\ &\leq \frac{1}{s_0} \int_0^\infty e^{-2s_0 t} v^{(1)}(t) v(t) dt \leq 0. \end{aligned}$$

Thus

$$\int_0^\infty e^{-2s_0 t} v^2(t) dt \leq v^2(0) \int_0^\infty e^{-2s_0 t} dt = \frac{1}{2s_0} v^2(0). \quad \square$$

COROLLARY 2.6. Let \mathbf{E} satisfy (5)–(6) and suppose \mathbf{E}^h satisfies (7) and (8) with a general \mathbf{G} . Let $\mathbf{G} = \nabla \times \Phi + \nabla p$, $p \in H_0^1(\Omega)$. Suppose that ∇p and $\nabla \times \Phi$ satisfy all the assumptions of Theorems 4.1 and 4.2, respectively. In addition suppose all the regularity assumptions of Theorems 4.1–4.2 are satisfied. Then for $0 \leq t \leq T$ and $0 < \epsilon \leq \epsilon_0$ there exists a constant $C = C(\epsilon, T)$ such that

$$\|(\mathbf{E} - \mathbf{E}^h)(t)\|_0 \leq Ch^{k+1-\epsilon},$$

where C also depends on the constants given in Theorems 4.1 and 4.2.

DEFINITION 2.7. Let S be an isolated invariant set with isolating neighborhood N . An index pair for S is a pair of compact sets (N_1, N_0) with $N_0 \subset N_1 \subset N$ such that:

- (i) $cl(N_1 \setminus N_0)$ is an isolating neighborhood for S .
- (ii) N_i is positively invariant relative to N for $i = 0, 1$, i.e., given $x \in N_i$ and $x \cdot [0, t] \subset N$, then $x \cdot [0, t] \subset N_i$.
- (iii) N_0 is an exit set for N_1 , i.e. if $x \in N_1$, $x \cdot [0, \infty) \not\subset N_1$, then there is a $T \geq 0$ such that $x \cdot [0, T] \subset N_1$ and $x \cdot T \in N_0$.

2.1. Numerical experiments. We conducted numerical experiments in computing inexact Newton steps for discretizations of a *modified Bratu problem*, given by

$$(2.4) \quad \begin{aligned} \Delta w + ce^w + d \frac{\partial w}{\partial x} &= f \quad \text{in } D, \\ w &= 0 \quad \text{on } \partial D, \end{aligned}$$

where c and d are constants. The actual Bratu problem has $d = 0$ and $f \equiv 0$. It provides a simplified model of nonlinear diffusion phenomena, e.g., in combustion and semiconductors, and has been considered by Glowinski, Keller, and Rheinhardt [?], as well as by a number of other investigators; see [?] and the references therein. See also

problem 3 by Glowinski and Keller and problem 7 by Mittelmann in the collection of nonlinear model problems assembled by Moré [?]. The modified problem (??) has been used as a test problem for inexact Newton methods by Brown and Saad [?].

In our experiments, we took $D = [0, 1] \times [0, 1]$, $f \equiv 0$, $c = d = 10$, and discretized (??) using the usual second-order centered differences over a 100×100 mesh of equally spaced points in D . In GMRES(m), we took $m = 10$ and used fast Poisson right preconditioning as in the experiments in §2. The computing environment was as described in §2. All computing was done in double precision.

FIG. 2.1. \log_{10} of the residual norm versus the number of GMRES(m) iterations for the finite difference methods.

In the first set of experiments, we allowed each method to run for 40 iterations, starting with zero as the initial approximate solution, after which the limit of residual norm reduction had been reached. The results are shown in Fig. ?? . In Fig. ?? , the top curve was produced by method FD1. The second curve from the top is actually a superposition of the curves produced by methods EHA2 and FD2; the two curves are visually indistinguishable. Similarly, the third curve from the top is a superposition of the curves produced by methods EHA4 and FD4, and the fourth curve from the top, which lies barely above the bottom curve, is a superposition of the curves produced by methods EHA6 and FD6. The bottom curve was produced by method A.

In the second set of experiments, our purpose was to assess the relative amount of computational work required by the methods which use higher-order differencing to reach comparable levels of residual norm reduction. We compared pairs of methods EHA2 and FD2, EHA4 and FD4, and EHA6 and FD6 by observing in each of 20 trials the number of iterations, number of F -evaluations, and run time required by each method to reduce the residual norm by a factor of ϵ , where for each pair of methods ϵ was chosen to be somewhat greater than the limiting ratio of final to initial residual norms obtainable by the methods. In these trials, the initial approximate solutions were obtained by generating random components as in the similar experiments in §2. We note that for every method, the numbers of iterations and F -evaluations required before termination did not vary at all over the 20 trials. The iteration counts, numbers of F -evaluations, and means and standard deviations of the run times are given in Table ?? .

TABLE 2.1

Statistics over 20 trials of GMRES(m) iteration numbers, F -evaluations, and run times required to reduce the residual norm by a factor of ϵ . For each method, the number of GMRES(m) iterations and F -evaluations was the same in every trial.

Method	ϵ	Number of Iterations	Number of F -Evaluations	Mean Run Time (Seconds)	Standard Deviation
EHA2	10^{-10}	26	32	47.12	.1048
FD2	10^{-10}	26	58	53.79	.1829
EHA4	10^{-12}	30	42	56.76	.1855
FD4	10^{-12}	30	132	81.35	.3730
EHA6	10^{-12}	30	48	58.56	.1952
FD6	10^{-12}	30	198	100.6	.3278

FIG. 2.2. \log_{10} of the residual norm versus the number of GMRES(m) iterations for $c = d = 10$ with fast Poisson preconditioning. Solid curve: Algorithm EHA; dotted curve: FDP method; dashed curve: FSP method.

In our first set of experiments, we took $c = d = 10$ and used right preconditioning with a fast Poisson solver from FISHPACK [?], which is very effective for these fairly small values of c and d . We first started each method with zero as the initial approximate solution and allowed it to run for 40 iterations, after which the limit of residual norm reduction had been reached. Figure ?? shows plots of the logarithm of the Euclidean norm of the residual versus the number of iterations for the three methods. We note that in Fig. ?? and in all other figures below, the plotted residual norms were not the values maintained by , but rather were computed as accurately as possible “from scratch.” That is, at each iteration, the current approximate solution was formed and its product with the coefficient matrix was subtracted from the right-hand side, all in double precision. It was important to compute the residual norms in this way because the values maintained by become increasingly untrustworthy as the limits of residual norm reduction are neared; see [?]. It is seen in Fig. ?? that Algorithm EHA achieved the same ultimate level of residual norm reduction as the FDP method and required only a few more iterations to do so.

In our second set of experiments, we took $c = d = 100$ and carried out trials analogous to those in the first set above. No preconditioning was used in these experiments, both because we wanted to compare the methods without preconditioning and because the fast Poisson preconditioning used in the first set of experiments is not cost effective for these large values of c and d . We first allowed each method to run for 600 iterations, starting with zero as the initial approximate solution, after which the limit of residual norm reduction had been reached.

Acknowledgments. The author thanks the anonymous authors whose work largely constitutes this sample file. He also thanks the INFO-Tex mailing list for the valuable indirect assistance he received.

REFERENCES

- [1] VIVETTE GIRAULT AND PIERRE-ARNAUD RAVIART *Finite Element Method for Navier-Stokes Equations: Theory and Algorithms*, Springer-Verlag, Berlin, Germany, 1986.
- [2] FINITE ELEMENTS: THEORY, FAST SOLVERS, AND APPLICATIONS IN SOLID MECHANICS, Cambridge University Press, 2007.
- [3] R. A. BRUALDI AND B. L. SHADER, *On sign-nonsingular matrices and the conversion of the permanent into the determinant*, in Applied Geometry and Discrete Mathematics, The Victor Klee Festschrift, P. Gritzmann and B. Sturmfels, eds., American Mathematical Society, Providence, RI, 1991, pp. 117–134.
- [4] J. DREW, C. R. JOHNSON, AND P. VAN DEN DRIESSCHE, *Strong forms of nonsingularity*, Linear Algebra Appl., 162 (1992), to appear.
- [5] P. M. GIBSON, *Conversion of the permanent into the determinant*, Proc. Amer. Math. Soc., 27 (1971), pp. 471–476.
- [6] V. KLEE, R. LADNER, AND R. MANBER, *Signsolvability revisited*, Linear Algebra Appl., 59 (1984), pp. 131–157.
- [7] K. MUROTA, *LU-decomposition of a matrix with entries of different kinds*, Linear Algebra Appl., 49 (1983), pp. 275–283.
- [8] O. AXELSSON, *Conjugate gradient type methods for unsymmetric and inconsistent systems of linear equations*, Linear Algebra Appl., 29 (1980), pp. 1–16.
- [9] P. N. BROWN AND Y. SAAD, *Hybrid Krylov methods for nonlinear systems of equations*, SIAM J. Sci. Statist. Comput., 11 (1990), pp. 450–481.
- [10] R. S. DEMBO, S. C. EISENSTAT, AND T. STEIHAUG, *Inexact Newton methods*, SIAM J. Numer. Anal., 19 (1982), pp. 400–408.
- [11] S. C. EISENSTAT, H. C. ELMAN, AND M. H. SCHULTZ, *Variational iterative methods for non-symmetric systems of linear equations*, SIAM J. Numer. Anal., 20 (1983), pp. 345–357.
- [12] H. C. ELMAN, *Iterative methods for large, sparse, nonsymmetric systems of linear equations*, Ph.D. thesis, Department of Computer Science, Yale University, New Haven, CT, 1982.
- [13] R. GLOWINSKI, H. B. KELLER, AND L. RHEINHART, *Continuation-conjugate gradient methods for the least-squares solution of nonlinear boundary value problems*, SIAM J. Sci. Statist. Comput., 6 (1985), pp. 793–832.
- [14] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Second ed., The Johns Hopkins University Press, Baltimore, MD, 1989.
- [15] J. J. MORÉ, *A collection of nonlinear model problems*, in Computational Solutions of Nonlinear Systems of Equations, E. L. Allgower and K. Georg, eds., Lectures in Applied Mathematics, Vol. 26, American Mathematical Society, Providence, RI, 1990, pp. 723–762.
- [16] Y. SAAD, *Krylov subspace methods for solving large unsymmetric linear systems*, Math. Comp., 37 (1981), pp. 105–126.
- [17] Y. SAAD AND M. H. SCHULTZ, *GMRES: A generalized minimal residual method for solving nonsymmetric linear systems*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 856–869.
- [18] P. N. SWARZTRAUBER AND R. A. SWEET, *Efficient FORTRAN subprograms for the solution of elliptic partial differential equations*, ACM Trans. Math. Software, 5 (1979), pp. 352–364.
- [19] H. F. WALKER, *Implementation of the GMRES method using Householder transformations*, SIAM J. Sci. Statist. Comput., 9 (1988), pp. 152–163.
- [20] ———, *Implementations of the GMRES method*, Computer Phys. Comm., 53 (1989), pp. 311–320.