

## SAMPLE FILE FOR EQUADIFF 2017 L<sup>A</sup>T<sub>E</sub>X MACRO PACKAGE\*

PAUL DUGGAN<sup>†</sup> AND VARIOUS A. U. THORS<sup>‡</sup>

**Abstract.** An example of EQUADIFF 2017 L<sup>A</sup>T<sub>E</sub>X macros is presented.

**Key words.** sign-nonsingular matrix, LU-factorization, indicator polynomial

**AMS subject classifications.** 15A15, 15A09, 15A23

**1. Introduction and examples.** This paper presents a sample file for the use of EQUADIFF 2017's L<sup>A</sup>T<sub>E</sub>X macro package. This paper also serves as an example of EQUADIFF 2017'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.

**1.1. Sample text.** Let  $S = [s_{ij}]$  ( $1 \leq i, j \leq n$ ) be a  $(0, 1, -1)$ -matrix of order  $n$ . Then  $S$  is a *sign-nonsingular matrix* (SNS-matrix) provided that each real matrix with the same sign pattern as  $S$  is nonsingular. There has been considerable recent interest in constructing and characterizing SNS-matrices [1], [4]. There has also been interest in strong forms of sign-nonsingularity [2]. 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  $S$  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.

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

$$(1.1) \quad \int_a^b \left( \sum_i E_i B_{i,k,x}(t) \right) \left( \sum_j F_j B_{j,l,y}(t) \right) dt,$$

$$(1.2) \quad \int_a^b f(t) \left( \sum_i E_i B_{i,k,x}(t) \right) dt,$$

---

\*This work was supported by Grant No.: xxyy .

<sup>†</sup>Composition Department, Society for Industrial and Applied Mathematics, 3600 Univeristy City Science Center, Philadelphia, Pennsylvania, 19104-2688 ([duggan@siam.org](mailto:duggan@siam.org)).

<sup>‡</sup>Various Affiliations, supported by various foundation grants.

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 (1.1) is a special case of (1.2) where  $f(t)$  is a spline.

There are five different methods for calculating (1.1) 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 (1.2). The first four methods will be touched on and the last will be discussed at length.

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

$$(1.3) \quad \langle (A_1, B_1), (A_2, B_2) \rangle := \langle A_1, A_2 \rangle + \langle B_1, B_2 \rangle,$$

we calculate the Fréchet derivative of  $F$  as follows:

$$(1.4) \quad \begin{aligned} F'(U, V)(H, K) &= \langle R(U, V), H\Sigma V^T + U\Sigma K^T - P(H\Sigma V^T + U\Sigma K^T) \rangle \\ &= \langle R(U, V), H\Sigma V^T + U\Sigma K^T \rangle \\ &= \langle R(U, V)V\Sigma^T, H \rangle + \langle \Sigma^T U^T R(U, V), K^T \rangle. \end{aligned}$$

In the middle line of (1.4) we have used the fact that the range of  $R$  is always perpendicular to the range of  $P$ . The gradient  $\nabla F$  of  $F$ , therefore, may be interpreted as the pair of matrices:

$$(1.5) \quad \nabla F(U, V) = (R(U, V)V\Sigma^T, R(U, V)^T U\Sigma) \in R^{m \times m} \times R^{n \times n}.$$

Another array of equations

$$(1.6) \quad g(U, V) = \left( \frac{R(U, V)V\Sigma^T U^T - U\Sigma V^T R(U, V)^T}{2} U, \right. \\ \left. \frac{R(U, V)^T U\Sigma V^T - V\Sigma^T U^T R(U, V)}{2} V \right).$$

Thus, the vector field

$$(1.7) \quad \frac{d(U, V)}{dt} = -g(U, V)$$

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 (2.1) 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 2.1 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  $\psi$  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  $J_{m,p}$ .

DEFINITION 2.4. *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 [11], as well as by a number of other investigators; see [11] and the references therein. See also problem 3 by Glowinski and Keller and problem 7 by Mittelman in the collection of nonlinear model problems assembled by Moré [13]. The modified problem (2.4) has been used as a test problem for inexact Newton methods by Brown and Saad [7].

In our experiments, we took  $D = [0, 1] \times [0, 1]$ ,  $f \equiv 0$ ,  $c = d = 10$ , and discretized (2.4) using the usual second-order centered differences over a  $100 \times 100$  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.

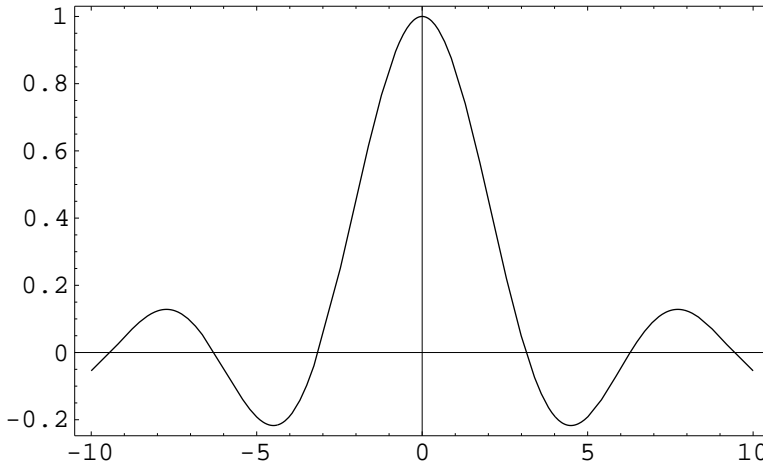


FIG. 2.1. Graph of the function  $\sin(x)/x$ .

In the first set of experiments, we allowed each method to run for 40 GMRES( $m$ ) iterations, starting with zero as the initial approximate solution, after which the limit

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  $\epsilon$ . For each method, the number of GMRES( $m$ ) iterations and  $F$ -evaluations was the same in every trial.

Method	$\epsilon$	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

of residual norm reduction had been reached. The results are shown in Fig. 2.1. In Fig. 2.1, 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 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 GMRES( $m$ ) 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] 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.
- [2] J. DREW, C. R. JOHNSON, AND P. VAN DEN DRIESSCHE, *Strong forms of nonsingularity*, Linear Algebra Appl., 162 (1992), to appear.
- [3] P. M. GIBSON, *Conversion of the permanent into the determinant*, Proc. Amer. Math. Soc., 27 (1971), pp. 471–476.
- [4] V. KLEE, R. LADNER, AND R. MANBER, *Signsolvability revisited*, Linear Algebra Appl., 59 (1984), pp. 131–157.
- [5] K. MUROTA, *LU-decomposition of a matrix with entries of different kinds*, Linear Algebra Appl., 49 (1983), pp. 275–283.
- [6] O. AXELSSON, *Conjugate gradient type methods for unsymmetric and inconsistent systems of linear equations*, Linear Algebra Appl., 29 (1980), pp. 1–16.
- [7] P. N. BROWN AND Y. SAAD, *Hybrid Krylov methods for nonlinear systems of equations*, SIAM J. Sci. Statist. Comput., 11 (1990), pp. 450–481.
- [8] R. S. DEMBO, S. C. EISENSTAT, AND T. STEIHAUG, *Inexact Newton methods*, SIAM J. Numer. Anal., 19 (1982), pp. 400–408.
- [9] 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.

- [10] 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.
- [11] 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.
- [12] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Second ed., The Johns Hopkins University Press, Baltimore, MD, 1989.
- [13] 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.
- [14] Y. SAAD, *Krylov subspace methods for solving large unsymmetric linear systems*, Math. Comp., 37 (1981), pp. 105–126.
- [15] 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.
- [16] 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.
- [17] H. F. WALKER, *Implementation of the GMRES method using Householder transformations*, SIAM J. Sci. Statist. Comput., 9 (1988), pp. 152–163.
- [18] ———, *Implementations of the GMRES method*, Computer Phys. Comm., 53 (1989), pp. 311–320.