

**MATRIX DIFFERENTIAL EQUATIONS:  
A CONTINUOUS REALIZATION PROCESS  
FOR LINEAR ALGEBRA PROBLEMS**

MOODY T. CHU \*

**Abstract.**

Many mathematical problems, such as existence questions, are studied by using an appropriate realization process, either iteratively or continuously. In this article differential equation techniques are used as a special continuous realization process for linear algebra problems. The matrix differential equations are cast in fairly general frameworks of which special cases have been found to be closely related to important numerical algorithms. The main thrust is to study the dynamics of various isospectral flows. This approach has potential applications ranging from new development of numerical algorithms to theoretical solution of open problems. Various aspects of the recent development and application in this direction are reviewed in this article.

**1. Introduction.**

Continuous realization methods are based on the idea of connecting two abstract problems through a mathematical bridge. Usually one of the abstract problems is a make-up whose solution is trivial while the other is the real problem whose solution is difficult to find. The bridge, if it exists, is regarded as a continuous path in the problem space. Following the path means deforming the underlying abstract problem mathematically. It is hoped that by following the path, the obvious solution will systematically be deformed into the solution that we are seeking for.

In applying a continuous realization method, two basic tasks should be carried out first since they are most accountable for the success:

1. One needs to establish a mathematical theory that can ensure the existence of bridge connecting the two abstract problems.
2. One needs to develop a numerical algorithm that can effectively follow the path.

The bridge usually takes the form as an integral curve of an ordinary differential equation describing how the problem data are transformed from the simple system to the more complicated system. The numerical algorithm thus should be an efficient ODE solver.

Depending upon how the bridge is constructed, continuous realization methods appear in different forms. One of the best known continuous realization methods in the literature perhaps is the so called homotopy method [1, 2, 34, 44, 60]. The philosophy behind the homotopy method is quite straightforward. We use the homotopy method to demonstrate the idea of continuation as follows: Suppose the original problem is to solve a nonlinear equation

$$(1) \quad f(x) = 0$$

where  $f : R^n \rightarrow R^n$  is a continuously differentiable function. We consider the homotopy function  $H : R^{n+1} \rightarrow R^n$  defined by

$$(2) \quad H(x, t) = f(x) - tf(x_0)$$

where  $x_0$  is an arbitrarily fixed point in  $R^n$ . Clearly, the solution to the equation  $H(x, 1) = 0$  is  $x_0$  and the solution to the equation  $H(x, 0) = 0$  is the same as that to (1). By tracing the solution of the equation  $H(x, t) = 0$  while the value of  $t$  is gradually changed from 1 to 0, we hope we will be led to a solution of (1). The idea is mathematically appealing because the zero set

---

\* department of Mathematics, North Carolina State University, Raleigh, North Carolina 27695-8205. This research was supported in part by National Science Foundation under grant DMS-9006135.

$H^{-1}(0) = \{(x, t) | H(x, t) = 0\}$  is indeed a smooth curve in  $R^{n+1}$ , provided  $0 \in R^n$  is a regular value for  $H$  [41, 56]. This smooth curve is what we refer to “the bridge”. Upon differentiation, we see that the homotopy curve is characterized by the initial value problem:

$$\begin{aligned}
 f'(x) \frac{dx}{ds} - \frac{1}{t} f(x) \frac{dt}{ds} &= 0 \\
 x(0) &= x_0 \\
 t(0) &= 1 \\
 \left(\frac{dx}{ds}\right)^2 + \left(\frac{dt}{ds}\right)^2 &= 1
 \end{aligned}
 \tag{3}$$

where  $s$  represents the arc length. The system (3) can be solved numerically by many available software packages.

What is not clear in (3) (and is crucial in all continuous realization methods) is whether the curve will ever reach the level  $t = 0$ . Properties of  $f$  and the selection of  $x_0$  must be taken into account in order that the bridge really makes the desired connection. Numerous applications together with their very own special homotopy functions have been studied in the literature. Far from being complete, we simply mention, for example, [13] for eigenvalue problems, [34] for nonlinear programming problem, [60] for physics applications and boundary value problems, [49, 50] for polynomial systems. Some special purpose curve tracing techniques can be found in [59] for HOMPACT, [48] for eigenvalue problems, [44] for pseudo-arc length technique, [53] for parameterized nonlinear systems and [50] for polynomial systems.

The idea of continuation can also be motivated through iterative methods. It is a well known fact that iterative methods have played very significant roles in a variety of ways for solving important mathematical problems. In our context, iterative methods may as well be regarded as discrete realization processes. It has been observed that many of the iterative methods used in numerical analysis may be regarded as the discrete realization of certain continuous dynamical systems. Suppose, for example, that the Jacobian matrix  $f'(x)$  in (3) is invertible. Then the differential equation may be written as

$$\frac{dx}{ds} = \frac{dt}{ds} \frac{1}{t} (f'(x))^{-1} f(x)
 \tag{4}$$

With the appropriate step size chosen, it is clear that one Euler step applied to the differential equation (4) is equivalent to one regular iteration of the classical Newton method applied to the equation (1) [44]. In a recent paper [12], we have already reviewed the development in the continuous realization of some of the basic iterative methods, including the Newton method, the QR algorithm, the RQI method, the SVD algorithm, the QZ algorithm and so on.

In this article, we want to address another facet of continuous realization methods. We shall examine how the continuation idea can be used to study linear algebra problems, including some of those where no numerical algorithm is available and some of those where the existence theory is yet to be settled. Our approaches described below are entirely different from those discussed in [12], yet the spirit of continuous realization is manifested clearly throughout the settings. In other words, unlike those discussed in [12], the dynamical systems to be considered are not constructed artificially from attempts to model existing iterative methods. Rather, the construction of bridges are monitored by the values of certain specified functions through which we hope certain desired properties will eventually be realized.

Most of the matrix differential equations discussed in this article are in the same basic form of

$$\frac{dX}{dt} = [X, k(X)]
 \tag{5}$$

where  $X = X(t) \in R^{n \times n}$ ,  $t \in R$ ,  $k$  stands for a certain matrix-valued operator acting on the matrix  $X$ , and  $[A, B] := AB - BA$  denotes the Lie bracket. Of particular importance is the case when  $k$  is skew-symmetric. Under appropriately formulated  $k$ , the equation (5) has found applications in eigenvalue or singular value problems, spectrally or singular-value constrained least squares approximation problems, inverse Toeplitz or non-negative eigenvalue problems, nearest normal matrix problem, quadratic programming problems, and simultaneous reduction problems. This paper summarizes some of the recent contributions in these aspects.

In addition, we think the approach by continuous realization methods might have the following advantages:

- There are many well-developed classical results for continuous dynamical systems. The study of continuous system might shed critical insights into the understanding of the dynamics of the corresponding discrete methods.
- In contrast to the local properties for some discrete methods, the continuous approach usually offers a global method for solving the underlying problem.
- Some existence problems, seemingly impossible to be tackled by any conventional discrete methods, may be solved by formulating a special differential equation that ensure a specific task is taking place.
- Continuous realization sometimes unifies different discrete methods as special cases of its discretization and often gives rise to the design of new numerical algorithms.
- In a sense a continuation process means a spontaneous evolution of a problem. Thus it has potential applications in numerical analysis, control theory, signal processing, matrix theory, and mathematical programming. This opens up a new direction of applying numerical ODE techniques, although the differential systems resulted from continuous realization present immediate challenge to most current ODE methods. In particular, we would like to have an ODE solver that can effectively approximate the asymptotically stable attractors of a dynamical system [45]. Also, matrix differential equations are especially suitable for integration on a massively data-parallel computing system, such as the Connection Machine. Thus matrix differential equations may be used as large-scale benchmark problems for testing parallel ODE techniques [20]. Conversely, parallel ODE techniques may benefit the numerical solution of matrix differential equations.

We shall see that the solution  $X(t)$  of (5) represents, in a sense, the evolution of  $X(0)$  under a smooth change of coordinate systems [3]. From this viewpoint, we think the ideas conveyed in this paper should be of interest to readers from different research fields.

## 2. Isospectral Flows.

In this section we consider some general properties of the differential equation (5).

Let  $\mathcal{G}(n)$  denote the Lie group [23, 41, 56] of all nonsingular matrices in  $R^{n \times n}$ . Associated with any given matrix  $X_0 \in R^{n \times n}$ , we define an isospectral surface

$$(6) \quad \mathcal{M}(X_0) := \{Z^{-1}XZ \mid Z \in \mathcal{G}(n)\}.$$

We note that every element in  $\mathcal{M}(X_0)$  is similar to  $X_0$ . Suppose  $Z(t)$ , with  $Z(0) = 1$ , represents a differentiable curve on the manifold  $\mathcal{G}(n)$ . Then

$$(7) \quad X(t) := Z(t)^{-1}X_0Z(t)$$

defines a differentiable curve, with  $X(0) = X_0$ , on the surface  $\mathcal{M}(X_0)$ . Upon differentiation, it is easy to see the curve  $X(t)$  is the solution of the initial value problem

$$(8) \quad \begin{aligned} \frac{dX(t)}{dt} &= [X(t), k(t)] \\ X(0) &= X_0 \end{aligned}$$

with  $k(t)$  defined by

$$(9) \quad k(t) := Z(t)^{-1} \frac{dZ(t)}{dt}.$$

Conversely, given any 1-parameter family of matrices  $k(t)$  in  $R^{n \times n}$ , it can be proved that the solution of (8) can be written in the form (7) where  $Z(t)$  satisfies

$$(10) \quad \begin{aligned} \frac{dZ(t)}{dt} &= Z(t)k(t) \\ Z(0) &= I. \end{aligned}$$

Henceforth, we shall call the system (10) the *dual problem* of (8).

With different choices of  $k(t)$ , the differential equation in (8) defines different isospectral curves, all emanating from the same initial value  $X_0$ . Obviously, the asymptotic behavior of  $X(t)$ , if there is any, on the surface  $\mathcal{M}(X_0)$  will be determined by that of the corresponding  $Z(t)$  on the manifold  $\mathcal{G}(n)$ , and vice versa.

A special case of the above discussion is particularly important. We may replace the group  $\mathcal{G}(n)$  by the subgroup  $\mathcal{O}(n)$  of all orthogonal matrices and repeat the above argument. Without causing ambiguity, we shall use the same notation  $\mathcal{M}(X_0)$  to represent

$$(11) \quad \mathcal{M}(X_0) = \{Q^T X_0 Q \mid Q \in \mathcal{O}(n)\}.$$

We note that in the equation (10),  $Q(t) \in \mathcal{O}(n)$  if and only if  $k(t)$  is skew-symmetric. We note also that  $\|X(t)\|_2 = \|X_0\|_2$  so long as  $k(t)$  is defined and remains to be skew-symmetric. It is important to recognize that the tangent space of the manifold  $\mathcal{O}(n)$  at any orthogonal matrix  $Q$  is given by [23, 41, 56]

$$(12) \quad T_Q \mathcal{O}(n) = Q\mathcal{S}(n)^\perp$$

where  $\mathcal{S}(n)$  is the subspace of all symmetric matrices in  $R^{n \times n}$  and  $\mathcal{S}(n)^\perp$  is the orthogonal complement of  $\mathcal{S}(n)$  under the Frobenius inner product

$$(13) \quad \langle A, B \rangle := \text{trace}(AB^T) = \sum_{i,j} a_{ij}b_{ij}.$$

Since  $\mathcal{S}(n)^\perp$  is the subspace of all skew-symmetric matrices, the dual problem (10) indeed defines a flow on the manifold  $\mathcal{O}(n)$ .

The differential equation (5) is another important special case of (8) in which  $k(t) = k(X(t))$ . That is, (5) is an autonomous system. If, in addition,  $k(t)$  is also skew-symmetric, then  $X(t)$  exists and is bounded for all  $t \in (-\infty, \infty)$ . Also, because of the relationship (7), the curve  $Q(t)$  on the manifold  $\mathcal{O}(n)$  is self-determined by the dual initial value problem

$$(14) \quad \begin{aligned} \frac{dQ(t)}{dt} &= Q(t)k(Q(t)^T X_0 Q(t)) \\ Q(0) &= I \end{aligned}$$

which is independent of  $X(t)$ .

When  $k$  is skew-symmetric, the dual system (10) has an interesting physical interpretation in terms of rigid body motions [37]. We demonstrate the simple case when  $n = 3$ . Suppose we write the operator  $k(t)$  as

$$(15) \quad k(t) = \begin{bmatrix} 0 & \omega_3(t) & -\omega_2(t) \\ -\omega_3(t) & 0 & \omega_1(t) \\ \omega_2(t) & -\omega_1(t) & 0 \end{bmatrix}$$

and define

$$(16) \quad \omega(t) := [\omega_1(t), \omega_2(t), \omega_3(t)]^T.$$

Also we write the transpose of  $Q(t)$  in columns

$$(17) \quad Q(t)^T = [p_1(t), p_2(t), p_3(t)].$$

Then the dual system (10) implies, for  $i = 1, 2, 3$ ,

$$(18) \quad \begin{aligned} \frac{dp_i}{dt} &= \omega \times p_i \\ p_i(0) &= e_i \end{aligned}$$

where  $e_i$  is the  $i$ -th standard unit basis vector in  $R^3$  and  $\times$  denotes the usual cross-product in  $R^3$ . The differential equation in (18) describes the linear velocity of the vector  $p_i(t)$  and the vector  $\omega(t)$  may be interpreted as the angular velocity of the motion. As a whole, the dual system then describes the rotation of an orthogonal coordinate system about the origin. If the initial matrix  $X_0$  is symmetric and positive definite, then  $X(t)$  stays to be symmetric and positive definite. We may thus interpret  $X(t)$  as the moment of inertia tensor of a rigid body motion. In this case, the total kinetic energy of motion is given by [37]

$$(19) \quad T(t) := \frac{1}{2} \omega(t)^T X(t) \omega(t).$$

If the kinetic energy is dissipated to zero, then the motion eventually stops.

The idea of rigid body motion can be generalized to higher  $n$ . The differential system (10) with skew-symmetric  $k$ , therefore, represents a smooth change of coordinate systems and the solution  $X(t)$  of (8) is simply the same transformation as  $X(0)$  expressed in different coordinate systems.

### 3. QR-type Algorithms.

It is a well known fact that every matrix in  $R^{n \times n}$  can be factored as the product of an orthogonal matrix and an upper triangular matrix [38, 43]. Such a decomposition is called the QR decomposition. Based on the QR decomposition, the unshifted QR algorithm generates a sequence of matrices  $\{A_k\}$  according to the scheme:

$$(20) \quad A_k = Q_k R_k \implies A_{k+1} := R_k Q_k$$

where  $Q_k R_k$  is the QR decomposition of  $A_k$ . This algorithm is very important in the computation of eigenvalues for  $A_0$  [61]. We also consider the matrix equation (known as the Toda lattice):

$$(21) \quad \begin{aligned} \frac{dX}{dt} &= [X, \Pi_0(X)] \\ X(0) &= X_0 \end{aligned}$$

where  $\Pi_0(X) := (X^-) - (X^-)^T$ ,  $X^-$  is the strictly lower triangular part of  $X$ . Recently it has been shown [55] that the sequence  $\{\exp(X(k))\}$  obtained by sampling  $X(t)$  at integer times corresponds exactly to the sequence  $\{A_k\}$  generated by (20) if  $A_0 = \exp(X_0)$ . Therefore, the convergence properties of the QR algorithm may be understood by studying the dynamics of (21) [16, 24, 57].

In this section we shall show that the system (21) is a special case of a more general setting from which more iterative processes may arise.

Suppose the space  $R^{n \times n}$  is split as the direct sum of two subspaces  $\mathcal{V}_1$  and  $\mathcal{V}_2$ . Let  $P_1$  and  $P_2$  represent the natural projection mappings from  $R^{n \times n}$  into  $\mathcal{V}_1$  and  $\mathcal{V}_2$ , respectively. Consider the initial value problem:

$$(22) \quad \begin{aligned} \frac{dX}{dt} &= [X, P_1(X)] \\ X(0) &= X_0. \end{aligned}$$

Since  $[X(t), X(t)] = 0$ , the solution of (21) also satisfies the initial value problem:

$$(23) \quad \begin{aligned} \frac{dX}{dt} &= [P_2(X), X] \\ X(0) &= X_0. \end{aligned}$$

Correspondingly, from the discussion in the preceding section, there exist dual problems

$$(24) \quad \begin{aligned} \frac{dZ_1}{dt} &= Z_1 P_1(X) \\ Z_1(0) &= I \end{aligned}$$

and

$$(25) \quad \begin{aligned} \frac{dZ_2}{dt} &= P_2(X) Z_2 \\ Z_2(0) &= I \end{aligned}$$

whose solutions satisfy

$$(26) \quad X(t) = Z_1(t)^{-1} X_0 Z_1(t) = Z_2(t) X_0 Z_2(t)^{-1}.$$

Furthermore, the matrices  $Z_1(t)$  and  $Z_2(t)$  are related in a special way [14]:

**THEOREM 3.1.** *Suppose  $X(t)$ ,  $Z_1(t)$  and  $Z_2(t)$  exist on the interval  $[0, T]$ . Then*

$$(27) \quad \exp(tX_0) = Z_1(t) Z_2(t)$$

$$(28) \quad \exp(tX(t)) = Z_2(t) Z_1(t)$$

for all  $t \in [0, T]$ .

If we set  $t = 1$  in Theorem 3.1, then it follows that

$$(29) \quad \exp(X(0)) = Z_1(1)Z_2(1)$$

$$(30) \quad \exp(X(1)) = Z_2(1)Z_1(1).$$

Since the differential equation (22) is autonomous, it follows that the “swapping multiplication” relationship observed in (29) and (30) will hold at every feasible integer time. In this sense, the dual systems (24) and (25) give rise to an abstract matrix factorization (of the matrix  $\exp(X(t))$ ) and the isospectral flow (22) gives rise to a specific QR-like iterative process.

The Toda lattice corresponds to the special splitting of  $R^{n \times n}$  where  $\mathcal{V}_1$  is the subspace of all skew-symmetric matrices and  $\mathcal{V}_2$  is the subspace of all upper triangular matrices. Such a splitting is a Lie algebra decomposition of  $R^{n \times n}$  and therefore there corresponds a Lie group decomposition of  $\mathcal{G}(n)$  which is the QR decomposition. Other kinds of Lie algebra decompositions and the associated isospectral flows are discussed in [25, 51, 52, 57, 58].

Our approach has the advantage that only subspace decomposition of  $R^{n \times n}$  is involved, which should be much easier to manipulate than Lie algebra decomposition [58]. Nonetheless, we have proved that the time-1 mapping of the solution  $X(t)$  of (22) still enjoys a QR-type algorithm. With special choices of the subspaces (and, hence, the projections), We are able to unify several matrix decomposition algorithms and give rise to a number of new decompositions.

As an example, the following theorem generalizes the well-known Schur theorem [43] in the sense that one can zero out *any* pattern of off-diagonal elements of a symmetric matrix by orthogonal similarity transformation [14]:

**THEOREM 3.2.** *Suppose  $X_0$  is symmetric. Let  $\Delta$  be an arbitrary subset of the index set  $\{(i, j) | 1 \leq j < i \leq n\}$ . For each  $X \in R^{n \times n}$ , define  $\hat{X}$  so that*

$$(31) \quad \hat{x}_{ij} := \begin{cases} x_{ij} & \text{if } (i, j) \in \Delta \\ 0 & \text{otherwise} \end{cases}$$

and define

$$(32) \quad P_1(X) := \hat{X} - \hat{X}^T.$$

Then the solution  $X(t)$  to the system (22) with  $P_1$  defined by (32) is defined for all  $t$ , stays to be symmetric, and converges to a limit point as  $t \rightarrow \infty$ . Furthermore,  $x_{ij}(t) \rightarrow 0$  if  $(i, j) \in \Delta$ .

As another example, the following theorem suggests how to change a non-symmetric matrix into a stair-case matrix by using orthogonal similarity transformations [14].

**THEOREM 3.3.** *Suppose  $X_0$  has only simple eigenvalues. Let  $\Delta$  be the collection of indices corresponding to any block strictly lower triangular matrix. For any  $X \in R^{n \times n}$ , define  $\hat{X}$  and  $P_1(X)$  as in (31) and (32). Then the same conclusion as in Theorem 3.2 holds.*

A special application of Theorem 3.3 is worth mentioning. We conjecture that the Hamiltonian eigenvalue problem should enjoy a QR-type algorithm that preserves the Hamiltonian structure [8, 10]. A matrix  $X \in R^{2n \times 2n}$  is Hamiltonian if and only if it is of the form

$$(33) \quad X = \begin{bmatrix} A & N \\ K & -A^T \end{bmatrix}$$

where  $K, N \in R^{n \times n}$  are symmetric. Suppose we define

$$(34) \quad P_1(X) = \begin{bmatrix} 0 & -K \\ K & 0 \end{bmatrix}.$$

Then it is easy to see that  $[X, P_1(X)]$  is also Hamiltonian. With this in mind, starting with a Hamiltonian matrix  $X_0$ , we define the dynamical system (22) with  $P_1$  given by (34). Then the solution  $Z(t)$  to the corresponding dual problem (24) is both orthogonal and symplectic, and the solution  $X(t) = Z_1(t)^T X_0 Z_1(t)$  to (22) remains Hamiltonian for all  $t$ . Furthermore, by Theorem 3.3, we know that  $P_1(X) \rightarrow 0$  as  $t \rightarrow \infty$ . Unfortunately, the associated iterative algorithm is not explicitly known due to lacking knowledge of the structure of the corresponding  $Z_2(t)$ . As a matter of fact, to find an efficient numerically stable iterative process that also preserves the Hamiltonian structure is still an open problem. See Bunse-Gerstner et al. [9] and the many references contained therein.

A similar approach that considers the linkage between QR-type algorithms and solutions to the Yang-Baxter equation can be found in [47, 30]. In particular, the relationship between an explicit iterative scheme associated with a group decomposition of the symplectic group and a dynamical system of the form (22) that preserves the Hamiltonian structure is established. Although the iterative scheme is explicit and is proved to be convergent, it does not involve any orthogonal transformation and, hence, is possibly unstable.



#### 4. Projected Gradient Flows.

In this section we shall restrict our attention to the subspace  $\mathcal{S}(n)$ . Let  $\mathcal{V}$  be either a single matrix  $V$  or a subspace in  $\mathcal{S}(n)$ . For any  $X \in \mathcal{S}(n)$ , the projection of  $X$  into  $\mathcal{V}$  is denoted as  $P(X)$ . If  $\mathcal{V}$  is a single matrix, then  $P(X) \equiv V$ ; otherwise, the projection is taken with respect to the Frobenius inner product. We shall consider another special case of (5), i.e., we shall consider equation of the form:

$$(35) \quad \begin{aligned} \frac{dX}{dt} &= [X, [X, P(X)]] \\ X(0) &= X_0 \end{aligned}$$

with  $X_0 \in \mathcal{S}(n)$ .

We first note that the system (35) evolves in the space  $\mathcal{S}(n)$  since  $X \in \mathcal{S}(n)$  implies

$$(36) \quad k(X) := [X, P(X)]$$

is skew-symmetric and  $\frac{dX}{dt} \in \mathcal{S}(n)$ , and vice versa. It is also clear that  $X(t)$  is orthogonally similar to  $X_0$ .

The system (35) is derived from the following minimization problem:

$$(37) \quad \begin{aligned} \text{Minimize} \quad & F(X) := \frac{1}{2} \|X - P(X)\|^2 \\ \text{Subject to} \quad & X \in \mathcal{M}(X_0) \end{aligned}$$

where  $\|\cdot\|$  stands for the Frobenius matrix norm. Literally, problem (37) is minimizing the distance [36] between the two sets  $\mathcal{M}(X_0)$  and  $\mathcal{V}$ . From this prospect, we may thus use (35) to continuously realize the solution to, for example, the following linear algebra problems:

**(Problem A)** Given a real symmetric matrix  $N$ , find a least squares approximation of  $N$  that is still symmetric but has a prescribed set of eigenvalues  $\{\lambda_1, \dots, \lambda_n\}$ . In this setting, we choose  $V \equiv N$  and  $X_0 = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ .

**(Problem B)** Given a set of real numbers  $\{\lambda_1, \dots, \lambda_n\}$ , construct a symmetric Toeplitz matrix that has the prescribed set as its spectrum. In this setting, we choose  $\mathcal{V}$  to be the subspace  $\mathcal{T}$  of all symmetric Toeplitz matrices and  $X_0$  to be any matrix orthogonally similar to  $\text{diag}\{\lambda_1, \dots, \lambda_n\}$ .

**(Problem C)** Given a matrix  $A \in \mathcal{S}(n)$ , find its eigenvalues. In this setting, we choose  $\mathcal{V}$  to be the subspace of all diagonal matrices and  $X_0 = A$ .

Problem (37) is equivalent to

$$(38) \quad \begin{aligned} \text{Minimize} \quad & G(Q) := \frac{1}{2} \langle Q^T X_0 Q - P(Q^T X_0 Q), Q^T X_0 Q - P(Q^T X_0 Q) \rangle \\ \text{Subject to} \quad & Q \in \mathcal{O}(n). \end{aligned}$$

It is easy to see that with respect to the Frobenius inner product, the gradient of  $G$  at  $Q \in \mathcal{O}(n)$  is the matrix

$$(39) \quad \nabla G(Q) = 2X_0 Q(Q^T X_0 Q - P(Q^T X_0 Q)).$$

We observe that

$$(40) \quad R^{n \times n} = T_Q \mathcal{O}(n) \oplus T_Q \mathcal{O}(n)^\perp = Q\mathcal{S}(n)^\perp \oplus Q\mathcal{S}(n).$$

Therefore, the projection  $g(Q)$  of  $\nabla G(Q)$  onto the manifold  $\mathcal{O}(n)$  can be calculated explicitly [15]:

$$\begin{aligned} g(Q) &= Q \left\{ \frac{1}{2} (Q^T \nabla G(Q) - \nabla G(Q)^T Q) \right\} \\ (41) \quad &= Q [P(Q^T X_0 Q), Q^T X_0 Q]. \end{aligned}$$

It is now obvious that the dual problem of (35):

$$\begin{aligned} \frac{dQ}{dt} &= Qk(Q^T X_0 Q) \\ (42) \quad Q(0) &= I \end{aligned}$$

with  $k$  defined by (36) signifies a steepest descent flow on  $\mathcal{O}(n)$  for problem (38). Equivalently, (35) defined a descent flow on  $\mathcal{M}(X_0)$  for problem (37).

Because (35) (or (42)) is a gradient flow, the function  $F(X(t))$  (or  $G(Q(t))$ ) is a natural Lyapunov function. The Lyapunov function can be used to characterize the dynamics of the flow. By using a Lyapunov function, it is possible to derive an effective, strictly stable multistep method for approximating the  $\omega$ -limit set. See, for example, [29], and [45]. This direction certainly is worth further investigation.

The system (35) applied to Problem A becomes

$$\begin{aligned} \frac{dX}{dt} &= [X, [X, N]] \\ (43) \quad X(0) &= X_0. \end{aligned}$$

The dynamics of (43) has been studied independently by Chu [15, 22] and Brockett [6, 7]. The main results are as follows:

**THEOREM 4.1.** *Suppose both  $N$  and  $X_0$  have distinct eigenvalues. Let the eigenvalues of  $N$  and  $X_0$  be ordered as  $\mu_1 < \dots < \mu_n$  and  $\lambda_1 < \dots < \lambda_n$ , respectively. Then as  $t \rightarrow \infty$ , the solution  $X(t)$  of (43) converges to the unique limit point*

$$(44) \quad \hat{X} = \lambda_1 q_1 q_1^T + \dots + \lambda_n q_n q_n^T$$

where  $q_1, \dots, q_n$  are the normalized eigenvectors of  $N$  corresponding respectively to  $\mu_1, \dots, \mu_n$ .

**THEOREM 4.2.** *Suppose  $N$  is a diagonal matrix with distinct eigenvalues (but  $X_0$  may have repeated eigenvalues). Then as  $t \rightarrow \infty$ , the solution  $X(t)$  of (43) converges to a diagonal matrix whose elements are similarly ordered as those in  $N$ .*

At the first glance, it is quite amazing that the  $\omega$ -limit point of (43) (and, hence, the solution to Problem A) can be expressed explicitly. After further consideration, we find that Theorem 4.1 may be regarded as a reproof of the well known Wielandt-Hoffman Theorem [43, 15]. In fact, we have proved that the bound in Wielandt-Hoffman Theorem is sharp.

Recently it is also observed [6] and [7] that, due to the special ‘‘sorting property’’ of (43), the continuous realization idea can even be applied to solve data matching problem and a variety of generic combinatorial optimizations. We demonstrate one application to the linear programming problem. Consider the LP problem:

$$\begin{aligned} \text{Maximize} \quad & c^T x \\ (45) \quad \text{Subject to} \quad & Ax \leq b \end{aligned}$$

where  $c, x \in R^n$  and we shall assume the feasible set  $\mathcal{P} := \{x | x \in R^n, Ax \leq b\}$  is a convex polytope with vertices at  $a_1, \dots, a_p \in R^n$ . It is a well known fact that one of the values  $\mu_i := c^T a_i$ ,  $i = 1, \dots, p$  will be the optimum for problem (45). To sort out this particular vertex, we define  $T \in R^{n \times p}$  as

$T := [a_1, \dots, a_p]$ . Obviously  $T$  maps the standard simplex  $\mathcal{S} := \{d \mid d \in \mathbb{R}^p, d_i \geq 0, \sum_{i=1}^p d_i = 1\}$  onto the given polytope  $\mathcal{P}$ . Let  $X_0 := \text{diag}\{1, 0, \dots, 0\} \in \mathbb{R}^{p \times p}$  and let  $N := \text{diag}\{\mu_1, \dots, \mu_p\}$ . Then, by Theorem 4.2, the corresponding isospectral flow  $X(t)$  converges to a diagonal matrix  $\hat{X}$ . The elements of  $\hat{X}$  must be a permutation of those of  $X_0$  and must be arranged in an order similar to that of  $N$ . By identifying the index corresponding to the value 1 in  $\hat{X}$ , we locate the optimal vertex.

In applying the system (35) to Problem B, the projection  $P$  can easily be calculated from

$$(46) \quad P(X) = \sum_{k=1}^n \langle X, E_k \rangle E_k$$

where  $\{E_1, \dots, E_n\}$  is an orthogonal basis of  $\mathcal{T}$ . Denoting  $E_k$  by  $(e_{ij}^{(k)})$ , then one natural basis is

$$(47) \quad e_{ij}^{(k)} = \begin{cases} 1/\sqrt{2(n-k+1)} & \text{if } 1 < k \leq n \text{ and } |i-j| = k-1 \\ 1/\sqrt{n} & \text{if } k=1 \text{ and } i=j \\ 0 & \text{otherwise.} \end{cases}$$

To our knowledge, the existence question of a solution to the inverse Toeplitz eigenvalue problem when  $n \geq 5$  is still an open problem [46, 26]. Yet our descent flow formulation offers a numerical method for computing the solution. We should point out, however, that the differential system (35) applied to Problem B may have asymptotically stable  $\omega$ -limit points other than those that are in  $\mathcal{T}$ . If this happens, we should change to a different course of integration by starting with a different initial value. Another approach to circumvent this difficult will be discussed in a later section.

The system (35) applied to Problem C becomes

$$(48) \quad \begin{aligned} \frac{dX}{dt} &= [X, [X, \text{diag}(X)]] \\ X(0) &= A \end{aligned}$$

where  $\text{diag}(X)$  denotes the diagonal matrix whose elements are those along the diagonal of  $X$ . Recall that the system (48) is derived by minimizing the sum of squares of the off-diagonal elements of  $X$ . Thus the system (48) may be regarded as a continuous analog of the classical Jacobi method for the eigenvalue problem [38]. Although the differential equation itself has many equilibrium points, most of them are not stable. In fact, it can be proved [27] that

**THEOREM 4.3.** *As  $t \rightarrow \infty$ , the solution  $X(t)$  of (48) converges to a diagonal matrix. Thus, if  $A$  has distinct eigenvalues, then on the manifold  $\mathcal{M}(A)$  there are exactly  $n!$  asymptotically stable equilibrium points.*

We have just discussed three possible applications of (35). Indeed, the liberty of choosing the projected space  $\mathcal{V}$  suggests that the system (35) might have broader applications in other areas. It is interesting to see that for different choices of  $\mathcal{V}$ , the system (35) generates different descent flows evolving on the manifold  $\mathcal{M}(X_0)$ . In the next section, we shall apply the same idea to consider the simultaneous reduction problems.

## 5. Simultaneous Reduction Problems.

Undoubtedly, orthogonal similarity transformations  $Q^T A Q$  and orthogonal equivalence transformations  $Q^T A Z$  with  $Q, Z \in \mathcal{O}(n)$  play important roles in the computation of eigenvalues and singular values, respectively, for a general matrix  $A \in R^{n \times n}$  [38, 61]. These transformations reduce  $A$  to some simpler forms. Simultaneous reduction of more than one matrix, on the other hand, finds applications in areas like control theory or mechanics for the identification of multivariable systems or for the study of small oscillations about a stable equilibrium [43]. A few theoretical results concerning the simultaneous diagonalization of two symmetric (or hermitian) matrices can be found in [43]. The reduction problem for more than two or for more general matrices are considerably more difficult both in theory and in computation. In this section we shall discuss another special case of (5) that induces an easy but versatile reduction procedure.

Suppose  $A_i \in R^{n \times n}$ ,  $i = 1, \dots, p$  are the matrices under consideration to be reduced by orthogonal similar transformations. For each  $i$ , let  $\mathcal{V}_i$  denote the subspace of all matrices having the specified form to which  $A_i$  is supposed to be reduced. These subspaces need not to be the same. Given any  $A \in R^{n \times n}$ , let  $P_i(A)$  denote the projection of  $A$  into the subspace  $\mathcal{V}_i$  with respect to the Frobenius inner product. For convenience, we also define the residual operator

$$(49) \quad \alpha_i(Q) := Q^T A_i Q - P_i(Q^T A_i Q).$$

We shall use the same idea as in problem (38) to solve the problem:

$$(50) \quad \begin{aligned} \text{Minimize} \quad & F(Q) := \frac{1}{2} \sum_{i=1}^p \|\alpha_i(Q)\|^2 \\ \text{Subject to} \quad & Q \in \mathcal{O}(n). \end{aligned}$$

That is, while moving along the orthogonal similarity orbit of the given matrices  $A_1, \dots, A_p$ , we want to minimize the total distance between the point  $Q^T A_i Q$  and the subspace  $\mathcal{V}_i$  for all  $i$ .

The gradient of  $F$  is given by [19]

$$(51) \quad \nabla F(Q) = \sum_{i=1}^p (A_i^T Q \alpha_i(Q) + A_i Q \alpha_i^T(Q)).$$

By using (40), the projection of  $\nabla F(Q)$  onto the manifold  $\mathcal{O}(n)$  is calculated to be

$$(52) \quad g(Q) = Q \sum_{i=1}^p \frac{[Q^T A_i Q, \alpha_i^T(Q)] - [Q^T A_i Q, \alpha_i^T(Q)]^T}{2}.$$

Therefore, the differential equation

$$(53) \quad \frac{dQ}{dt} = Q k(Q)$$

with

$$(54) \quad k(Q) := - \sum_{i=1}^p \frac{[Q^T A_i Q, \alpha_i^T(Q)] - [Q^T A_i Q, \alpha_i^T(Q)]^T}{2}.$$

defines a steepest descent vector field on the manifold  $\mathcal{O}(n)$  for problem (50). Note that (53) is again in the form of (10). Together with the initial value  $Q(0) = I$ , the equation (53) is the dual problem of the system

$$(55) \quad \begin{aligned} \frac{dX_i}{dt} &= [X_i, k(Q)] \\ &= \left[ X_i, \sum_{j=1}^p \frac{[X_j, P_j^T(X_j)] - [X_j, P_j^T(X_j)]^T}{2} \right] \\ X_i(0) &= A_i \end{aligned}$$

where  $X_i(t) := Q(t)^T A_i Q(t)$ , and  $i = 1, \dots, p$ .

The solution  $\{X_1(t), \dots, X_p(t)\}$  of (55) represents a continuous evolution from  $\{A_1, \dots, A_p\}$  in realizing the prescribed reduction forms while members are influenced by each others in such a way that the total distance  $F(Q)$  is monotonically decreased. We note that such a continuous realization process has the advantages that the desired form to which matrices are reduced can be almost arbitrary, and that if a desired form is not attainable then the limit point of the differential equation gives a way of measuring the distance from the best reduced matrices to the nearest matrices that have the desired form.

Similarly, suppose  $A_i \in R^{m \times n}$ ,  $i = 1, \dots, p$  are the matrices being considered by using orthogonal equivalence transformations. Define

$$(56) \quad \beta_i(Q, Z) := Q^T A_i Z - P_i(Q^T A_i Z)$$

and consider the problem

$$(57) \quad \begin{aligned} \text{Minimize} \quad & G(Q, Z) := \frac{1}{2} \sum_{i=1}^p \|\beta_i(Q, Z)\|^2 \\ \text{Subject to} \quad & Q \in \mathcal{O}(m) \\ & Z \in \mathcal{O}(n). \end{aligned}$$

By introducing the product topology on  $R^{m \times m} \times R^{n \times n}$  where the induced inner product is defined by

$$(58) \quad \langle (A_1, B_1), (A_2, B_2) \rangle := \langle A_1, A_2 \rangle + \langle B_1, B_2 \rangle,$$

we know that [19]

$$(59) \quad \nabla G(Q, Z) = \left( \sum_{i=1}^p A_i Z \beta_i^T(Q, Z), \sum_{i=1}^p A_i^T Q \beta_i(Q, Z) \right).$$

Define

$$(60) \quad k_1(Q, Z) := - \sum_{i=1}^p \frac{Q^T A_i Z \beta_i^T(Q, Z) - (Q^T A_i Z \beta_i^T(Q, Z))^T}{2},$$

$$(61) \quad k_2(Q, Z) := - \sum_{i=1}^p \frac{Z^T A_i^T Q \beta_i(Q, Z) - (Z^T A_i^T Q \beta_i(Q, Z))^T}{2}.$$

(62)

Note  $k_1 \in R^{m \times m}$  and  $k_2 \in R^{n \times n}$  are both skew-symmetric. Since

$$(63) \quad T_{(Q, Z)} \mathcal{O}(m) \times \mathcal{O}(n) = QS(m)^\perp \times ZS(n)^\perp,$$

using the same principle as in (40), we find the projection of  $\nabla G(Q, Z)$  onto the manifold  $\mathcal{O}(m) \times \mathcal{O}(n)$  is given by

$$(64) \quad g(Q, Z) = (-Qk_1(Q, Z), -Zk_2(Q, Z)).$$

Thus the system

$$(65) \quad \begin{aligned} \frac{d(Q, Z)}{dt} &= (Qk_1(Q, Z), Zk_2(Q, Z)) \\ (Q(0), Z(0)) &= (I_m, I_n) \end{aligned}$$

defines a steepest descent flow on  $\mathcal{O}(m) \times \mathcal{O}(n)$  and is the dual problem of the system

$$\begin{aligned}
\frac{dX_i}{dt} &= X_i k_2(Q, Z) - k_1(Q, Z) X_i \\
&= \sum_{j=1}^p \left\{ X_i \frac{X_j^T P_j(X_j) - P_j^T(X_j) X_j}{2} + \frac{P_j(X_j) X_j^T - X_j P_j^T(X_j)}{2} X_i \right\} \\
(66) \quad X_i(0) &= A_i
\end{aligned}$$

where  $X_i(t) := Q(t)^T A_i Z(t)$ , and  $i = 1, \dots, p$ . We note that (66) is not quite in the same form as (5).

In contrast to the approach of deriving (66), it is worthwhile to mention two other matrix differential systems:

$$\begin{aligned}
\frac{dX}{dt} &= X \Pi_0(X X^T) - \Pi_0(X^T X) X \\
(67) \quad X(0) &= A
\end{aligned}$$

and

$$\begin{aligned}
\frac{dX_1}{dt} &= X_1 \Pi_0(X_2^{-1} X_1) - \Pi_0(X_1 X_2^{-1}) X_1 \\
\frac{dX_2}{dt} &= X_2 \Pi_0(X_2^{-1} X_1) - \Pi_0(X_1 X_2^{-1}) X_2 \\
X_1 &= A_1 \\
(68) \quad X_2 &= A_2.
\end{aligned}$$

We have proved earlier that, just as the Toda lattice (21) models the QR algorithm, the system (67) models the SVD algorithm [18] for the  $A \in R^{m \times n}$ , and (68) models the QZ algorithm [17] for the matrix pencil  $(A_1, A_2) \in R^{n \times n} \times R^{n \times n}$ . Although (67) and (68) are quite similar to (66), the way they are derived has nothing to do with optimization of any objective function. On the other hand, suppose we take  $p = 1$  and  $\mathcal{V}$  to be the subspace of all diagonal matrices in  $R^{m \times n}$ . Then the differential equation in (66) becomes

$$(69) \quad \frac{dX}{dt} = X \frac{X^T \text{diag}(X) - (X^T \text{diag}(X))^T}{2} + \frac{\text{diag}(X) X^T - (\text{diag}(X) X^T)^T}{2} X,$$

which, in spirit, is a continuous analog of the Jacobi method for the singular value decomposition. The asymptotic stability property of (69) is similar to that of Theorem 4.3 [15].

## 6. Nearest Normal Matrix Problems.

Thus far, we have kept the discussion in the real context. But such a restriction is solely for the reason of facilitating the notion. With appropriate modifications, all the discussion can be generalized to the complex-valued case. Applications of such a generalization include, for example, the following problems:

**(Problem D)** Given a general matrix  $A \in R^{n \times n}$  and a set of eigenvalues  $\{\lambda_1 \pm i\nu_1, \dots, \lambda_q \pm i\nu_q, \lambda_{2q+1}, \dots, \lambda_n\}$  where  $\lambda_k, \nu_k$  are real numbers and  $\nu_k \neq 0$ , find a real normal matrix that has the prescribed set as its spectrum and best approximate  $A$  in the Frobenius norm.

**(Problem E)** Given an arbitrary matrix  $A \in C^{n \times n}$ , find its closest normal matrix in the Frobenius norm [42, 54].

Problem D can be reformulated as [35]:

$$(70) \quad \begin{array}{ll} \text{Minimize} & F(Q) := \frac{1}{2} \|Q^T \Lambda Q - A\|^2 \\ \text{Subject to} & Q \in \mathcal{O}(n) \end{array}$$

where

$$(71) \quad \Lambda := \left\{ \left[ \begin{array}{cc} \lambda_1 & \nu_1 \\ -\nu_1 & \lambda_1 \end{array} \right], \dots, \left[ \begin{array}{cc} \lambda_q & \nu_q \\ -\nu_q & \lambda_q \end{array} \right], \lambda_{2q+1}, \dots, \lambda_n \right\}$$

Note Problem D is not equivalent to the Wielandt-Hoffman Theorem in that the normal matrix  $Q^T \Lambda Q$  is only real-valued [43]. Using the same idea as before, we find that [22]

$$(72) \quad \begin{array}{l} \frac{dQ}{dt} = Q \frac{[Q^T \Lambda Q, A^T] - [Q^T \Lambda Q, A^T]^T}{2} \\ Q(0) = I \end{array}$$

defines a steepest descent flow for problem (70) and is the dual problem of

$$(73) \quad \begin{array}{l} \frac{dX}{dt} = \left[ X, \frac{[X, A^T] - [X, A^T]^T}{2} \right] \\ X(0) = \Lambda. \end{array}$$

Problem E can be reformulated as [35]:

$$(74) \quad \begin{array}{ll} \text{Minimize} & G(U, D) := \frac{1}{2} \|A - UDU^*\|^2 \\ \text{Subject to} & U \in \mathcal{U}(n) \\ & D \in \mathcal{D}(n) \end{array}$$

where  $\mathcal{U}(n)$  is the group of all unitary matrices in  $C^{n \times n}$  and  $\mathcal{D}(n)$  is the subspace of all diagonal matrices in  $C^{n \times n}$ . Obviously, for any given  $U \in \mathcal{U}(n)$ , the best  $D \in \mathcal{D}(n)$  that minimizes  $G(U, D)$  is  $D = \text{diag}(U^*AU)$ . Therefore, at a global minimum, problem (74) is equivalent to

$$(75) \quad \begin{array}{ll} \text{Minimize} & H(U) := \frac{1}{2} \|U^*AU - \text{diag}(U^*AU)\|^2 \\ \text{Subject to} & U \in \mathcal{U}(n). \end{array}$$

The closest normal matrix is thus characterized by the following theorem [54]:

THEOREM 6.1. *Let  $A \in C^{n \times n}$  and let  $Z = UDU^*$  with  $U \in \mathcal{U}(n)$  and  $D \in \mathcal{D}(n)$ . Then  $Z$  is the closest normal matrix to  $A$  in the Frobenius norm if and only if the unitary matrix  $U$  is a global minimizer of problem (75) and the diagonal matrix  $D = \text{diag}(U^*AU)$ .*

Problem (75) is very similar to Problem C. By identifying any complex matrix  $Z$  as a pair of real matrices  $(\Re Z, \Im Z)$  where  $\Re Z$  and  $\Im Z$  are the real and the imaginary part of  $Z$ , we introduce an inner product on  $C^{n \times n}$  by

$$(76) \quad \langle A, B \rangle_C := \langle \Re A, \Re B \rangle + \langle \Im A, \Im B \rangle .$$

Then resembling the techniques used in the preceding section, we can show that [19]

$$(77) \quad \begin{aligned} \frac{dU}{dt} &= Uk(U) \\ U(0) &= I \end{aligned}$$

with

$$(78) \quad k(U) := \frac{[U^*AU, \text{diag}(U^*AU)] - [U^*AU, \text{diag}(U^*AU)]^*}{2}$$

defined a steepest descent vector field on  $\mathcal{U}(n)$  for problem (75). It is now easy to see that the matrix  $W(t) := U(t)^*AU(t)$  satisfies the differential equation

$$(79) \quad \frac{dW}{dt} = [W, k(U)].$$

The closest normal matrix can then be constructed from the  $\omega$ -limit point of (77) according to Theorem 6.1.



## 7. Inverse Eigenvalue Problem.

Given  $A_0, \dots, A_n \in \mathcal{S}(n)$ , let  $\mathcal{A}$  denote the affine subspace consisting of all  $A(c) \in \mathcal{S}(n)$  where

$$(80) \quad A(c) := A_0 + \sum_{i=1}^n c_i A_i$$

and  $c := (c_1, \dots, c_n) \in R^n$ . The following is called an inverse eigenvalue problem [33]:

**(Problem F)** Given a set of real numbers  $\{\lambda_1, \dots, \lambda_n\}$ , find coefficients  $c_1, \dots, c_n$  such that  $A(c) \in \mathcal{A}$  has the prescribed set as its spectrum.

We shall assume, without loss of generality, that  $A_1, \dots, A_n$  are mutually orthonormal with respect to the Frobenius inner product (Obviously, we may apply the Gram-Schmidt orthogonalization process to achieve this if necessary, and then the two bases are related by an upper triangular matrix). We may also assume that  $A_0$  is perpendicular to all  $A_i$  for  $i = 1, \dots, n$ . Given any  $X \in \mathcal{S}(n)$ , it is easy to see that the distance between  $X$  and  $\mathcal{A}$  is given by

$$(81) \quad \text{dist}(X, \mathcal{A}) = \|X - (A_0 + P(X))\|$$

where  $P(X)$  is the projection of  $X$  onto the subspace spanned by  $A_1, \dots, A_n$ , and, therefore, is given by

$$(82) \quad P(X) = \sum_{i=1}^n \langle X, A_i \rangle A_i.$$

One approach of solving Problem F is to consider the following minimization problem:

$$(83) \quad \begin{array}{ll} \text{Minimize} & F(Q) := \frac{1}{2} \|Q^T \Lambda Q - A_0 - P(Q^T \Lambda Q)\|^2 \\ \text{Subject to} & Q \in \mathcal{O}(n). \end{array}$$

where  $\Lambda := \text{diag}\{\lambda_1, \dots, \lambda_n\}$ . It can be shown that

$$(84) \quad \nabla F(Q) = 2\Lambda Q \{Q^T \Lambda Q - A_0 - P(Q^T \Lambda Q)\}.$$

Therefore, the differential equation

$$(85) \quad \frac{dQ}{dt} = Qk(Q)$$

with

$$(86) \quad k(Q) := [Q^T \Lambda Q, A_0 + P(Q^T \Lambda Q)]$$

defines a steepest descent vector field on  $\mathcal{O}(n)$  for problem (83). Correspondingly, the matrix  $X(t) := Q(t)^T \Lambda Q(t)$  satisfies the equation

$$(87) \quad \frac{dX}{dt} = [X, [X, A_0 + P(X)]],$$

stays on the manifold  $\mathcal{M}(\Lambda)$ , and moves in the direction to minimize the distance between the two sets  $\mathcal{M}(\Lambda)$  and  $\mathcal{A}$ . Suppose  $X(t) \rightarrow \hat{X}$  as  $t \rightarrow \infty$  and suppose  $\hat{X}$  is also in  $\mathcal{A}$ . Then the coefficients needed in Problem F are determined from  $c_i = \langle \hat{X}, A_i \rangle$ .

Problem B, the inverse Toeplitz eigenvalue problem, is just a special case of Problem F in which  $\mathcal{A}$  is the linear subspace  $\mathcal{T}$  and  $A_k = E_k$  as defined in (47). By experimenting with (87) for the

inverse Toeplitz eigenvalue problem numerically, we have found that sometimes  $\mathcal{M}(\Lambda)$  unfortunately contains stable equilibrium points that are not Toeplitz [15, 27]. This fact does not cause any serious computational difficulty since we can easily change to another initial value and restart the integration. In the following, nevertheless, we wish to re-design the differential equation so that equilibrium points must be in  $\mathcal{A}$ .

To make sure  $X(t)$  stay on the isospectral surface  $\mathcal{M}(\Lambda)$ , we shall consider an initial value problem of the form

$$(88) \quad \begin{aligned} \frac{dX}{dt} &= [X, k(X)] \\ X(0) &= \Lambda \end{aligned}$$

where  $k$  is a mapping from  $\mathcal{S}(n)$  into  $\mathcal{S}(n)^\perp$ . Furthermore, we shall require  $k$  to be an annihilator of the affine subspace  $\mathcal{A}$ . That is, we want  $k$  to be such that

$$(89) \quad k(X) = 0 \text{ if and only if } X \in \mathcal{A}.$$

In view of the dimensions of the three spaces, the construction of a mapping  $k : \mathcal{S}(n) \rightarrow \mathcal{S}(n)^\perp$  with property (89) is possible.

We have observed earlier that  $\|X(t)\| = \|\Lambda\|$  for all  $t \in \mathbb{R}$ . Suppose all elements in  $\Lambda$  are distinct. Then  $[X, k(X)] = 0$  if and only if  $k(X)$  is a polynomial of  $X$  [35]. But then  $k(X) \in \mathcal{S}(n) \cap \mathcal{S}(n)^\perp = \{0\}$ . If condition (89) holds, then we find that all equilibria of (88) are necessarily in  $\mathcal{A}$ . A bounded flow necessarily has a non-empty invariant  $\omega$ -limit set [5]. If  $X(t) \rightarrow \hat{X}$  as  $t \rightarrow \infty$ , then  $\hat{X}$  is a solution of Problem F.

As an example, we now apply the above idea to the inverse Toeplitz eigenvalue problem. Since  $\mathcal{A} = \mathcal{T}$  is a linear subspace, we may require  $k : \mathcal{S}(n) \rightarrow \mathcal{S}(n)^\perp$  to be linear as well. One simple way of defining  $k$  is by

$$(90) \quad k_{ij} := \begin{cases} x_{i+1,j} - x_{i,j-1} & \text{if } 1 \leq i < j \leq n \\ 0 & \text{if } 1 \leq i = j \leq n \\ x_{i,j-1} - x_{i+1,j} & \text{if } 1 \leq j < i \leq n \end{cases}$$

where  $k_{ij}$  denotes the  $(i, j)$ -component of  $k(X)$ . It is obvious that  $\text{kernel}(k) = \mathcal{T}$ .

Let  $C \in \mathcal{S}(n)$  denote the matrix

$$(91) \quad C := \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & & \vdots \\ 0 & 1 & \ddots & & \\ \vdots & & & & 0 & 1 \\ 0 & \dots & & & 1 & 0 \end{bmatrix}.$$

Let

$$(92) \quad \mathcal{S}(n) = \mathcal{T} \oplus \mathcal{T}^C$$

denote any direct sum splitting of  $\mathcal{S}(n)$  with  $\mathcal{T}^C$  denoting the complementary subspace of  $\mathcal{T}$  in  $\mathcal{S}(n)$ . Another way of defining  $k$  is by

$$(93) \quad k(X) := [L(X), C]$$

where  $L(X)$  is the projection of  $X$  onto  $\mathcal{T}^C$  along  $\mathcal{T}$ . Note  $k$  is still a linear map from  $\mathcal{S}(n)$  into  $\mathcal{S}(n)^\perp$ . If  $k(X) = 0$ , then

$$(94) \quad L(X) = \sum_{k=1}^{n-1} \gamma_k C^k$$

since  $C$  has distinct eigenvalues [35]. Observe that if any diagonal of  $X$  is constant, then the corresponding diagonal of  $L(X)$  is zero. Observe also that for each fixed  $j = 1, \dots, n$ , the only matrix among the powers  $C, C^2, \dots, C^{n-j}$  such that the  $(n-j+1)$ -th diagonal is not entirely zero is  $C^{n-j}$ . Indeed, the elements there are all 1. By induction, therefore,  $\gamma_{n-j} = 0$  for all  $j = 1, \dots, n$ . Thus,  $k(X) = 0$  if and only if  $X \in \mathcal{T}$ .

We have experimented both (90) and (93) extensively with different spectral data for the inverse Toeplitz eigenvalue problem. We have found that the orbit *always* converges to an equilibrium point. Thus, we conjecture that the inverse Toeplitz eigenvalue problem is always solvable. What remains to be proved, however, is that the solution of (88) does converge to a single point in theory. No other invariant set such as limit cycles or strange attractors should occur [5]. In this way, we would have settled the existence question. Unfortunately, this convergence proof is not an easy problem either. Despite this theoretical difficulty, we suggest that following the solution flow of (88) is a feasible numerical method for solving the inverse eigenvalue problem. By shifting  $\Lambda$  by  $\Lambda + \sigma I$  with a sufficiently large  $\sigma \in R$  if necessary, we may assume  $\Lambda$  is positive definite. We have observed numerically that the total kinetic energy  $T(t)$  of rotation (19) is monotonically decreasing to zero. Therefore, it appears that  $T(t)$  could be used as a Lyapunov function. The proof, again, is not available at the present time.

The matrix differential equations (87) and (88) offer a new avenue of attacking the inverse eigenvalue problems. They are interesting because of their generality and versatility. There are still, however, many open areas that deserve further investigation. We hope this paper will stimulate some useful discussion either in the theoretical or in the numerical aspect.

### 8. Inverse Non-negative Eigenvalue Problem.

In this section we shall construct a matrix differential equation for another type of inverse eigenvalue problem that is different from Problem F:

**(Problem E)** Given a set of real values  $\{\lambda_1, \dots, \lambda_n\}$  that, by some means, is known a priori to be the spectrum of some non-negative matrices, find a symmetric non-negative matrix whose spectrum is precisely  $\{\lambda_1, \dots, \lambda_n\}$ .

Our approach is similar to the projected gradient flow discussed earlier — we want to minimize the Frobenius distance between the cone  $\pi_s(R_+^n)$  of symmetric non-negative matrices and the isospectral surface  $\mathcal{M}(\Lambda)$  of the given spectrum. The optimization problem is formed as follows:

$$\begin{aligned}
 & \text{Minimize} && F(Q, R) := \frac{1}{2} \|Q^T \Lambda Q - R * R\|^2 \\
 & \text{Subject to} && Q \in \mathcal{O}(n) \\
 (95) & && R \in \mathcal{S}(n)
 \end{aligned}$$

where  $\Lambda := \text{diag}\{\lambda_1, \dots, \lambda_n\}$  and  $*$  denotes the Hadamard product. In the space  $R^{n \times n} \times R^{n \times n}$  we shall use the induced Frobenius inner product defined in (58). Then the gradient of  $F$  in (95) is given by [21]

$$(96) \quad \nabla F(Q, R) = (2\Lambda Q(Q^T \Lambda Q - R * R), -2(Q^T \Lambda Q - R * R) * R).$$

It is obvious that the tangent space of  $\mathcal{O}(n) \times \mathcal{S}(n)$  at  $(Q, R)$  is given by

$$(97) \quad T_{(Q,R)}\mathcal{O}(n) \times \mathcal{S}(n) = Q\mathcal{S}(n)^\perp \times \mathcal{S}(n).$$

So the projection of  $\nabla F(Q, R)$  onto the manifold  $\mathcal{O}(n) \times \mathcal{S}(n)$  can be calculated. The initial value problem,

$$\begin{aligned}
 \frac{dQ}{dt} &= Q\{Q^T \Lambda Q(R * R) + (R * R)Q^T \Lambda Q\} \\
 \frac{dR}{dt} &= 2(Q^T \Lambda Q - R * R) * R \\
 Q(0) &= \Lambda \\
 (98) \quad R(0) &= \text{an arbitrary positive matrix,}
 \end{aligned}$$

therefore, defines a steepest descent flow on  $\mathcal{O}(n) \times \mathcal{S}(n)$  for problem (95).

Define

$$(99) \quad X(t) := Q(t)^T \Lambda Q(t)$$

and

$$(100) \quad Y(t) := R(t) * R(t).$$

Then it is easy to see that  $(X(t), Y(t))$  satisfies the system of equations

$$\begin{aligned}
 \frac{dX}{dt} &= [X, [X, Y]] \\
 (101) \quad \frac{dY}{dt} &= 4Y * (X - Y).
 \end{aligned}$$

Note that  $X(t)$  and  $Y(t)$  moves, respectively, in the isospectral surface  $\mathcal{M}(\Lambda)$  and the cone  $\pi_s(R_+^n)$  so as to reduce the distance  $G(t) := \|X(t) - Y(t)\|^2$ . Therefore,  $G(t)$  serves as a natural Lyapunov function. Recently, we are able to use center manifold theory [11] to study the structure of  $\omega$ -limit sets of (101) [21]:

**THEOREM 8.1.** *If  $(\hat{X}, \hat{X}) \in \mathcal{M}(\Lambda) \times \pi_s(R_+^n)$  ever becomes an  $\omega$ -limit point of an orbit  $(X(t), Y(t))$  of (101), then  $\lim_{t \rightarrow \infty} (X(t), Y(t)) = (\hat{X}, \hat{X})$ .*

So far as we know, most of the discussions in the literature are centered around establishing a sufficient or a necessary condition so that a given set of values is the spectrum of a non-negative matrix [4, 31]. Very few of these theoretical results are ready to be implemented to find the actual matrix [21]. By following the integral curve of (101), however, we have a numerical algorithm that systematically reducing the distance between  $\mathcal{M}(\Lambda)$  and  $\pi_s(R_+^n)$ . If these two sets do intersect, then of course the distance is zero. Otherwise, our approach still finds a matrix from  $\mathcal{M}(\Lambda)$  and a matrix from  $\pi_s(R_+^n)$  such that their distance is a local minimum. In the latter case, because  $\pi_s(R_+^n)$  is a convex set, the matrix from  $\pi_s(R_+^n)$  must lie on a facet of the cone, i.e., some of the components of the non-negative matrix is zero.

### 9. Quadratic Assignment Problem.

The quadratic assignment problem consists of [32, 39]:

$$(102) \quad \begin{array}{ll} \text{Minimize} & \langle C + ASB, S \rangle \\ \text{Subject to} & S \in \Pi \end{array}$$

where  $A, B, C \in R^{n \times n}$  are given matrices and  $\Pi$  is the set of all permutation matrices. Problem (102) is known to be an NP-hard problem.

In view of the success of the above discussions, we suggest the following continuous realization process for problem (102). First we relax (102) to the problem

$$(103) \quad \begin{array}{ll} \text{Minimize} & F(Q) := \langle C + AQB, Q \rangle \\ \text{Subject to} & Q \in \mathcal{O}(n) \end{array}$$

so that the idea of projected gradient can be applied. It can be shown that

$$(104) \quad \nabla F(Q) = C + A^T Q B + A Q B^T.$$

Therefore, the differential equation

$$(105) \quad \frac{dQ}{dt} = Q k(Q)$$

with

$$(106) \quad k(Q) := \frac{(C^T Q - Q^T C) + (B^T Q^T A Q - Q^T A^T Q B) + (B Q^T A^T Q - Q^T A Q B^T)}{2}$$

defines a steepest descent flow on  $\mathcal{O}(n)$  for problem (103). By tracing the integral curve of (105), a limit point  $\hat{Q}$  of (105) (and, hence a local optimizer of problem (103)) should be found. Because of the continuity of the objective function, we think the permutation matrix that is nearest to  $\hat{Q}$  should be a putative solution to problem (102). Thus, we solve the problem

$$(107) \quad \begin{array}{ll} \text{Minimize} & \|S - \hat{Q}\|^2 \\ \text{Subject to} & S \in \Pi. \end{array}$$

But problem (107) is equivalent to

$$(108) \quad \begin{array}{ll} \text{Maximize} & \sum_{i=1}^n \hat{q}_{i\sigma(i)} \\ \text{Subject to} & \sigma \in \Pi \end{array}$$

where  $\hat{q}_{ij}$  are the components of  $\hat{Q}$ . Problem (108) is a classical linear assignment problem and, hence, can be solved by many well developed techniques [40].

## 10. Conclusion.

Matrix differential equations by nature are complicated, since the components are coupled into nonlinear terms. Nonetheless, as we have demonstrated, there have been substantial advances in understanding some of the dynamics. For the time being, the numerical implementation is still very primitive. But most important of all, we think there are many opportunities where new algorithms may be developed from the realization process. It is hoped that this paper has conveyed some values of this idea.

## REFERENCES

- [1] E. Allgower, A survey of homotopy methods for smooth mappings, in *Numerical Solution of Nonlinear Equations*, Lecture Notes in Math., 878, 1980, 1-29.
- [2] E. Allgower and K Georg, Simplicial and continuation methods for approximating fixed points and solutions to systems of equations, *SIAM Review*, 22(1980), 28-85.
- [3] V. I. Arnold, *Geometrical Methods in the Theory of Ordinary Differential Equations*, 2nd Ed., Springer-Verlag, New York, 1988.
- [4] A. Berman and R. J. Plemmons, *Non-negative Matrices in the Mathematica Sciences*, Academic Press, New York, 1979.
- [5] F. Brauer and J. A. Nohel, *Qualitative Theory of Ordinary Differential Equations*, Benjamin, New York, 1969.
- [6] R. W. Brockett, Dynamical systems that sort lists, diagonalize matrices and solve linear programming problems, preprint, 1989.
- [7] R. W. Brockett, Least squares matching problems, *Linear Alg. Appl.*, 122/123/124(1989), 761-777.
- [8] A. Bunse-Gerstner, Matrix factorizations for symplectic QR-like methods, *Linear Alg. Appl.*, 83(1986), 49-77.
- [9] A. Bunse-Gerstner, R. Byers and V. Mehrmann, A chart of numerical method for structured eigenvalue problems, *Materialien LVII*, Universität Bielefeld, preprint, 1989.
- [10] R. Byers, A Hamiltonian QR Algorithm, *SIAM J. Sci. Stat. Comput.*, 7(1986), 212-229.
- [11] Carr, J., *Applications of Center Manifold Theory*, Applied Mathematical Sciences, 35, Springer-Verlag, Berlin, 1981.
- [12] M. T. Chu, On the continuous realization of iterative processes, *SIAM Review*, 30(1988), 375-387.
- [13] M. T. Chu, T. Y. Li, and T. Sauer, Homotopy method for general  $\lambda$ -matrix problems, *SIAM J. Matrix Anal. Appl.*, 9(1988), 528-536.
- [14] M. T. Chu and L. K. Norris, Isospectral flows and abstract matrix factorizations, *SIAM J. Numer. Anal.*, 25(1988), 1383-1391.
- [15] M. T. Chu and K. R. Driessel, The projected gradient method for least squares matrix approximations with spectral constraints, *SIAM J. Numer. Anal.*, to appear.
- [16] M. T. Chu, The generalized Toda flow, the QR algorithm and the center manifold theorem, *SIAM J. Alg. Disc. Meth.*, 5(1984), 187-210.
- [17] M. T. Chu, A continuous approximation to the generalized Schur decomposition, *Linear Alg. Appl.*, 78(1986), 119-132.
- [18] M. T. Chu, On a differential equation approach to the singular value decomposition of bi-diagonal matrices, *Linear Alg. Appl.*, 80(1986), 71-79.
- [19] M. T. Chu, A continuous Jacobi-like approach to the simultaneous reduction of real matrices, *Linear Alg. Appl.*, to appear in the special issue on "Matrix Canonical Forms".
- [20] M. T. Chu and H. Hamilton, Parallel solution of ODE's by multiblock methods, *SIAM J. Sci. Stat. Comput.*, 8(1987), 342-353.
- [21] M. T. Chu and K. R. Driessel, Constructing symmetric non-negative matrices with prescribed eigenvalues by differential equations, *SIAM J. Math. Anal.*, to appear.
- [22] M. T. Chu, Least squares approximation by real normal matrices with specified spectrum, *SIAM J. Matrix. Appl.*, to appear.



- [23] M. L. Curtis, *Matrix Groups*, Springer-Verlag, New York, 1979.
- [24] P. Deift, T. Nanda and C. Tomei, Differential equations for the symmetric eigenvalue problem, *SIAM J. Numer. Anal.*, 20(1983), 1-22.
- [25] J. Della-Dora, Numerical Linear Algorithms and Group Theory, *Linear Alg. Appl.*, 10(1975), 267-283.
- [26] P. Delsarte and Y. Genin, Spectral properties of finite Toeplitz matrices, *Proceedings of the 1983 International Symposium of Mathematical Theory of Networks and Systems*, Beer-Sheva, Israel, 194-213.
- [27] K. R. Driessel and M. T. Chu, Can real symmetric Toeplitz matrices have arbitrary real spectra?, *SIAM J. Matrix Anal. Appl.*, under revision.
- [28] K. R. Driessel, On isospectral gradient flow — solving matrix eigenproblems using differential equations, in *Inverse Problems*, J. R. Cannon and U Hornung, ed., ISNM 77, Birkhauser, 1986, 69-91.
- [29] T. Eirola and O. Nevanlinna, What do multistep methods approximate? *Numer. Math.*, 53(1988), 559-569.
- [30] L. E. Faibusovich, QR-type factorizations, the Yang-Baxter equations, and an eigenvalue problem of control theory, *Linear Alg. Appl.*, 122/123/124(1989), 943-971.
- [31] M. Fiedler, Eigenvalues of non-negative symmetric matrices, *Linear Alg. Appl.*, 9(1974), 119-142.
- [32] G. Finke, R. E. Burkard and F. Rendl, Quadratic assignment problems, *Annals Discrete Math.*, 31(1987), 61-82.
- [33] S. Friedland, J. Nocedal and M. L. Overton, The formulation and analysis of numerical methods for inverse eigenvalue problems, *SIAM J. Numer. Anal.*, 24(1987), 634-667.
- [34] C. B. Garcia and W. I. Zangwill, *Pathways to Solutions, Fixed Points, a Equilibria*, Prentice-Hall, Englewood Cliffs, NJ, 1981.
- [35] F. R. Gantmacher, *Matrix Theory*, Vol. 1 and 2, Chelsea, New York, 1959.
- [36] P. E. Gill, W. Murray and M. H. Wright, *Practical Optimization*, Academic Press, London, 1981.
- [37] H. Goldstein, *Classical Mechanics*, Addison-Wesley, Massachusetts, 1965.
- [38] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 2nd ed., Johns Hopkins, Baltimore, 1989.
- [39] S. W. Hadley, F. Rendl and H. Holkowitz, A new lower bound via elimination for the quadratic assignment problem, *Research Report CORR 89-5*, University of Waterloo, 1989.
- [40] P. L. Hammer, E. L. Johnson and B. H. Korte, *Discrete Optimization*, vol I and II, Elsevier, North-Holland, 1979.
- [41] S. Helgason, *Differential Geometry, Lie Groups and Symmetric Spaces*, Academic, New York, 1978.
- [42] N. J. Higham, Matrix nearness problems and applications, *Applications of Matrix Theory*, S. Barnett and M. J. C. Gover, ed., Oxford University Press, 1989, 1-27.
- [43] R. A. Horn and C. A. Johnson, *Matrix Analysis*, Cambridge University Press, Cambridge, 1985.
- [44] H. B. Keller, Global homotopies and Newton methods, in *Recent Advances Numerical Analysis*, C. de Boor and G. Golub, eds., Academic Press, New York, 1978, 73-94.
- [45] P. E. Kloeden and J. Lorenz, A note on multistep methods and attracting sets of dynamical systems, *Numer. Math.*, 56(1990), 667-673.
- [46] D. P. Laurie, A numerical approach to the inverse Toeplitz eigenproblem *SIAM J. Sci. Stat. Comput.*, 9(1988), 401-405.

- [47] L. C. Li and S. Parmentier, A new class of quadratic Poisson structures and the Yang-Baxter equation, *Math. Phys.*, 307(1988), 279-281.
- [48] T. Y. Li and N. H. Rhee, Homotopy algorithm for symmetric eigenvalue problems, *Numer. Math.*, 55(1989), 265-280.
- [49] T. Y. Li, T. Sauer and J. Yorke, Numerical solution of a class of deficient polynomial systems, *SIAM J. Numer. Anal.*, 24(1987), 435-451.
- [50] A. P. Morgan, *Solving Polynomial Systems Using Continuation for Scientific and Engineering Problems*, Prentice-Hall, Englewood Cliff, NJ, 1987.
- [51] Nakamura, Y., Fractional transformation group induced by QR factorization and linear prediction problems, *System and Control Letters*, 10(1988), 181-184.
- [52] Nakamura, Y., Group actions on linear predictors for non-stationary processes, *IMA J. Math. Control Inform.*, 5(1988), 69-75.
- [53] W. C. Rheinboldt, *Numerical Analysis of Parameterized Nonlinear Equations*, John Wiley and Sons, New York, 1986.
- [54] A. Ruhe, Closest normal matrix finally found!, *BIT*, 27(1987), 585-595.
- [55] W. W. Symes, The QR algorithm and scattering for the finite non-periodic Toda lattice, *Physica*, 4D(1982), 275-280.
- [56] F. Warner, *Foundations of Differentiable Manifolds and Lie Group*, Springer-Verlag, New York, 1983.
- [57] D. S. Watkins, Isospectral flows, *SIAM Rev.*, 26(1984), 379-391.
- [58] D. S. Watkins and L. Elsner, Self-similar flows, *Linear Alg. Appl.*, 110(1988), 213-242.
- [59] L. T. Watson, S. C. Billups and A. P. Morgan, HOMPACk: A suite of codes for globally convergent homotopy algorithms, *ACM Trans. Math. Software*, 13(1987), 281-310.
- [60] L. T. Watson, Numerical linear algebra aspects of globally convergent homotopy methods, *SIAM Review*, 28(1986), 529-545.
- [61] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.