

Adaptive Machine Learning Algorithms

Proofs of Convergence

Chanchal Chatterjee

14521 Weeth Drive

San Jose, CA 95124

Vwani P. Roychowdhury

6371C Boelter Hall,

UCLA, Los Angeles, CA 90095

August 30, 2021

1. Introduction to the Proofs of Convergence

1.1 Objective Functions for Adaptive Algorithms

In order to compute matrix functions by adaptive algorithms, we need a convenient methodology or strategy to derive them. One such approach is to identify an *objective or cost function* J , which is a function of a parameter matrix W and the asymptotic matrix A . Given an objective function $J(W;A)$, we minimize it to obtain W^* as:

$$W^* = \arg \min_W J(W;A) \quad (1.1)$$

where W^* is the desired matrix function of A . For example, if we are computing the inverse square root ($A^{-1/2}$) of A , then $W^* = A^{-1/2}$. A crucial step in this strategy is to identify an objective function $J(W;A)$ such that its minimizer W^* is the desired matrix function of A . Obtaining such objective functions offers several benefits as described below.

Derive Adaptive Algorithms

Since the minimizer W^* of the objective function $J(W;A)$ is the desired matrix function of A , we can obtain several adaptive algorithms to compute W^* by applying standard optimization techniques [Luenberger 84] to the objective function $J(W;A)$. Examples of optimization techniques are: (1) gradient descent, (2) steepest descent, (3) conjugate direction, (4) Newton-Raphson, and (5) Recursive Least Squares (RLS). An adaptive algorithm for a data sequence $\{A_k\}$ can be obtained by replacing A with A_k , although a rigorous convergence proof is necessary for this substitution.

For example, the standard gradient descent algorithm for minimizing $J(W;A)$ is:

$$W_{k+1} = W_k - \eta_k \nabla_W J(W_k;A), \quad (1.2)$$

where η_k is a positive gain sequence, and W_0 is a starting matrix. In order to obtain an adaptive algorithm that uses the online data sequence $\{A_k\}$ instead of its asymptotic value A , we modify (1.2) as:

$$W_{k+1} = W_k - \eta_k \nabla_W J(W_k;A_k). \quad (1.3)$$

While the convergence of (1.2) is well-established, the convergence of (1.3) is not guaranteed by classical optimization theory. Hence, we use an alternative method due to Stochastic Approximation Theory [Kushner&Clark 78; Ljung 77,78,84,92; Benveniste *et al.* 90] to prove the convergence of (1.3).

Speedup of Adaptive Algorithms

The availability of the objective function $J(W;A)$ allows us to speed up the basic gradient descent stochastic approximation algorithm (1.3). By applying speedup techniques in optimization theory such as steepest descent, conjugate direction, Newton-Raphson and recursive least squares, we can enhance the speed of the stochastic approximation algorithm (1.3). Details of these methods for principal component analysis are given in Chapter 6.

Convergence Analysis

Although the adaptive algorithms are derived by following standard optimization techniques on the objective function $J(W;A)$, their derivations do not constitute a proof of convergence. For example, algorithm (1.3) is obtained from (1.2). Although the convergence of (1.2) is guaranteed by standard optimization theory, the convergence of (1.3) should be proven rigorously. Hence, it is important to provide a convergence analysis for our adaptive algorithms.

In simple terms, our derivations of the algorithms from objective function $J(W;A)$ show that the descent direction of J is the same as the average evolution direction of the adaptive algorithms. We still need to show that the global minimum of the objective function $J(W;A)$ is the desired matrix function of A . Hence, a study of the landscape of the objective function is necessary to determine the convergence of the stochastic approximation algorithm (1.3). The objective function offers an energy or Lyapunov function for such analysis.

The stationary points W^* of the objective function $J(W;A)$ are given by the *Kuhn-Tucker conditions*:

$$\nabla_W J(W^*; A) = 0. \quad (1.4)$$

In almost all our algorithms, this equation (1.4) is same as equating the ODE (see (1.7) below) to 0. Some of these stationary points may be stable equilibrium or local minimum points, whereas some may be unstable equilibrium or saddle points. In order to determine the stable equilibrium points, we compute the *Hessian* H of the objective function at the equilibrium points W^* :

$$H = \nabla_{WW}^2 J(W^*; A). \quad (1.5)$$

Stable equilibrium points that are local minima have positive definite Hessians, whereas unstable equilibrium points have indefinite Hessians, and local maxima have negative definite Hessians.

1.2 Common Methodology for Derivations and Convergence Proofs

The literature for adaptive algorithms for matrix computation offers a wide range of techniques (including ad hoc methods) and various types of convergence procedures. In this study, we present a *common methodology* to derive and prove the convergence of our adaptive algorithms.

In the stochastic approximation (adaptive) algorithm (1.3), for each time instant k , let W_k be an instantaneous estimate of the desired matrix function W of the asymptotic matrix A . Thus, if we are computing the inverse square root ($A^{-1/2}$) of A , then W_k is an instantaneous estimate of $A^{-1/2}$ at time instant k . For each algorithm, we offer an update rule for W_k for each new observation A_k , such that W_k converges to the desired matrix function of A . Steps for the derivation and convergence analysis of each adaptive algorithm is discussed below:

1. We first present an *Objective Function* $J(W; A_k)$ such that the minimizer W^* of J is the desired matrix function of the asymptotic data matrix A .
2. Derive an *Adaptive Update Rule* for matrix W by applying the gradient descent technique on the objective function $J(W; A_k)$. Note that other methods of nonlinear optimization such as steepest descent, conjugate direction, Newton-Raphson or Recursive Least Squares (RLS) can also be used on the objective function. The adaptive gradient descent update rule is:

$$W_{k+1} = W_k - \eta_k \nabla_W J(W_k, A_k) = W_k + \eta_k h(W_k, A_k), \quad (1.6)$$

where the function $h(W_k, A_k)$ follows certain continuity and regularity properties, and η_k is a decreasing gain sequence.

3. In order to prove the convergence of W_k in (1.6) to the desired matrix function of A , we employ the well-known theory of *Stochastic Approximation* [Kushner&Clark 78; Ljung 77,78,84,92; Benveniste *et al.* 90]. In most instances, we assume a stationary data sequence $\{A_k\}$ in order to prove the convergence of the adaptive algorithm by stochastic approximation theory, although practical implementations of the same algorithms on non-stationary sequences can be achieved and yield good results.
4. As an intermediate step of the Stochastic Approximation convergence analysis, we obtain an *Ordinary Differential Equation (ODE)* $dW(t)/dt = \lim_{k \rightarrow \infty} E[h(W, A_k)]$, where $W(t)$ is the continuous time counterpart of W_k . The stable stationary solution of the ODE is a convergence point of the adaptive algorithm (1.6). In solving the ODE, we use the well-known solution for the continuous time *Riccati Differential Equation* [Anderson&Moore

90]. If it is difficult to solve the ODE, we study the landscape of the corresponding Lyapunov function $J(W;A)$. We discuss its stable stationary points and show the convergence properties of the ODE for initial states close to the stable points.

5. In order to estimate the *rate of convergence* of the stochastic approximation algorithm (1.6), we compute the time constants (τ) for some of our algorithms. If we can solve the ODE, we get a solution for $W(t, W(0))$ as a function of $e^{-t/\tau}$. Comparatively smaller time constants (τ) indicate faster convergence of the adaptive algorithm.

In the following two sections, we briefly describe the Stochastic Approximation Theory and a solution to the Riccati Differential Equation.

1.3 Stochastic Approximation Theory

An important property of the algorithms discussed here is that they are stochastic approximation algorithms i.e., they are adaptive and of the type (1.6). In contrast to conventional processes, our data arrives in temporal succession i.e., in vector or matrix sequences $\{\mathbf{x}_k\}$ or $\{A_k\}$, instead of vectors or matrices \mathbf{x} or A respectively. For every data sample \mathbf{x}_k or A_k , we update the estimates of the targeted parameters W_k , and we require that the estimates converge strongly to the desired solution. We also require that the statistical procedure keep pace with the incoming data so that, at any instant, the estimates fully reflect all of the currently available data. Such stochastic approximation procedures are also useful when a given estimate has to adapt to small changes in the data (e.g., a few incoming samples). Thus, if the features are computed with conventional methods from some initial samples, then these estimates can be used as starting values for the stochastic approximation procedure when a few additional samples are available. In this situation, the adaptive techniques allow direct updating in a computationally inexpensive way.

An important advantage of stochastic approximation algorithms is that there is a significant amount of theory associated with them, which can be used to analyze and prove their convergence. We use the stochastic approximation theory due to Ljung [Ljung 77,78,84,92; Benveniste *et al.* 90] which deals with stationary sequences. An alternative proof by a similar approach due to Kushner and Clark [Kushner&Clark 78] can also be used. In a somewhat looser language, stochastic approximation theory states the following:

1. Matrix W_k can converge only to stable stationary points of the Ordinary Differential Equation (ODE):

$$\frac{dW}{dt} = \lim_{k \rightarrow \infty} E[h(W, A_k)], \quad (1.7)$$

where $W(t)$ is the continuous time counterpart of W_k with t denoting continuous time.

2. If W_k belongs to the domain of attraction of a stable stationary point W^* of the ODE infinitely often with probability one (w.p.1), then W_k converges w.p.1 to W^* as $k \rightarrow \infty$.
3. The trajectories of the ODE are the “asymptotic paths” of W_k generated by (1.6). The convergence proof requires the following steps:
 - a. defining a set of assumptions,
 - b. finding the stable stationary points of the ODE, and
 - c. showing that W_k visits the domain of attraction of a stable stationary point infinitely often.

Assumptions

In order to prove the convergence of (1.6), we use Theorem 1 of Ljung [Ljung 77; Benveniste *et al.* 90]. The following is a general set of assumptions for the convergence proof of algorithm (1.6):

Assumption (A1.1). The sequence $\{\mathbf{x}_k\}$ consists of real random vectors, where each \mathbf{x}_k is uniformly bounded with probability one (w.p.1); i.e., $\|\mathbf{x}_k\| < \alpha < \infty$, and $\lim_{k \rightarrow \infty} E[\mathbf{x}_k \mathbf{x}_k^T] = A$ where A is positive definite.

Assumption (A1.2). The gain sequence $\{\eta_k \in \mathbb{R}^+\}$ is decreasing such that $\sum_{k=0}^{\infty} \eta_k = \infty$, $\sum_{k=0}^{\infty} \eta_k^r < \infty$ for some $r > 1$, and $\lim_{k \rightarrow \infty} \sup(\eta_k^{-1} - \eta_{k-1}^{-1}) < \infty$.

Assumption A2.1 is reasonable for most practical implementations, where $\{\mathbf{x}_k\}$ are kept bounded either by deliberate measures or automatically. Methods to keep $\{\mathbf{x}_k\}$ bounded are discussed in Ljung [Ljung 77]. The physical meaning of A2.2 can be described as follows. Condition $\eta_k \rightarrow 0$ allows the process to settle down in the limit whereas, $\sum_{k=0}^{\infty} \eta_k = \infty$ insures that there is enough corrective action to avoid stopping short of the solution. Conditions $\sum_{k=0}^{\infty} \eta_k^r < \infty$ and $\lim_{k \rightarrow \infty} \sup(\eta_k^{-1} - \eta_{k-1}^{-1}) < \infty$ guarantee that the variance of the accumulated noise is finite so that we can correct for the effect of noise. Assumption A2.2 holds for $\eta_k = ck^{-\delta}$ where $c \in \mathbb{R}^+$, and $0 < \delta \leq 1$. The choice of $\delta = 1$ is a leading case. Another alternative used in this study extensively is $\eta_k = c_1(k+c_2)^{-1}$, where $c_1, c_2 \in \mathbb{R}^+$.

In the literature for stochastic approximation proofs, there are many assumptions that are usually made on the statistical properties of $\{\mathbf{x}_k\}$, such as statistically independent and independent and identically distributed (i.i.d). However, Ljung [Ljung 77; Benveniste *et al.* 90] permits far less restrictive choices for $\{\mathbf{x}_k\}$. Specifically, we can assume that $\{\mathbf{x}_k\}$ is generated by a linear structure:

$$\tilde{\mathbf{x}}_k = A(W_k)\tilde{\mathbf{x}}_{k-1} + B(W_k)\mathbf{e}_k \text{ and } \mathbf{x}_k = C(W_k)\tilde{\mathbf{x}}_k, \quad (1.8)$$

or a nonlinear variant:

$$\tilde{\mathbf{x}}_k = g(k, W_k, \tilde{\mathbf{x}}_{k-1}, \mathbf{e}_k) \text{ and } \mathbf{x}_k = h(k, W_k, \tilde{\mathbf{x}}_{k-1}) \quad (1.9)$$

has also been postulated [Ljung 77]. Here $\{\mathbf{e}_k\}$ is a uniformly bounded sequence of independent (not necessarily stationary or with zero means) random vectors. These structures are treated at length in [Ljung 77,92; Benveniste *et al.* 90]. In light of these models, we state the following general assumption for $\{\mathbf{x}_k\}$:

Assumption (A2.3). Sequence $\{\mathbf{x}_k\}$ is generated by (1.8) or (1.9) satisfying the stability and regularity conditions of Ljung [Ljung 77].

The main assumptions of Theorem 1 of Ljung [Ljung 77] are:

- L1.** The function $h(W, A)$ is continuously differentiable with respect to W and A . The derivatives are, for fixed W and A , bounded in k .
- L2.** The so-called *mean vector field* $\bar{h}(W) = \lim_{k \rightarrow \infty} E[h(W, A_k)]$ exists and is regular; i.e., locally Lipschitz. The expectation is with respect to the distribution of A_k for a fixed value of W .

Formulation and Solution of the ODE

We modify the results given by Ljung [Ljung 77] in Theorem 1 to suit our adaptive algorithms in the following Theorem:

Theorem 1.1. Let A2.1-A2.3 hold. Let W^* be a locally asymptotically stable (in the sense of Liapunov) solution for the ordinary differential equation (ODE):

$$\frac{dW}{dt} = \bar{h}(W) = \lim_{k \rightarrow \infty} E[h(W, A_k)] \quad (1.10)$$

with domain of attraction $D(W^*)$. If there is a compact subset $S \subset D(W^*)$ such that $W_k \in S$ infinitely often, then we have $W_k \rightarrow W^*$ with probability one as $k \rightarrow \infty$.

Proof. The proof can be obtained by using Ljung's Theorem 1 [Ljung 77]. ■

The solution for the ODE (1.10) varies from problem to problem. However, there are two common methods of analyzing the ODE. These are:

1. **Global Analysis:** Formulate the ODE as a Riccati equation [Anderson&Moore 90] and use its well-known solution, whereby we obtain the continuous time solution $W(t, W_0)$, where W_0 is the initial condition at $t=0$. This is discussed in detail in Section 1.4. We then obtain the asymptotically stable solution $W^* = \lim_{t \rightarrow \infty} W(t, W_0)$ where W^* is the desired matrix function of A . Conditions for the existence of the solution for $t \geq 0$ gives us the domain of attraction $D(W^*)$. This explicit solution is useful in studying the global convergence properties of the ODE (1.10) for initial states far from the equilibrium solutions.
2. **Local Analysis:** When it is difficult to directly solve the ODE, we analyze the solution in two steps. We first compute all the equilibrium points of the ODE from the equation $\bar{h}(W) = 0$. We then prove that all these equilibrium points of the ODE are unstable equilibrium points, except for the desired matrix function W^* , which is the stable equilibrium point. This analysis allows us to study convergence properties of the ODE (1.10) for initial states close to the stable equilibrium points.

We use both these methods (as needed) to solve the ODE corresponding to the adaptive algorithms presented here.

W_k Visits S Infinitely Often

Although the above analyses give us the limiting values of the ODE, they are not sufficient for the complete convergence proof. One must, in addition, prove that the adaptive algorithm (1.6) is stable; i.e., the weight matrix W_k must remain bounded on some realistic conditions. Such boundedness condition is also necessary for W_k to visit a compact subset S of the domain of attraction of W^* infinitely often. In practical implementation, we can hard-limit the entries of W_k so that their magnitudes remain below a certain limit ρ and thus within a compact region A . An alternative analytical approach also exists. It turns out that, for most of our algorithms, there exists a uniform upper bound for η_k such that W_k is uniformly bounded. The final convergence of the adaptive algorithm (1.6) is guaranteed by Theorem 1 of Ljung and stated in Theorem 2.1.

1.4 Solution of Riccati Differential Equation

Consider the Riccati differential equation (RDE) [Anderson&Moore 90] with time-varying coefficient matrices:

$$\frac{dP(t)}{dt} = PF + F^T P - PGP + Q, \quad P(0) = P_0. \quad (1.11)$$

Associated with the RDE is the linear differential equation:

$$\begin{bmatrix} \dot{X} \\ \dot{Y} \end{bmatrix} = \begin{bmatrix} -F & G \\ Q & F^T \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix}, \quad \begin{bmatrix} X(0) \\ Y(0) \end{bmatrix} = \begin{bmatrix} I \\ P_0 \end{bmatrix}. \quad (1.12)$$

Some important properties regarding (1.11) and (1.12) are summarized in the Theorem below:

Theorem 2.2. Consider two initial value problems (1.11) and (1.12). Then, the solution of (1.11) exists on $[0, T)$ if and only if $X(t)$ is nonsingular on $[0, T)$. Moreover, the solution to (1.11) is unique and is given by:

$$P(t) = Y(t)X(t)^{-1}. \quad (1.13)$$

Exponential Formula for Time-Invariant Problem

Suppose that F , G , and Q are constant matrices or scalars. One can define the so-called *Hamiltonian* matrix H as:

$$H = \begin{bmatrix} -F & G \\ Q & F^T \end{bmatrix}. \quad (1.14)$$

It has no imaginary eigenvalue, given detectability and stabilizability, and if λ is an eigenvalue, so is $-\lambda$. Thus, there exists a real M such that

$$H = M \begin{bmatrix} \Lambda_1 & 0 \\ 0 & \Lambda_2 \end{bmatrix} M^{-1} \quad (1.15)$$

and Λ_1, Λ_2 are real Jordan matrices such that the real parts of all eigenvalues are respectively negative and positive. It follows that:

$$\begin{bmatrix} X(t) \\ Y(t) \end{bmatrix} = M \begin{bmatrix} e^{\Lambda_1 t} & 0 \\ 0 & e^{\Lambda_2 t} \end{bmatrix} M^{-1} \begin{bmatrix} I \\ P_0 \end{bmatrix} = \begin{bmatrix} M_{11}e^{\Lambda_1 t} & M_{12}e^{\Lambda_2 t} \\ M_{21}e^{\Lambda_1 t} & M_{22}e^{\Lambda_2 t} \end{bmatrix} \begin{bmatrix} L \\ RL \end{bmatrix}$$

where L is an unimportant matrix and

$$R = -[M_{22} - P_0 M_{12}]^{-1} [M_{21} - P_0 M_{11}]. \quad (1.16)$$

Moreover,

$$P(t, P_0) = [M_{21} + M_{22}e^{\Lambda_2 t} \text{Re}^{-\Lambda_1 t}]^{-1} [M_{11} + M_{12}e^{\Lambda_2 t} \text{Re}^{-\Lambda_1 t}]. \quad (1.17)$$

Evaluating the Asymptotic Solution

It follows from (1.17) and because $e^{\Lambda_2 t}$ and $e^{-\Lambda_1 t}$ decay to 0 as $t \rightarrow \infty$, that

$$\lim_{t \rightarrow \infty} P(t, P_0) = \bar{P} = M_{21} M_{11}^{-1}. \quad (1.18)$$

Notice that the limit is approached at an exponential rate equal to twice the smallest real part of any eigenvalue of Λ_2 and is independent of the boundary condition of P_0 .

Proofs of Chapter 2

2.1 Convergence Proof for Adaptive Stationary Mean

The objective function $J(\mathbf{w}_k; \mathbf{x}_k)$ whose minimizer \mathbf{w}^* is the asymptotic mean $\mathbf{m} = \lim_{k \rightarrow \infty} E[\mathbf{x}_k]$ is:

$$J(\mathbf{w}_k; \mathbf{x}_k) = \|\mathbf{x}_k - \mathbf{w}_k\|^2. \quad (2.1)$$

The gradient of $J(\mathbf{w}_k; \mathbf{x}_k)$ with respect to \mathbf{w}_k is:

$$(1/2) \nabla_{\mathbf{w}_k} J(\mathbf{w}_k; \mathbf{x}_k) = -(\mathbf{x}_k - \mathbf{w}_k). \quad (2.2)$$

From the gradient in (2.2), we obtain the adaptive gradient descent algorithm

$$\mathbf{m}_k = \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i = \mathbf{m}_{k-1} + \frac{1}{k} (\mathbf{x}_k - \mathbf{m}_{k-1}). \quad (2.3)$$

where $\eta_k = 1/k$. The ODE corresponding to (2.3) is:

$$d\mathbf{w}(t)/dt = \mathbf{m} - \mathbf{w}, \text{ where } \mathbf{w}_0 = \mathbf{w}(0). \quad (2.4)$$

The solution of (2.4) is:

$$\mathbf{w}(t) = \mathbf{m} + (\mathbf{w}_0 - \mathbf{m})e^{-t} \rightarrow \mathbf{m}, \text{ as } t \rightarrow \infty. \quad (2.5)$$

The domain of attraction $D(\mathbf{w}^*) = \{\mathfrak{R}^n\}$. The time constant τ for convergence is 1.

2.2 Convergence Proof for Adaptive Stationary Correlation

The objective function $J(W; A)$ whose minimizer W^* is the asymptotic correlation matrix $A = \lim_{k \rightarrow \infty} E[\mathbf{x}_k \mathbf{x}_k^T]$ is:

$$J(W; A) = \text{tr}((A - W)^T (A - W)) \quad (2.6)$$

The gradient of $J(W; A)$ with respect to W is:

$$(1/2) \nabla_W J(W; A) = -(A - W). \quad (2.7)$$

From the gradient in (2.7), we obtain the following adaptive gradient descent algorithm:

$$A_k = \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^T = A_{k-1} + \frac{1}{k} (\mathbf{x}_k \mathbf{x}_k^T - A_{k-1}). \quad (2.8)$$

for $\eta_k=1/k$. The ODE corresponding to (2.8) is:

$$dW/dt = A - W, \quad (2.9)$$

which is a Riccati differential equation (2.9) with $F=-I/2$, $Q=A$, $G=0$, and Hamiltonian matrix H :

$$H = \begin{bmatrix} -I/2 & 0 \\ A & I/2 \end{bmatrix} = \begin{bmatrix} 0 & I \\ I & A \end{bmatrix} \begin{bmatrix} -I/2 & 0 \\ 0 & I/2 \end{bmatrix} \begin{bmatrix} 0 & I \\ I & A \end{bmatrix}^{-1}. \quad (2.10)$$

We have $R = (W_0 - A)^{-1}$, where W_0 is the initial value of W at $t=0$. From (2.33), we obtain:

$$W(t, W_0) = A + (W_0 - A)e^{-t} \rightarrow A \text{ as } t \rightarrow \infty. \quad (2.11)$$

The domain of attraction $D(W^*) = \{\mathbb{R}^{n \times n}, A=A^T\}$. The time constant is 1, and the rate of convergence is independent of the eigen-structure of A .

2.3 Adaptive Normalized Mean Algorithm

The most obvious choice for adaptive normalized mean algorithm is to use (2.3) and normalize each \mathbf{m}_k . However, a more efficient algorithm can be obtained from the following cost function whose minimizer \mathbf{w}^* is the asymptotic normalized mean $\mathbf{m}/\|\mathbf{m}\|$, where $\mathbf{m} = \lim_{k \rightarrow \infty} E[\mathbf{x}_k]$:

$$J(\mathbf{w}_k; \mathbf{x}_k) = \|\mathbf{x}_k - \mathbf{w}_k\|^2 + \alpha (\mathbf{w}_k^T \mathbf{w}_k - 1), \quad (2.12)$$

where α is a Lagrange multiplier that enforces the constraint that the mean is normalized. The gradient of $J(\mathbf{w}_k; \mathbf{x}_k)$ with respect to \mathbf{w}_k is:

$$(1/2) \nabla_{\mathbf{w}_k} J(\mathbf{w}_k; \mathbf{x}_k) = -(\mathbf{x}_k - \mathbf{w}_k) + \alpha \mathbf{w}_k. \quad (2.13)$$

Multiplying (2.13) by \mathbf{w}_k^T , and applying the constraint $\mathbf{w}_k^T \mathbf{w}_k = 1$, we obtain:

$$\alpha = \mathbf{w}_k^T \mathbf{x}_k - 1. \quad (2.14)$$

Using this α in (2.14), we obtain the adaptive gradient descent algorithm for normalized mean:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta_k (\mathbf{x}_k - \mathbf{w}_k^T \mathbf{x}_k \mathbf{w}_k), \quad (2.15)$$

where η_k follows assumption A1.2. The ODE corresponding to (2.15) is:

$$d\mathbf{w}/dt = \mathbf{m} - \mathbf{w}^T \mathbf{m} \mathbf{w}. \quad (2.16)$$

Theorem 2.3. For the ODE (2.16), the point $\mathbf{w}^* = \mathbf{m}/\|\mathbf{m}\|$ is (uniformly) asymptotically stable. The domain of attraction of \mathbf{w}^* is $D(\mathbf{w}^*) = \{\mathbb{R}^n\}$.

Proof. Define $v = \mathbf{m}^T \mathbf{w}$. From (2.16), we get $dv/dt = \|\mathbf{m}\|^2 - v^2$, whose solution is:

$$v(t) = \|\mathbf{m}\| \left[\frac{(\|\mathbf{m}\| + v(0)) - (\|\mathbf{m}\| - v(0))e^{-2\|\mathbf{m}\|t}}{(\|\mathbf{m}\| + v(0)) + (\|\mathbf{m}\| - v(0))e^{-2\|\mathbf{m}\|t}} \right] \rightarrow \|\mathbf{m}\| \text{ as } t \rightarrow \infty. \quad (2.17)$$

The domain of attraction $D(\mathbf{w}^*)$ is:

$$D(\mathbf{w}^*) = \left\{ \left(\|\mathbf{m}\| + v(0) \right) e^{2\|\mathbf{m}\|t} + \left(\|\mathbf{m}\| - v(0) \right) \neq 0 \quad \forall t \geq 0 \right\}.$$

Since $(v(0) - \|\mathbf{m}\|) < (v(0) + \|\mathbf{m}\|)$, the above inequality is valid for all $t \geq 0$. Thus, the domain of attraction $D(\mathbf{w}^*) = \{\mathfrak{R}^n\}$. The time constant is $1/(2\|\mathbf{m}\|)$. Clearly, $v = \mathbf{m}^T \mathbf{w} \rightarrow \|\mathbf{m}\|$ as $t \rightarrow \infty$.

We next define $u = \mathbf{w}^T \mathbf{w}$. Then, from (2.16) and $\mathbf{m}^T \mathbf{w} \rightarrow \|\mathbf{m}\|$, we obtain $du/dt = 2\|\mathbf{m}\|(1-u)$. The solution of this ODE is:

$$u(t) = 1 + (u(0) - 1)e^{-2\|\mathbf{m}\|t} \rightarrow 1 \text{ as } t \rightarrow \infty$$

The domain of attraction $D(\mathbf{w}^*) = \{\mathfrak{R}^n\}$. The time constant is $1/(2\|\mathbf{m}\|)$. Clearly, $u = \mathbf{w}^T \mathbf{w} \rightarrow 1$ as $t \rightarrow \infty$. Considering that $\mathbf{m}^T \mathbf{w} \rightarrow \|\mathbf{m}\|$ and $\mathbf{w}^T \mathbf{w} \rightarrow 1$ as $t \rightarrow \infty$, we conclude $\mathbf{w}(t) \rightarrow \mathbf{m}/\|\mathbf{m}\|$ as $t \rightarrow \infty$. ■

In order to prove the convergence of the adaptive algorithm (2.15), we also need to prove that \mathbf{w}_k is bounded. For this, we determine an upper bound of η_k such that \mathbf{w}_k is bounded for all k . The following Theorem gives the result without proof.

Theorem 2.4. For the adaptive algorithm (2.39) let A2.1 and A2.2 hold. Then there exists a uniform upper bound of η_k such that \mathbf{w}_k is almost surely uniformly bounded. Furthermore, if $\|\mathbf{x}_k\| \leq \alpha$ (Assumption A2.1) and θ is the almost sure upper bound of $\|\mathbf{w}_k\|$, then $\|\mathbf{w}_{k+1}\| \leq \|\mathbf{w}_k\|$ if:

$$\eta_k \leq \frac{2}{\alpha\theta}. \quad \blacksquare$$

The convergence of (2.15) is now a direct corollary of the above theorems and an application of Lemma 2.1

2.4 Convergence Proof of Adaptive Median

Given a sequence $\{\mathbf{x}_k\}$, its asymptotic median μ satisfies the following:

$$\lim_{k \rightarrow \infty} P(\mathbf{x}_k \geq \mu) = \lim_{k \rightarrow \infty} P(\mathbf{x}_k < \mu) = 0.5, \quad (2.18)$$

where $P(E)$ is the probability measure of event E , and $0 \leq P(E) \leq 1$.

Adaptive median algorithm:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta_k \operatorname{sgn}(\mathbf{x}_k - \mathbf{w}_k), \quad (2.19)$$

The ODE corresponding to (2.19) is:

$$\frac{d\mathbf{w}(t)}{dt} = \lim_{k \rightarrow \infty} E[\operatorname{sgn}(\mathbf{x}_k - \mathbf{w})]. \quad (2.20)$$

Clearly, a stationary point \mathbf{w}^* of (2.20) is:

$$\lim_{k \rightarrow \infty} P(\mathbf{x}_k \geq \mathbf{w}^*) = \lim_{k \rightarrow \infty} P(\mathbf{x}_k < \mathbf{w}^*), \quad (2.21)$$

which satisfies the condition (2.18) of asymptotic median μ . Besides, we can consider the objective function $J(\mathbf{w}_k; \mathbf{x}_k)$

$$J(\mathbf{w}_k; \mathbf{x}_k) = \|\mathbf{x}_k - \mathbf{w}_k\| \quad (2.22)$$

as an energy function, and it has been proven [Bickel&Doksum 77] that this energy function is minimized for $\mathbf{w}^* = \mu$.

2.5 Brief Review of Optimization Theory

Most of the methods described in this section are from [Luenberger 84].

Unconstrained Minimization Problem

We first consider the unconstrained optimization problem:

$$\text{Minimize } f(\mathbf{x}) \text{ subject to } \mathbf{x} \in \mathfrak{R}^n. \quad (2.23)$$

Definition 2.1. If a real-valued function f is continuous on \mathfrak{R} , and has continuous partial derivatives of order p , we write $f \in C^p$.

Definition 2.2. A point $\mathbf{x}^* \in \mathfrak{R}^n$ is said to be a *local minimum point* of f over \mathfrak{R}^n if there is an $\varepsilon > 0$ such that $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for all $\mathbf{x} \in \mathfrak{R}^n$ within a distance ε of \mathbf{x}^* . If $f(\mathbf{x}) > f(\mathbf{x}^*)$ for all $\mathbf{x} \in \mathfrak{R}^n$, $\mathbf{x} \neq \mathbf{x}^*$, within a distance ε of \mathbf{x}^* , then \mathbf{x}^* is said to be the *strict local minimum point* of f over \mathfrak{R}^n .

Theorem 2.5. (First Order Necessary Conditions). Let $f \in C^1$ be a function of \mathbf{x} on \mathfrak{R}^n . If $\mathbf{x}^* \in \mathfrak{R}^n$ is a local minimum point of f over \mathfrak{R} , then $\nabla_{\mathbf{x}} f(\mathbf{x}^*) = 0$.

Theorem 2.6. (Second Order Sufficient Conditions). Let $f \in C^2$ be a function of \mathbf{x} on \mathfrak{R}^n , and let $\mathbf{x}^* \in \mathfrak{R}^n$. Suppose in addition that (1) $\nabla_{\mathbf{x}} f(\mathbf{x}^*) = 0$, and (2) $\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*)$ is positive definite, then \mathbf{x}^* is the strict local minimum point of f .

We next describe the gradient-based methods for iteratively solving the unconstrained minimization problem (2.23).

1. Gradient Descent Method: $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla_{\mathbf{x}} f(\mathbf{x}_k)$,
where $\alpha > 0$ is a small scalar constant.
2. Steepest Descent Method:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla_{\mathbf{x}} f(\mathbf{x}_k),$$

where the scalar $\alpha_k \geq 0$ is obtained by minimizing $f(\mathbf{x}_k - \alpha \nabla_{\mathbf{x}} f(\mathbf{x}_k))$ for α .

3. Conjugate Direction Method:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k,$$

where $\alpha_k \geq 0$ is obtained by minimizing $f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ for α , and

$$\mathbf{d}_{k+1} = -\nabla_{\mathbf{x}} f(\mathbf{x}_{k+1}) + \beta_k \mathbf{d}_k,$$

where β_k is obtained by one of several methods, where $\mathbf{g}_k = \nabla_{\mathbf{x}} f(\mathbf{x}_k)$:

$$\text{Hestenes-Stiefel: } \beta_k = \mathbf{g}_{k+1}^T (\mathbf{g}_{k+1} - \mathbf{g}_k) / \mathbf{d}_k^T (\mathbf{g}_{k+1} - \mathbf{g}_k),$$

$$\text{Polak-Ribiere: } \beta_k = \mathbf{g}_{k+1}^T (\mathbf{g}_{k+1} - \mathbf{g}_k) / \mathbf{g}_k^T \mathbf{g}_k,$$

$$\text{Fletcher-Reeves: } \beta_k = \mathbf{g}_{k+1}^T \mathbf{g}_{k+1} / \mathbf{g}_k^T \mathbf{g}_k,$$

$$\text{Powell: } \beta_k^i = \max[0, \mathbf{g}_{k+1}^T (\mathbf{g}_{k+1} - \mathbf{g}_k) / \mathbf{g}_k^T \mathbf{g}_k].$$

This method accelerates the typically slow convergence of the gradient or steepest descent methods.

4. Newton's Method: $\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla_{\mathbf{xx}}^2 f(\mathbf{x}_k)]^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}_k).$

Here we assume that the Hessian $H_k = \nabla_{\mathbf{xx}}^2 f(\mathbf{x}_k)$ is nonsingular.

5. Quasi-Newton Method:

Although there are a variety of quasi-Newton methods [Luenberger 84], we consider a simple case where the Hessian H_k is recursively updated by a rule like

$$A_k^{-1} = \frac{k}{\beta(k-1)} \left(A_{k-1}^{-1} - \frac{A_{k-1}^{-1} \mathbf{x}_k \mathbf{x}_k^T A_{k-1}^{-1}}{\beta(k-1) + \mathbf{x}_k^T A_{k-1}^{-1} \mathbf{x}_k} \right)$$

for a new sample \mathbf{x}_k . We use the Sherman-Morrison formula in Section 2.11 to create an update rule for H_k^{-1} similar to those in Section 2.7. We then replace H_k^{-1} in Newton's method by this update rule.

Constrained Minimization Problem

We consider the constrained optimization problem:

$$\text{Minimize } f(\mathbf{x}) \in C^2 \text{ subject to } \mathbf{h}(\mathbf{x}) = \mathbf{0}, \quad (2.24)$$

where and $\mathbf{h} = (h_1, h_2, \dots, h_m) \in C^2$, and $\mathbf{x} \in \mathbb{R}^n$.

1. Lagrange Multiplier Method:

We introduce the Lagrangian associated with the constrained problem as:

$$l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}), \text{ where } \boldsymbol{\lambda} \in \mathbb{R}^m, \quad (2.25)$$

which reduces the constrained problem to an unconstrained one. By the first order necessary conditions (Theorem 2.5), if $\mathbf{x}^* \in \mathfrak{R}^n$ is a local minimum point of f over \mathfrak{R} subject to $\mathbf{h}(\mathbf{x})=\mathbf{0}$, then

$$\nabla_{\mathbf{x}} l(\mathbf{x}^*, \boldsymbol{\lambda}) = \nabla_{\mathbf{x}} f(\mathbf{x}^*) + \boldsymbol{\lambda}^T \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}^*) = \mathbf{0} \text{ and } \nabla_{\boldsymbol{\lambda}} l(\mathbf{x}^*, \boldsymbol{\lambda}) = \mathbf{h}(\mathbf{x}^*) = \mathbf{0}. \quad (2.26)$$

By the second order sufficient conditions (Theorem 2.6), if we denote by M the tangent plane $M = \{\mathbf{y} : \mathbf{y}^T \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}^*) = \mathbf{0}\}$, then if the matrix:

$$\nabla_{\mathbf{xx}}^2 l(\mathbf{x}^*, \boldsymbol{\lambda}) = \nabla_{\mathbf{xx}}^2 f(\mathbf{x}^*) + \boldsymbol{\lambda}^T \nabla_{\mathbf{xx}}^2 \mathbf{h}(\mathbf{x}^*), \quad (2.27)$$

is positive definite on M , then \mathbf{x}^* is the strict local minimum point of f subject to $\mathbf{h}(\mathbf{x})=\mathbf{0}$.

2. Penalty Method:

Here the unconstrained problem is:

$$\text{Minimize } f(\mathbf{x}) + \mu \|\mathbf{h}(\mathbf{x})\|^2, \text{ for some large } \mu > 0. \quad (2.28)$$

Here $\|\mathbf{h}(\mathbf{x})\|^2$ is a penalty function used in literature.

3. Augmented Lagrangian Method:

This method can be viewed as a combination of the lagrangian and penalty function methods. The unconstrained problem is:

$$\text{Minimize } f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) + \mu \|\mathbf{h}(\mathbf{x})\|^2, \text{ for some large } \mu > 0. \quad (2.29)$$

2.6 Matrix Operations

We use several matrix operations in deriving and analyzing our adaptive algorithms. Given a real, symmetric matrix $A \in \mathfrak{R}^{n \times n}$, where $A = [a_{ij}]$, let the eigen-decomposition of A be $A = \Phi \Lambda \Phi^T$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the diagonal eigenvalue matrix with $\lambda_1 \geq \dots \geq \lambda_n$, and $\Phi \in \mathfrak{R}^{n \times n}$ is the eigenvector matrix. Some matrix operations are:

1. *Transpose*: $\text{Transpose}(A) = A^T$,
2. *Trace*: $\text{tr}[A] = \text{tr}[A^T] = \sum_{i=1}^n a_{ii} = \sum_{i=1}^n \lambda_i$,
3. *Determinant*: $\det(A) = \prod_{i=1}^n \lambda_i$,
4. *Frobenius Norm*: $\|A\|_F^2 = \text{tr}[AA^T] = \text{tr}[A^T A] = \lambda_1^2 + \dots + \lambda_n^2$,
5. *Euclidean Norm*: $\|A\| = \lambda_1$,
6. *Sherman-Morrison Formula*: $[A + \mathbf{a}\mathbf{b}^T]^{-1} = A^{-1} - \frac{A^{-1}\mathbf{a}\mathbf{b}^T A^{-1}}{1 + \mathbf{b}^T A^{-1}\mathbf{a}}$.

For positive definite square matrix, A : $\text{tr}(\log(A)) = \log(\det(A))$

2.7 Derivatives of Matrix Functions

In formulating the adaptive algorithms from a scalar objective function, we use gradients of these functions with respect to a matrix. Some useful gradients are listed below. We assume that $C, D \in \mathfrak{R}^{n \times n}$ are symmetric and positive definite, $A \in \mathfrak{R}^{n \times n}$ is symmetric, and $X \in \mathfrak{R}^{n \times n}$.

1. $\nabla_X \text{tr}[X] = I$,
2. $\nabla_X \text{tr}[e^X] = e^{X^T}$,
3. $\nabla_X \text{tr}[CX] = \nabla_X \text{tr}[X^T C] = C$,
4. $\nabla_X \text{tr}[CX^T] = \nabla_X \text{tr}[XC] = C$,
5. $\nabla_X \text{tr}[XCX^T] = \nabla_X \text{tr}[CX^T X] = \nabla_X \text{tr}[X^T XC] = 2XC$,
6. $\nabla_X \text{tr}[CXX^T] = \nabla_X \text{tr}[X^T CX] = \nabla_X \text{tr}[XX^T C] = 2CX$,
7. $\nabla_X \text{tr}[XX^T XX^T] = \nabla_X \text{tr}[X^T XX^T X] = 4XX^T X$,
8. $\nabla_X \text{tr}[CXCX] = \nabla_X \text{tr}[X^T CX^T C] = 2C^2 X$,
9. $\nabla_X \text{tr}[CXX^T C] = \nabla_X \text{tr}[X^T C^2 X] = 2CX^T C$,
10. $\nabla_X \text{tr}[XCX^T XCX^T] = 4XCX^T XC$,
11. $\nabla_X \text{tr}[X^T CXX^T CX] = 4CXX^T CX$,
12. $\nabla_X \text{tr}[XAX^T C] = \nabla_X \text{tr}[AX^T CX] = \nabla_X \text{tr}[X^T CXA] = 2CXA$,
13. $\nabla_X \text{tr}[X^T XX^T CX] = \nabla_X \text{tr}[XX^T CXX^T] = \nabla_X \text{tr}[X^T CXX^T X] = 2XX^T CX + 2CXX^T X$,
14. $\nabla_X \text{tr}[\log(X^T CX)] = 2CX(X^T CX)^{-1}$,
15. $\nabla_X \text{tr}[\log(X^T CXD)] = 2CXD(X^T CX)^{-1} D^{-1}$.

Proofs of Chapter 3

3.2 Adaptive Square Root Algorithm – Method 1

Let $\{\mathbf{x}_k \in \mathfrak{R}^n\}$ be a sequence of data vectors, whose online data correlation matrix $A_k \in \mathfrak{R}^{n \times n}$ is given by:

$$A_k = \frac{1}{k} \sum_{i=1}^k \beta^{k-i} \mathbf{x}_i \mathbf{x}_i^T. \quad (3.1)$$

Here \mathbf{x}_k is an observation vector at time k , and $0 < \beta \leq 1$ is a forgetting factor used for non-stationary sequences. If the data is stationary, the asymptotic correlation matrix A is:

$$A = \lim_{k \rightarrow \infty} E[A_k]. \quad (3.2)$$

Objective Function

Following the methodology described in Section 2.3, we present the algorithm by first showing an objective function J , whose minimum with respect to matrix W gives us the square root of the asymptotic data correlation matrix A . The objective function is:

$$J(W) = \|A - W^T W\|_F^2. \quad (3.3)$$

The gradient of $J(W)$ with respect to W is:

$$\nabla_W J(W) = -4W(A - W^T W). \quad (3.4)$$

Adaptive Algorithm

From the gradient in (3.4), we obtain the following adaptive gradient descent algorithm:

$$W_{k+1} = W_k - \eta_k (1/4) \nabla_W J(W_k; A_k) = W_k + \eta_k (W_k A_k - W_k W_k^T W_k), \quad (3.5)$$

where η_k follows assumption A2.2 in Chapter 2. There are several models for generating the random sequence $\{A_k\}$ from observations $\{\mathbf{x}_k\}$. We can represent A_k simply as its instantaneous value $\mathbf{x}_k \mathbf{x}_k^T$ or by its recursive formulae in (2.21, 2.23).

Solution of the Ordinary Differential Equation (ODE)

We use Stochastic Approximation Theory described in Chapter 2.4 to prove the convergence of (3.5). From (3.3) and (3.5), we obtain the ODE below:

$$\frac{dW}{dt} = -\frac{1}{4} \nabla_W J(W; A) = WA - WW^T W. \quad (3.6)$$

Let us define $P = W^T W$. Then from (3.6), we obtain:

$$\frac{dP}{dt} = AP + PA - 2P^2. \quad (3.7)$$

which is a Riccati differential equation (2.12) with $F=A$, $G=2I$, $Q=0$. Let $A=\Phi\Lambda\Phi^T$ be the eigen-decomposition of A . Here $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the diagonal eigenvalue matrix with $\lambda_1 \geq \dots \geq \lambda_n > 0$, and $\Phi \in \mathbb{R}^{n \times n}$ is the eigenvector matrix. The Hamiltonian matrix H is:

$$H = \begin{bmatrix} -A & 2I \\ 0 & A \end{bmatrix} = \begin{bmatrix} \Phi\Lambda^{-1} & \Phi \\ \Phi & 0 \end{bmatrix} \begin{bmatrix} \Lambda & 0 \\ 0 & -\Lambda \end{bmatrix} \begin{bmatrix} 0 & \Phi^T \\ \Phi^T & -\Lambda^{-1}\Phi^T \end{bmatrix}. \quad (3.8)$$

We have $R = \Phi^T (P_0^{-1} - A^{-1}) \Phi$, where $P_0 = W_0^T W_0$ is the initial condition of the ODE at $t = 0$. Clearly, R exists if P_0 is nonsingular, and A is nonsingular by assumption A2.1. Using the standard solution [Anderson&Moore 90] for the RDE (3.7), we obtain a trajectory $P(t, P_0)$ from (2.18) as:

$$P(t, P_0) = (A^{-1} + e^{-At} (P_0^{-1} - A^{-1}) e^{-At})^{-1} = W(t, W_0)^T W(t, W_0). \quad (3.9)$$

It follows from (3.9) that:

$$W(t, W_0) = U(t) (A^{-1} + e^{-At} ((W_0^T W_0)^{-1} - A^{-1}) e^{-At})^{-1/2}. \quad (3.10)$$

As $t \rightarrow \infty$, we have $W(t) \rightarrow UA^{1/2}$ where $U(t)$ is an orthonormal matrix, and $A^{1/2} = \Phi \Lambda^{1/2} \Phi^T$, i.e., the *symmetric positive definite* square root of A . The domain of attraction $D(W^*)$ is:

$$D(W^*) = \{ \det((W_0^T W_0)^{-1} + e^{At} A^{-1} e^{At} - A^{-1}) \neq 0 \quad \forall t \geq 0 \}. \quad (3.11)$$

The domain attraction can be further simplified to:

$$\det((W_0^T W_0)^{-1} + (e^{2At} - I) A^{-1}) \neq 0 \quad \forall t \geq 0. \quad (3.12)$$

Since $\det((e^{2At} - I) A^{-1}) \geq 0 \quad \forall t \geq 0$, a sufficient condition for (3.12) is:

$$\det(W_0^T W_0) > 0 \text{ or } \det(W_0) \neq 0 \quad \forall t \geq 0. \quad (3.13)$$

The time constant for (3.10) is A^{-1} . The rate of convergence is dependent on the eigen-structure of A , i.e., larger eigenvalues (Λ) of A lead to faster convergence.

Boundedness of W_k

For the complete convergence proof of (3.5) according to Theorem 2.1 (See Section 2.4.2), we need to find conditions under which W_k is bounded. For this, we determine an upper bound of η_k such that W_k is bounded for all k . From (3.5), we obtain a sufficient condition for $\|W_{k+1}\|_F < \|W_k\|_F$ as:

$$\eta_k < \frac{2 \text{tr}(W_k^T W_k (W_k^T W_k - A_k))}{\text{tr}(W_k^T W_k (W_k^T W_k - A_k)^2)}. \quad (3.14)$$

Thus, if $\|W_k\|_F^2 = \text{tr}(W_k^T W_k) \leq \alpha$ and $\|A_k\| < \beta$ (see Assumption A2.1), then

$$\eta_k < \frac{2}{(\alpha + \beta)} \quad (3.15)$$

Convergence of adaptive algorithm (3.5) is a direct application of Theorem 1.

3.3 Adaptive Square Root Algorithm – Method 2

Objective Function

The objective function $J(W)$, whose minimum with respect to W gives us the square root of A is:

$$J(W) = \|A - WW^T\|_F^2. \quad (3.16)$$

The gradient of $J(W)$ with respect to W is:

$$\nabla_W J(W) = -4(A - WW^T)W. \quad (3.17)$$

Adaptive Algorithm

We obtain the following adaptive gradient descent algorithm for square root of A :

$$W_{k+1} = W_k - \eta_k (1/4) \nabla_W J(W_k; A_k) = W_k + \eta_k (A_k W_k - W_k W_k^T W_k), \quad (3.18)$$

where η_k follows the assumption A2.2 in Chapter 2.

Convergence Analysis

From (3.18), we obtain the ODE below:

$$\frac{dW}{dt} = -\frac{1}{4} \nabla_W J(W; A) = AW - WW^T W. \quad (3.19)$$

Let us define $P = WW^T$. Then from (3.19), we obtain:

$$\frac{dP}{dt} = AP + PA - 2P^2. \quad (3.20)$$

which is the same Riccati differential equation as (3.7). Hence the solution for $P(t)$ in (3.9) applies. The solution for $W(t)$ is:

$$W(t, W_0) = \left(A^{-1} + e^{-At} \left((W_0^T W_0)^{-1} - A^{-1} \right) e^{-At} \right)^{-1/2} U(t). \quad (3.21)$$

As $t \rightarrow \infty$, we have $W(t) \rightarrow A^{1/2} U$ where $U(t)$ is an orthonormal matrix, and $A^{1/2} = \Phi \Lambda^{1/2} \Phi^T$, i.e., the *symmetric positive definite* square root of A . The domain of attraction $D(W^*)$ is the same as (3.11-3.13). The time constant is A^{-1} , and the rate of convergence is dependent on the eigen-structure of A .

3.4 Adaptive Square Root Algorithm – Method 3

Adaptive Algorithm

Following the adaptive algorithms (3.5) and (3.18), we now present an algorithm for the computation of a symmetric positive definite square root of A :

$$W_{k+1} = W_k + \eta_k (A_k - W_k^2), \quad (3.22)$$

where η_k follows the assumption A2.2 in Chapter 2 and W_k is symmetric.

Convergence Analysis

From (3.22), we obtain the ODE below:

$$\frac{dW}{dt} = A - W^2, \quad (3.23)$$

which is a Riccati differential equation (2.12) with $F=0$, $G=I$, $Q=A$. The Hamiltonian matrix H is:

$$H = \begin{bmatrix} 0 & I \\ A & 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{2}\Phi\Lambda^{-\frac{1}{2}} & \frac{1}{2}\Phi\Lambda^{-\frac{1}{2}} \\ \frac{1}{2}\Phi & -\frac{1}{2}\Phi \end{bmatrix} \begin{bmatrix} \Lambda^{\frac{1}{2}} & 0 \\ 0 & -\Lambda^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{2}\Phi\Lambda^{-\frac{1}{2}} & \frac{1}{2}\Phi\Lambda^{-\frac{1}{2}} \\ \frac{1}{2}\Phi & -\frac{1}{2}\Phi \end{bmatrix}^{-1}. \quad (3.24)$$

We have $R = \Phi^T (I + W_0 A^{-\frac{1}{2}})^{-1} (I - W_0 A^{-\frac{1}{2}}) \Phi$, where W_0 is the initial condition of the ODE at $t=0$.

Clearly R exists if $(I + W_0 A^{-\frac{1}{2}})$ is positive definite, i.e., $W_0 \neq -A^{\frac{1}{2}}$. From (2.18), we obtain:

$$W(t, W_0) = \left(I - e^{-A^{\frac{1}{2}}t} \Phi R \Phi^T e^{-A^{\frac{1}{2}}t} \right) \left(I + e^{-A^{\frac{1}{2}}t} \Phi R \Phi^T e^{-A^{\frac{1}{2}}t} \right)^{-1} A^{\frac{1}{2}}. \quad (3.25)$$

In other words,

$$W(t, W_0) = \left(I - e^{-A^{\frac{1}{2}}t} (I + W_0 A^{-\frac{1}{2}})^{-1} (I - W_0 A^{-\frac{1}{2}}) e^{-A^{\frac{1}{2}}t} \right) \left(I + e^{-A^{\frac{1}{2}}t} (I + W_0 A^{-\frac{1}{2}})^{-1} (I - W_0 A^{-\frac{1}{2}}) e^{-A^{\frac{1}{2}}t} \right)^{-1} A^{\frac{1}{2}} \quad (3.26)$$

As $t \rightarrow \infty$, we have $W(t) \rightarrow A^{\frac{1}{2}} = \Phi \Lambda^{\frac{1}{2}} \Phi^T$, i.e., the symmetric positive definite square root of A . The domain of attraction $D(W^*)$ is:

$$D(W^*) = \left\{ W = W^T, W_0 \neq -A^{\frac{1}{2}}, \det(e^{At} + (I + W_0 A^{-\frac{1}{2}})^{-1} (I - W_0 A^{-\frac{1}{2}})) \neq 0 \quad \forall t \geq 0 \right\}. \quad (3.27)$$

The domain of attraction can be simplified as:

$$\det(e^{At} + I + W_0 A^{-\frac{1}{2}}(e^{At} - I)) \neq 0 \quad \forall t \geq 0. \quad (3.28)$$

Since $\det(e^{At} \pm I) \geq 0 \quad \forall t \geq 0$, a sufficient condition for (3.28) is:

$$\det(W_0) > 0 \quad \forall t \geq 0. \quad (3.29)$$

The rate of convergence of (26) is dependent on the eigen-structure of A .

3.5 Adaptive Inverse Square Root Algorithm – Method 1

Objective Function

The objective function $J(W)$, whose minimizer W^* gives us the inverse square root of A is:

$$J(W) = \left\| I - W^T A W \right\|_F^2. \quad (3.30)$$

The gradient of $J(W)$ with respect to W is:

$$\nabla_W J(W) = -4AW(I - W^T AW). \quad (3.31)$$

Adaptive Algorithm

From the gradient in (3.31), we obtain the following adaptive gradient descent algorithm:

$$W_{k+1} = W_k - \eta_k (1/4) A_k^{-1} \nabla_W J(W_k; A_k) = W_k + \eta_k (W_k - W_k W_k^T A_k W_k). \quad (3.32)$$

Convergence Analysis

From (3.31) and (3.32), we obtain the ODE below:

$$\frac{dW}{dt} = -\frac{1}{4} A^{-1} \nabla_W J(W; A) = W - WW^T AW. \quad (3.33)$$

Let us define $P = W^T AW$. Then from (3.33), we obtain:

$$\frac{dP}{dt} = 2P - 2P^2. \quad (3.34)$$

which is a Riccati differential equation (2.12) with $F=I$, $G=2I$, $Q=0$. The Hamiltonian matrix H is:

$$H = \begin{bmatrix} -I & 2I \\ 0 & I \end{bmatrix} = \begin{bmatrix} I & I \\ I & 0 \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & -I \end{bmatrix} \begin{bmatrix} 0 & I \\ I & -I \end{bmatrix}. \quad (3.35)$$

We have $R = (P_0^{-1} - I)$, where $P_0 = W_0^T AW_0$ is the initial condition of the ODE at $t=0$. Clearly, R exists if P_0 is nonsingular. From (2.18), we obtain:

$$P(t, P_0) = \left(I + e^{-\Lambda t} (P_0^{-1} - I) e^{-\Lambda t} \right)^{-1} = W(t, W_0)^T AW(t, W_0) = \left(A^{1/2} W(t, W_0) \right)^T \left(A^{1/2} W(t, W_0) \right). \quad (3.36)$$

It follows from (3.36) that:

$$W(t, W_0) = A^{-1/2} U(t) \left(I + e^{-\Lambda t} \left((W_0^T AW_0)^{-1} - I \right) e^{-\Lambda t} \right)^{-1/2}. \quad (3.37)$$

As $t \rightarrow \infty$, we have $W(t) \rightarrow A^{-1/2} U$ where $U(t)$ is an orthonormal matrix, and $A^{-1/2} = \Phi \Lambda^{-1/2} \Phi^T$, i.e., the symmetric positive definite inverse square root of A . The domain of attraction $D(W^*)$ is:

$$D(W^*) = \left\{ \det(e^{2\Lambda t} - I + (W_0^T AW_0)^{-1}) \neq 0 \quad \forall t \geq 0 \right\}. \quad (3.38)$$

Since $\det(e^{2\Lambda t} - I) \geq 0 \quad \forall t \geq 0$, a sufficient condition for (3.38) is:

$$\det(W_0^T AW_0) > 0 \quad \forall t \geq 0. \quad (3.39)$$

The time constant for (37) is A^{-1} , and the rate of convergence is dependent on the eigen-structure of A .

Boundedness of W_k

For the complete convergence proof of (3.32) according to Theorem 2.1 (See Section 2.4.2), we need to find conditions under which W_k is bounded. For this, we determine an upper bound of

η_k such that W_k is bounded for all k . From (3.32), we obtain a sufficient condition for $\|W_{k+1}\|_F < \|W_k\|_F$ as:

$$\eta_k < \frac{2\text{tr}(W_k^T W_k (W_k^T A_k W_k - I))}{\text{tr}(W_k^T W_k (W_k^T A_k W_k - I)^2)}. \quad (3.40)$$

Thus, if $\|W_k\|_F^2 = \text{tr}(W_k^T W_k) \leq \alpha$ and $\|A_k\| < \beta$ (see Assumption A2.1), then

$$\eta_k < \frac{2}{(\alpha\beta + 1)}. \quad (3.41)$$

The convergence of adaptive algorithm (32) is a direct application of Theorem 1.

3.6 Adaptive Inverse Square Root Algorithm – Method 2

Objective Function

The objective function $J(W)$, whose minimum with respect to W gives us the inverse square root of A is:

$$J(W) = \|I - WAW^T\|_F^2. \quad (3.42)$$

The gradient of $J(W)$ with respect to W is:

$$\nabla_W J(W) = -4(I - WAW^T)WA. \quad (3.43)$$

Adaptive Algorithm

We obtain the following adaptive algorithm for inverse square root of A :

$$W_{k+1} = W_k - \eta_k (1/4) \nabla_W J(W_k; A_k) A_k^{-1} = W_k + \eta_k (W_k - W_k A_k W_k^T W_k), \quad (3.44)$$

where η_k follows the assumption A2.2 in Chapter 2.

Convergence Analysis

From (3.44), we obtain the ODE below:

$$\frac{dW}{dt} = W - WAW^T W. \quad (3.45)$$

Let us define $P = WAW^T$. Then from (3.45), we obtain:

$$\frac{dP}{dt} = 2P - 2P^2. \quad (3.46)$$

which is the a Riccati differential equation as (3.34). Hence the solution for $P(t)$ in (3.36) applies.

The solution for $W(t)$ is:

$$W(t, W_0) = \left(I + e^{-\Lambda t} \left((W_0 A W_0^T)^{-1} - I \right) e^{-\Lambda t} \right)^{-1/2} U(t) A^{-1/2}. \quad (3.47)$$

As $t \rightarrow \infty$, we have $W(t) \rightarrow U A^{-1/2}$ where $U(t)$ is an orthonormal matrix, and $A^{-1/2} = \Phi \Lambda^{-1/2} \Phi^T$. The domain of attraction $D(W^*)$ is the same as (3.38).

The time constant for (47) is A^{-1} , and the rate of convergence is dependent on the eigen-structure of A .

3.7 Adaptive Inverse Square Root Algorithm – Method 3

Adaptive Algorithm

By extending the adaptive algorithms (3.32) and (3.44), we now present an adaptive algorithm for the computation of a symmetric positive definite inverse square root of A :

$$W_{k+1} = W_k + \eta_k (I - W_k A W_k), \quad (3.48)$$

where η_k follows the assumption A2.2 in Chapter 2 and W_k is symmetric.

Convergence Analysis

From (3.46), we obtain the ODE below:

$$\frac{dW}{dt} = I - W A W, \quad (3.49)$$

which is a Riccati differential equation (2.12) with $F=0$, $G=A$, $Q=I$. The Hamiltonian matrix H as:

$$H = \begin{bmatrix} 0 & A \\ I & 0 \end{bmatrix} = \begin{bmatrix} \Phi \Lambda^{1/2} & \Phi \Lambda^{1/2} \\ \Phi & -\Phi \end{bmatrix} \begin{bmatrix} \Lambda^{1/2} & 0 \\ 0 & -\Lambda^{1/2} \end{bmatrix} \begin{bmatrix} 1/2 \Lambda^{-1/2} \Phi^T & 1/2 \Phi^T \\ 1/2 \Lambda^{-1/2} \Phi^T & -1/2 \Phi^T \end{bmatrix}. \quad (3.50)$$

We have $R = \Phi^T (I + W_0 A^{1/2})^{-1} (I - W_0 A^{1/2}) \Phi$, where W_0 is the initial condition of the ODE at $t=0$.

Clearly R exists if $(I + W_0 A^{1/2})$ is positive definite, i.e., $W_0 \neq -A^{-1/2}$. From (2.18), we obtain:

$$W(t, W_0) = \left(I - e^{-A^{1/2} t} \Phi R \Phi^T e^{-A^{1/2} t} \right) \left(I + e^{-A^{1/2} t} \Phi R \Phi^T e^{-A^{1/2} t} \right)^{-1} A^{-1/2}. \quad (3.51)$$

In other words,

$$W(t, W_0) = \left(I - e^{-A^{1/2} t} (I + W_0 A^{1/2})^{-1} (I - W_0 A^{1/2}) e^{-A^{1/2} t} \right) \left(I + e^{-A^{1/2} t} (I + W_0 A^{1/2})^{-1} (I - W_0 A^{1/2}) e^{-A^{1/2} t} \right)^{-1} A^{-1/2} \quad (3.52)$$

As $t \rightarrow \infty$, we have $W(t) \rightarrow A^{-1/2} = \Phi \Lambda^{-1/2} \Phi^T$, i.e., the symmetric positive definite inverse square root of A . The domain of attraction $D(W^*)$ is:

$$D(W^*) = \left\{ W = W^T, W_0 \neq -A^{-1/2}, \det((e^{At} + I) + W_0 A^{1/2} (e^{At} - I)) \neq 0 \quad \forall t \geq 0 \right\}. \quad (3.53)$$

Once again, since $\det(e^{At} \pm I) \geq 0 \quad \forall t \geq 0$, a sufficient condition for (3.53) is:

$$\det(W_0) > 0 \quad \forall t \geq 0. \quad (3.54)$$

The rate of convergence of (52) is dependent on the eigen-structure of A .

Proofs of Chapter 4

Theorem 4.3. Let A4.1 and A4.2 hold, and let $\mathbf{w}(t) = \sum_{i=1}^n a_i(t)\phi_i$ be solutions for the ODEs (4.16)-(4.18) in terms of the entire orthonormal set of eigenvectors $\{\phi_1, \phi_2, \dots, \phi_n\}$ of A . Then for any initial condition $\mathbf{w}(0) = \mathbf{w}_0 \in \mathbb{R}^n$, $\mathbf{w}_0^T \phi_1 \neq 0$, the solutions for the coefficients $a_i(t)$ for $t \geq 0$ for the RQ, OJAN, and LUO algorithms are:

$$\text{RQ:} \quad \frac{a_i(t)}{a_1(t)} = \frac{a_i(0)}{a_1(0)} e^{\frac{-(\lambda_1 - \lambda_i)t}{\|\mathbf{w}_0\|^2}} \quad \text{for } i=2, \dots, n \quad (4.23)$$

$$\text{OJAN:} \quad \frac{a_i(t)}{a_1(t)} = \frac{a_i(0)}{a_1(0)} e^{-(\lambda_1 - \lambda_i)t} \quad \text{for } i=2, \dots, n \quad (4.24)$$

$$\text{LUO:} \quad \frac{a_i(t)}{a_1(t)} = \frac{a_i(0)}{a_1(0)} e^{-\|\mathbf{w}_0\|^2 (\lambda_1 - \lambda_i)t} \quad \text{for } i=2, \dots, n \quad (4.25)$$

and

$$\text{RQ:} \quad a_1(t) = \pm \sqrt{\frac{c_0 \|\mathbf{w}_0\|^2}{e^{-2\lambda_1 \|\mathbf{w}_0\|^2 t} + c_0}}, \quad (4.26)$$

$$\text{OJAN:} \quad a_1(t) = \pm \sqrt{\frac{c_0 \|\mathbf{w}_0\|^2}{e^{-2\lambda_1 t} + c_0}}, \quad (4.27)$$

$$\text{LUO:} \quad a_1(t) = \pm \sqrt{\frac{c_0 \|\mathbf{w}_0\|^2}{e^{-2\lambda_1 \|\mathbf{w}_0\|^2 t} + c_0}}, \quad \text{where } c_0 = \frac{a_1(0)^2}{\|\mathbf{w}_0\|^2 - a_1(0)^2}. \quad (4.28)$$

The points $\pm \|\mathbf{w}_0\| \phi_1$ are (uniformly) asymptotically stable. The domain of attraction of $\|\mathbf{w}_0\| \phi_1$ is $D(\|\mathbf{w}_0\| \phi_1) = \{\mathbf{w} \in \mathbb{R}^n \mid \mathbf{w}^T \phi_1 > 0\}$ and that of $-\|\mathbf{w}_0\| \phi_1$ is $D(-\|\mathbf{w}_0\| \phi_1) = \{\mathbf{w} \in \mathbb{R}^n \mid \mathbf{w}^T \phi_1 < 0\}$.

Proof. We show the result for ODE (4.21) and the remaining ODEs are similar. Let $\mathbf{w}(t) = \sum_{i=1}^n a_i(t)\phi_i$ be solutions for the ODE (4.21) in terms of the entire orthonormal set of eigenvectors $\{\phi_1, \phi_2, \dots, \phi_n\}$ of A . Substituting this $\mathbf{w}(t)$ in (4.21) and multiplying on the left by ϕ_i^T , we obtain:

$$\frac{da_i}{dt} = \frac{1}{\|\mathbf{w}_0\|^2} \left(a_i \lambda_i - a_i \frac{\sum_{i=1}^n a_i^2 \lambda_i}{\|\mathbf{w}_0\|^2} \right). \quad (4.A.1)$$

Let $b_i = a_i / a_1$ for $i=2, \dots, n$. Then, we can write the (4.A.1) as:

$$\frac{db_i}{dt} = -b_i \left(\frac{\lambda_1 - \lambda_i}{\|\mathbf{w}_0\|^2} \right), \quad (4.A.2)$$

which gives us the solution in (4.23). Since we know from (4.23) that $a_i(t) \rightarrow 0$ for $i=2, \dots, n$ as $t \rightarrow \infty$, we obtain from (4.A.1):

$$\frac{da_1}{dt} = \frac{1}{\|\mathbf{w}_0\|^2} \left(a_1 \lambda_1 - \frac{a_1^3 \lambda_1}{\|\mathbf{w}_0\|^2} \right), \quad (4.A.3)$$

which gives us the solution in (4.26). If $a_1(0) \neq 0$ then $a_1(t) \neq 0$ for all t . Thus the sign of $a_1(t)$ is determined by the sign of $a_1(0) = \mathbf{w}(0)^T \boldsymbol{\phi}_1$. ■

Theorem 4.5. Let A4.1 hold. Then there exists a uniform upper bound for η_k such that \mathbf{w}_k is uniformly bounded w.p.1. Furthermore, if $\|\mathbf{w}_k\|^2 \leq \alpha + 1$, then $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2$ if:

$$\eta_k < \frac{2(\alpha + 1)}{\alpha}. \quad (4.44)$$

Proof. Let ρ be the principal eigenvector of A_k and θ the corresponding (largest) eigenvalue. We multiply (4.40) on the left by ρ^T and we define $\mathbf{v}_k = \rho^T \mathbf{w}_k$. Then from (4.40), we get:

$$\mathbf{v}_{k+1} = \mathbf{v}_k + \eta_k \left(\frac{\theta \mathbf{v}_k}{\mathbf{w}_k^T A_k \mathbf{w}_k} - \mathbf{v}_k \right). \quad (4.A.4)$$

Taking the norm of the above equation and using $\|\mathbf{v}_{k+1}\|^2 \leq \|\mathbf{v}_k\|^2$, we get from (4.A.4):

$$\eta_k < 2 \left/ \left(1 - \frac{\theta}{\mathbf{w}_k^T A_k \mathbf{w}_k} \right) \right. \quad (4.A.5)$$

Using $\mathbf{w}_k^T A_k \mathbf{w}_k \leq \theta \|\mathbf{w}_k\|^2$, we get from (4.A.5):

$$\eta_k \leq \frac{2}{1 - \left(\frac{1}{\|\mathbf{w}_k\|^2} \right)}, \text{ which implies } \eta_k < \frac{2(\alpha + 1)}{\alpha}. \quad \blacksquare$$

Theorem 4.13. Let A4.1 and A4.2 hold. The critical points of $E(\mathbf{w})$ in (4.70) are $\mathbf{0}$ and $\pm \boldsymbol{\phi}_i$ ($i=1, \dots, n$). The two critical points $\pm \boldsymbol{\phi}_1$ are global minimum points of $E(\mathbf{w})$, and $E(\mathbf{w})$ has no other local minimum point. The point $\mathbf{0}$ is a local maximum. In addition, $\pm \boldsymbol{\phi}_i$ ($i=2, \dots, n$) are saddle points of $E(\mathbf{w})$.

Proof. Substituting $\mathbf{w} = \sum_{i=1}^n a_i(t) \boldsymbol{\phi}_i$ in the ODE (4.69) (where $\{\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_n\}$ are orthonormal eigenvectors of A), and multiplying it on the left by $\boldsymbol{\phi}_k^T$ ($k=1, \dots, n$), and then equating it to 0, we get:

$$a_k \left(-(2\lambda_k + \mu) + \sum_{r=1}^n a_r^2 (\lambda_r + \lambda_k + \mu) \right) = 0 \quad \text{for } k=1, \dots, n, \quad (4.A.6)$$

which gives us:

$$a_k = 0 \text{ or } \sum_{r=1}^n a_r^2 (\lambda_r + \lambda_k + \mu) = 2\lambda_k + \mu \quad \text{for } k=1, \dots, n. \quad (4.A.7)$$

Rewriting the second equation of (4.A.7), we obtain:

$$\sum_{r=1, r \neq k}^n a_r^2 (\lambda_r + \lambda_k + \mu) + (a_k^2 - 1)(2\lambda_k + \mu) = 0 \quad \text{for } k=1, \dots, n. \quad (4.A.8)$$

Since A is positive definite (i.e., $\lambda_i > 0$ by Assumption A4.1), $\mu > 0$, and since λ_1 unit multiplicity (Assumption A4.2), from (4.A.8) we obtain at most one non-zero $a_k = \pm 1$. Thus, the equilibrium points of $E(\mathbf{w})$ are $\mathbf{w} = d_{(1)} \phi_{(1)}$, where $\phi_{(1)}$ is a permutation of the eigenvectors ϕ_1, \dots, ϕ_n , and $d_{(1)} = 0$ or ± 1 .

Now let us assume that $d_{(1)} = 0$, i.e. $\mathbf{w} = 0$. Then, from (4.70), $E(0) = \mu/2$. Next, we perturb \mathbf{w} by $\delta \phi_1$. We observe that $E(\delta \phi_1) - E(0) = (2\lambda_1 + \mu)(\delta^4 - 2\delta^2) < 0$ for $\delta < \sqrt{2}$. Clearly, the energy function $E(\mathbf{w})$ decreases. We next prove that if $\phi_{(1)} \neq \pm \phi_1$, then the critical points are unstable equilibrium points of $E(\mathbf{w})$. We obtain the Hessian of $E(\mathbf{w})$ with respect to \mathbf{w} as

$$\nabla_{\mathbf{w}}^2 E(\mathbf{w}) = A\mathbf{w}^T \mathbf{w} + 2\mathbf{w}\mathbf{w}^T A + \mathbf{w}^T A \mathbf{w} I + 2A\mathbf{w}\mathbf{w}^T + \mu(\mathbf{w}^T \mathbf{w} - 1)I + 2\mu\mathbf{w}\mathbf{w}^T - 2A. \quad (4.A.9)$$

We note that

$$\nabla_{\mathbf{w}}^2 E(\pm \phi_1) = -A + \lambda_1 I + 4\lambda_1 \phi_1 \phi_1^T + 2\mu \phi_1 \phi_1^T, \quad (4.A.10)$$

whose eigenvectors are ϕ_1, \dots, ϕ_n , and eigenvalues are $2(2\lambda_1 + \mu)$ for ϕ_1 , and $(\lambda_1 - \lambda_r)$ for ϕ_r ($r > 1$). Clearly, $\nabla_{\mathbf{w}}^2 E(\pm \phi_1)$ is positive definite. On the other hand, we observe that

$$\nabla_{\mathbf{w}}^2 E(\pm \phi_r) = -A + \lambda_r I + 4\lambda_r \phi_r \phi_r^T + 2\mu \phi_r \phi_r^T \text{ for } r > 1. \quad (4.A.11)$$

The eigenvectors of $\nabla_{\mathbf{w}}^2 E(\pm \phi_r)$ are ϕ_1, \dots, ϕ_n , and the eigenvalues are $-(\lambda_1 - \lambda_r)$ for ϕ_1 , and $2(2\lambda_r + \mu)$ for ϕ_r ($r > 1$). Clearly, $\nabla_{\mathbf{w}}^2 E(\pm \phi_r)$ for $r > 1$ is an indefinite matrix. Thus, $\pm \phi_1$ is a stable equilibrium point of the energy function $E(\mathbf{w})$, whereas $\pm \phi_r$ for $r > 1$ are unstable equilibrium points. ■

Theorem 5.1. For the ordinary differential equation (5.20), let A5.1 and A5.3 hold. Then $W = \Phi$ PDU are equilibrium points of (5.20), where $D = [D_1 | 0]^T \in \mathbb{R}^{n \times p}$ with $D_1 \in \mathbb{R}^{p \times p}$ is diagonal with elements $d_i = 0$ or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$, $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $U \in \mathbb{R}^{p \times p}$ is an arbitrary rotation matrix, i.e., $U^T U = U U^T = I_p$.

Proof. From (5.20), we need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$A W - \mu W (W^T W - I_p) = 0. \quad (5.21)$$

The trivial solution is $W = 0$. We next assume that $W \neq 0$. Let $W = QDU$ be the singular value decomposition of W , where $Q \in \mathbb{R}^{n \times n}$ and $U \in \mathbb{R}^{p \times p}$ are orthonormal, and $D \in \mathbb{R}^{n \times p}$ is diagonal. Replacing QDU for W in (5.21) and defining $B = Q^T A Q$, we get from (5.21):

$$BD - \mu D(D^T D - I_p) = 0. \quad (5.22)$$

Let $D = [D_1 | 0]^T$ where D_1 is a $p \times p$ diagonal matrix with diagonal elements d_i for $i=1, \dots, p$. Let

$B = \begin{bmatrix} B_1 & B_2 \\ B_2^T & B_3 \end{bmatrix}$ be a partition of B where $B_1 \in \mathbb{R}^{p \times p}$. From (5.22) we get:

$$(B_1 - \mu(D_1^2 - I_p))D_1 = 0 \text{ and } B_2^T D_1 = 0. \quad (5.23)$$

From (5.23), we get:

$$D_1 = 0 \text{ or } B_1 = \mu(D_1^2 - I_p) \text{ and } B_2 = 0. \quad (5.24)$$

From (5.24), we conclude that B_1 is diagonal, and $d_i^2 = 1 + (b_i / \mu)$ for $i=1, \dots, p$, where b_i are the diagonal elements of B_1 . Since $B = Q^T A Q$ is diagonal, where Q is orthonormal ($Q^T Q = Q Q^T = I_n$), the columns of Q are the n orthonormal eigenvectors of A , i.e., $Q = \Phi P$, where $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix. Also, $B = P^T \Lambda P$, i.e., a permutation of Λ . Therefore, $b_{(i)} = \lambda_{(i)}$ for $i=1, \dots, p$. Combining all results, we get $W = \Phi P D U$, where $D = [D_1 | 0]^T \in \mathbb{R}^{n \times p}$ and $D_1 \in \mathbb{R}^{p \times p}$ is a diagonal matrix with elements $d_i = 0$ or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$ for $i=1, \dots, p$. ■

Theorem 5.2. Let A5.1 and A5.3 hold. Then $W = \Phi D U$, where $D = [D_1 | 0]^T \in \mathbb{R}^{n \times p}$, $D_1 = \text{diag}(d_1, \dots, d_p) \in \mathbb{R}^{p \times p}$, $d_i = \pm \sqrt{1 + (\lambda_i / \mu)}$ for $i=1, \dots, p$, and $U \in \mathbb{R}^{p \times p}$ is an arbitrary rotation matrix, are stable equilibrium points of the ODE (5.20) and strict global minimum points of the objective function (5.18). In addition, $W = \Phi P D U$, where $d_i = 0$ for $i \leq p$ or $P \neq I$, are unstable equilibrium points of the ODE (5.20).

Proof. From (5.18), the energy function $E(W)$ for the PF Homogeneous adaptive algorithm is:

$$E(W) = -\text{tr}(W^T A W) + \frac{\mu}{2} \text{tr}((W^T W - I_p)^2). \quad (5.25)$$

From Theorem 5.1, $W = \Phi P D U = \Psi U$, where $\Psi = \Phi P D$. Then from (5.A.5):

$$E(W) = \sum_{i=1}^p \left(-\Psi_i^T A \Psi_i + \mu \left(\sum_{j=1, j \neq i}^p (\Psi_i^T \Psi_j)^2 + (\Psi_i^T \Psi_i - 1)^2 \right) \right).$$

Here $\Psi_i = d_{(i)} \phi_{(i)}$ is the i^{th} column of Ψ for $i=1, \dots, p$, and $d_{(i)} = 0$ or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$. We first prove that $d_{(i)} = 0$ is an unstable equilibrium point of $E(W)$. We perturb Ψ_i by $\delta \phi_i$, for $\delta > 0$. Then

$$E(\Psi_i = \delta \phi_i) - E(\Psi_i = 0) = -\delta^2 \lambda_i + \mu(\delta^2 - 1)^2 < 0$$

$$\text{for } (\lambda + 2\mu - \sqrt{\lambda(\lambda + 4\mu)}) / 2\mu < \delta^2 < (\lambda + 2\mu + \sqrt{\lambda(\lambda + 4\mu)}) / 2\mu,$$

i.e., the energy function decreases. We next prove that $\Psi_i = d_r \phi_r$, $d_r = \pm \sqrt{1 + (\lambda_r / \mu)}$, ($r > p$) is an unstable equilibrium point of $E(W)$. We perturb Ψ_i by $\delta \phi_p$, i.e., $\Psi_i = d_r (\phi_r + \delta \phi_p) / \sqrt{(1 + \delta^2)}$. Then

$$E\left(d_r(\phi_r + \delta\phi_p)/\sqrt{(1+\delta^2)}\right) - E(d_r\phi_r) = -\delta^2 d_r^2 (\lambda_p - \lambda_r)/(1+\delta^2) < 0.$$

We next prove that $\psi_i = d_r \phi_r$, $d_r = \pm \sqrt{1 + (\lambda_r / \mu)}$, ($r \leq p$) is a stable equilibrium point of $E(W)$. We perturb ψ_i by $\delta\phi_s$, i.e., $\psi_i = d_r(\phi_r + \delta\phi_s)/\sqrt{(1+\delta^2)}$ for $s > p$. Then

$$E\left(d_r(\phi_r + \delta\phi_s)/\sqrt{(1+\delta^2)}\right) - E(d_r\phi_r) = \delta^2 d_r^2 (\lambda_r - \lambda_s)/(1+\delta^2) > 0,$$

i.e., the energy function increases. Thus, the columns of Ψ consisting of the first p orthonormal eigenvectors ϕ_i of A scaled by d_i are the stable minimum points of $E(W)$. ■

Theorem 5.3. *Let A5.1 and A5.3 hold. Then, all the equilibrium points of the ODE (5.28) are up to an arbitrary permutation of the eigenvectors of A weighted by 0 or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$, i.e., any point $W = [d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0$ or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$ is an equilibrium point of the ODE (5.28).*

Proof. We need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$AW - \mu W U^T (W^T W - I_p) = 0. \quad (5.29)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. We multiply (5.29) on the left by W^T , and define $G = W^T A W + (\mu/2) W^T W$, and $H = W^T W$. From (5.29), we obtain:

$$G = H U^T(H). \quad (5.30)$$

Since G is symmetric, $H U^T(H)$ is also symmetric. Since W is assumed to be nonzero, H has positive diagonal elements. From (5.30), we obtain that H and G are diagonal, i.e., both $W^T W$ and $W^T A W$ are diagonal. Then, W is of the form $W = \Phi P D$, where $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $D = [D_1 \mid 0]^T \in \mathbb{R}^{n \times p}$ where $D_1 \in \mathbb{R}^{p \times p}$ is a diagonal matrix with elements d_i for $i=1, \dots, p$. Substituting $W = \Phi P D$ in (5.29), we obtain $\Lambda_1 D_1 = \mu D_1 (D_1^2 - I_p)$, where Λ_1 is a $p \times p$ partition of the permuted eigenvector matrix Λ of A . We obtain:

$$D_1 = 0 \text{ or } D_1^2 = I_p + (\Lambda_1 / \mu). \quad (5.31)$$

From the second equation in (5.31), we obtain $d_i^2 = 1 + (\lambda_{(i)} / \mu)$. Combining all results, we conclude that $W = \Phi P D$, where $d_i = 0$ or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$. ■

Theorem 5.4. *Let A5.1 and A5.3 hold. Then, the points $W^* = [d_1 \phi_1 \ d_2 \phi_2 \ \dots \ d_p \phi_p]$, where $d_i = \pm \sqrt{1 + (\lambda_i / \mu)}$, are the strict global minimum points of the objective function (5.20) and stable equilibrium points of the ODE (5.28). In addition, the points $W = [d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0$ or $\phi_{(i)} \neq \phi_j$ for $i \in \{1, 2, \dots, p\}$ are unstable equilibrium points of the ODE (5.28).*

Proof. From (5.20), the energy function $E(W)$ for the PF Deflation adaptive algorithm is:

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where } E_i(\mathbf{w}_i) = -\mathbf{w}_i^T A \mathbf{w}_i + \mu \left(\sum_{j=1}^{i-1} (\mathbf{w}_j^T \mathbf{w}_i)^2 + \frac{1}{2} (\mathbf{w}_i^T \mathbf{w}_i - 1)^2 \right).$$

From Theorem 5.3, $\mathbf{w}_i = d_{(i)} \phi_{(i)}$, $d_{(i)} = 0, \pm \sqrt{1 + (\lambda_{(i)} / \mu)}$, is the i^{th} column of W for $i=1, \dots, p$. We first prove that $d_{(i)}=0$ is an unstable equilibrium point of $E(W)$. We perturb \mathbf{w}_i by $\delta \phi_i$, for $\delta > 0$. Then

$$E(\mathbf{w}_i = \delta \phi_i) - E(\mathbf{w}_i = 0) = -\delta^2 \lambda_i + (\mu/2)(\delta^2 - 1)^2 < 0$$

$$\text{for } (\lambda + \mu - \sqrt{\lambda(\lambda + 2\mu)}) / \mu < \delta^2 < (\lambda + \mu + \sqrt{\lambda(\lambda + 2\mu)}) / \mu.$$

We next prove that $\mathbf{w}_i = d_i \phi_i$, $d_i = \pm \sqrt{1 + (\lambda_i / \mu)}$, ($i \leq p$) is an unstable equilibrium point of $E(W)$. If $\mathbf{w}_i \neq d_i \phi_i$ then there exists a pair of columns $\mathbf{w}_r = d_{(r)} \phi_{(r)}$, and $\mathbf{w}_s = d_{(s)} \phi_{(s)}$, such that $r < s \leq p$ and $\lambda_{(r)} < \lambda_{(s)}$. We perturb \mathbf{w}_r by $\delta \phi_{(s)}$, i.e., $\mathbf{w}_r = d_{(r)} (\phi_{(r)} + \delta \phi_{(s)}) / \sqrt{1 + \delta^2}$. Then

$$E(\mathbf{w}_r = d_{(r)} (\phi_{(r)} + \delta \phi_{(s)}) / \sqrt{1 + \delta^2}) - E(\mathbf{w}_r = d_{(r)} \phi_{(r)}) = -\delta^2 d_{(r)}^2 (\lambda_{(s)} - \lambda_{(r)}) / (1 + \delta^2) < 0.$$

We next prove that $\mathbf{w}_i = d_i \phi_i$, $d_i = \pm \sqrt{1 + (\lambda_i / \mu)}$, ($i \leq p$) is a stable equilibrium point of $E(W)$. We perturb \mathbf{w}_i by $\delta \phi_s$, i.e., $\mathbf{w}_i = d_i (\phi_i + \delta \phi_s) / \sqrt{1 + \delta^2}$ for $s > i$. Then

$$E(\mathbf{w}_i = d_i (\phi_i + \delta \phi_s) / \sqrt{1 + \delta^2}) - E(\mathbf{w}_i = d_i \phi_i) = \delta^2 d_i^2 (\lambda_i - \lambda_s) / (1 + \delta^2) > 0,$$

i.e., the energy function increases. Thus, $W = [d_1 \phi_1 \dots d_p \phi_p]$, $d_i = \pm \sqrt{1 + (\lambda_i / \mu)}$, ($i \leq p$) are the stable equilibrium points of $E(W)$. ■

Theorem 5.5. Let A5.1 and A5.3 hold. Then, all the equilibrium points of the ODE (5.34) are up to an arbitrary permutation of the eigenvectors of A weighted by 0 or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$, i.e., any point $W = [d_{(1)} \phi_{(1)} \dots d_{(p)} \phi_{(p)}]$, where $d_{(i)} = 0$ or $\pm \sqrt{1 + (\lambda_{(i)} / \mu)}$, is an equilibrium point of the ODE (5.34).

Proof. We need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$A W C - \mu W C (W^T W - I_p) = 0. \quad (5.35)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. We multiply (5.35) on the left by W^T , and define $G = W^T A W + (\mu/2) W^T W$, and $H = W^T W$. From (5.35), we obtain:

$$G C = \mu H C H. \quad (5.36)$$

Since H , G , and $H C H$ are symmetric matrices, we conclude that $G C = C G$. Since C is diagonal with distinct diagonal elements, G is diagonal. Let $G C = C G = D \in \mathbb{R}^{p \times p}$ be a diagonal matrix. Then, from (5.36), we get:

$$D = \mu H C H. \quad (5.37)$$

Since W is assumed to be nonzero, the diagonal elements of $H=W^TW$ are positive. We conclude from (5.37) that H is diagonal. The rest of the proof is similar to Theorem 5.3 above. ■

Theorem 5.6. Let A5.1 and A5.3 hold. Then, the points $W^*=[d_1\phi_1 \ d_2\phi_2 \ \dots \ d_p\phi_p]$, where $d_i = \pm \sqrt{1+(\lambda_i/\mu)}$, are stable equilibrium points of the ODE (5.34) and strict global minimum points of the objective function in (5.32). In addition, the points $W=[d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0$ or $\phi_{(i)} \neq \phi_i$ for $i \in \{1,2,\dots,p\}$ are unstable equilibrium points of the ODE (5.34).

Proof. From (5.32), the energy function $E(W)$ for the PF Weighted adaptive algorithm is:

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where } E_i(\mathbf{w}_i) = -c_i \mathbf{w}_i^T A \mathbf{w}_i + \mu \left(\sum_{j=1, j \neq i}^p c_j (\mathbf{w}_j^T \mathbf{w}_i)^2 + \frac{c_i}{2} (\mathbf{w}_i^T \mathbf{w}_i - 1)^2 \right).$$

From Theorem 5.5, $\mathbf{w}_i = d_{(i)}\phi_{(i)}$, $d_{(i)} = 0$ or $\pm \sqrt{1+(\lambda_{(i)}/\mu)}$, is the i^{th} column of W for $i=1,\dots,p$. We first prove that $d_{(i)}=0$ is an unstable equilibrium point of $E(W)$. We perturb \mathbf{w}_i by $\delta\phi_i$. Then

$$\begin{aligned} E(\mathbf{w}_i = \delta\phi_i) - E(\mathbf{w}_i = 0) &= -\delta^2 c_i \lambda_i + (\mu/2) c_i (\delta^2 - 1)^2 < 0 \\ \text{for } (\lambda + \mu - \sqrt{\lambda(\lambda + 2\mu)})/\mu &< \delta^2 < (\lambda + \mu + \sqrt{\lambda(\lambda + 2\mu)})/\mu. \end{aligned}$$

The rest of the proof is similar to Theorem 5.4 above. ■

Theorem 5.7. For the ordinary differential equation (5.46), let A5.1 and A5.3 hold. Then $W=\Phi PDU$ are equilibrium points of (5.46), where $D=[D_1|0]^T \in \mathbb{R}^{n \times p}$ with $D_1 \in \mathbb{R}^{p \times p}$ is diagonal with elements $d_i = +1, -1$ or 0 , $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $U \in \mathbb{R}^{p \times p}$ is an arbitrary rotation matrix, i.e., $U^T U = U U^T = I_p$.

Proof. We need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$2AW - WW^T AW - AWW^T W - \mu W(W^T W - I_p) = 0. \quad (5.47)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. Let $W=QDV$ be the singular value decomposition of W , where $Q \in \mathbb{R}^{n \times n}$, and $V \in \mathbb{R}^{p \times p}$ are orthonormal, and $D \in \mathbb{R}^{n \times p}$ is diagonal. Replacing QDV for W in (5.47) and defining $B=Q^T A Q$, we get from (5.47):

$$2BD - DD^T BD - BDD^T D - \mu D(D^T D - I_p) = 0. \quad (5.48)$$

Let $D=[D_1|0]^T$ where D_1 is a $p \times p$ diagonal matrix with diagonal elements d_i for $i=1,\dots,p$. Let

$B = \begin{bmatrix} B_1 & B_2 \\ B_2^T & B_3 \end{bmatrix}$ be a partition of B where $B_1 \in \mathbb{R}^{p \times p}$. From (5.48) we get:

$$(2B_1 - D_1^2 B_1 - B_1 D_1^2 - \mu(D_1^2 - I_p)) D_1 = 0. \quad (5.49)$$

We conclude:

$$D_1 = 0 \text{ or } 2B_1 - D_1^2 B_1 - B_1 D_1^2 - \mu(D_1^2 - I_p) = 0. \quad (5.50)$$

From the second equation in (5.50), we obtain:

$$(2b_{ii} + \mu)(d_i^2 - 1) = 0 \text{ for } i=j \text{ and } b_{ij}(d_i^2 + d_j^2 - 2) = 0 \text{ for } i \neq j, \quad (5.51)$$

where $B_1 = [b_{ij}]$ for $i, j = 1, \dots, p$. Since $\mu > 0$, we satisfy both equations with $d_i = \pm 1$. Thus, $W = QDV$, where $d_i = \pm 1$ for $i = 1, \dots, p$. Clearly, $W^T W = I_p$. Then, from (5.47), we obtain:

$$AW = WW^T AW. \quad (5.52)$$

Let $W^T AW = U^T L U$ be the eigen-decomposition of $W^T AW$, where $U \in \mathbb{R}^{p \times p}$ is orthonormal, and $L \in \mathbb{R}^{p \times p}$ is diagonal. From (5.52), we get $A(WU^T) = (WU^T)L$, where WU^T is orthonormal. Then, the columns of WU^T consist of the p orthonormal eigenvectors of A , i.e., $WU^T = \Phi P D$, where $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $D = [D_1 | 0]^T \in \mathbb{R}^{n \times p}$ where $D_1 \in \mathbb{R}^{p \times p}$ is a diagonal matrix with elements $d_i = \pm 1$ for $i = 1, \dots, p$. Combining all results, we conclude that $W = \Phi P D U$, where U is orthonormal, and $d_i = 0$ or ± 1 . ■

Theorem 5.8. Let A5.1 and A5.3 hold. Then $W = \Phi D U$, where $D = [D_1 | 0]^T \in \mathbb{R}^{n \times p}$, $D_1 = \text{diag}(d_1, \dots, d_p) \in \mathbb{R}^{p \times p}$, $d_i = \pm 1$ for $i = 1, \dots, p$, and $U \in \mathbb{R}^{p \times p}$ is an arbitrary rotation matrix, are stable equilibrium points of the ODE (5.46) and strict global minimum points of the objective function (5.44). In addition, $W = \Phi P D U$, where $d_i = 0$ for $i \leq p$ or $P \neq I$, are unstable equilibrium points of the ODE (5.46).

Proof. From (5.44), the energy function $E(W)$ for the AL2 Homogeneous adaptive algorithm is:

$$E(W) = -2\text{tr}(W^T A W) + \text{tr}(W^T A W W^T W) + \frac{\mu}{2} \text{tr}((W^T W - I_p)^2). \quad (5.53)$$

From Theorem 5.7, $W = \Phi P D U = \Psi U$, where $\Psi \in \Phi P D$. Then from (5.53):

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where}$$

$$E_i(\mathbf{w}_i) = -2\mathbf{\Psi}_i^T A \mathbf{\Psi}_i + \mathbf{\Psi}_i^T A \mathbf{\Psi}_i \mathbf{\Psi}_i^T \mathbf{\Psi}_i + \sum_{j=1, j \neq i}^p \mathbf{\Psi}_j^T A \mathbf{\Psi}_i \mathbf{\Psi}_j^T \mathbf{\Psi}_i + \mu \left(\sum_{j=1, j \neq i}^p (\mathbf{\Psi}_j^T \mathbf{\Psi}_i)^2 + (\mathbf{\Psi}_i^T \mathbf{\Psi}_i - 1)^2 \right).$$

Here $\mathbf{\Psi}_i = d_{(i)} \phi_{(i)}$ is the i^{th} column of Ψ for $i = 1, \dots, p$, and $d_{(i)} = 0$ or ± 1 . We first prove that $d_{(i)} = 0$ is an unstable equilibrium point of $E(W)$. We perturb $\mathbf{\Psi}_i$ by $\delta \phi_i$. Then

$$E(\mathbf{\Psi}_i = \delta \phi_i) - E(\mathbf{\Psi}_i = 0) = -2\delta^2 \lambda_i + \delta^4 \lambda_i + \mu(\delta^2 - 1)^2 < 0$$

$$\text{for } (\lambda + \mu - \sqrt{\lambda(\lambda + \mu)})/(\lambda + \mu) < \delta^2 < (\lambda + \mu + \sqrt{\lambda(\lambda + \mu)})/(\lambda + \mu).$$

We next prove that $\mathbf{\Psi}_i = \pm \phi_r$, ($r > p$) is an unstable equilibrium point of $E(W)$. We perturb $\mathbf{\Psi}_i$ by $\delta \phi_p$, i.e., $\mathbf{\Psi}_i = \pm(\phi_r + \delta \phi_p)/\sqrt{1 + \delta^2}$. Then

$$E\left(\pm(\phi_r + \delta\phi_p)/\sqrt{(1+\delta^2)}\right) - E(\pm\phi_r) = -\delta^2(\lambda_p - \lambda_r)/(1+\delta^2) < 0.$$

We next prove that $\psi_i = \pm\phi_r$, ($r \leq p$) is a stable equilibrium point of $E(W)$. We perturb ψ_i by $\delta\phi_s$, i.e., $\psi_i = \pm(\phi_r + \delta\phi_s)/\sqrt{(1+\delta^2)}$ for $s > p$. Then

$$E\left(\pm(\phi_r + \delta\phi_s)/\sqrt{(1+\delta^2)}\right) - E(\pm\phi_r) = \delta^2(\lambda_r - \lambda_s)/(1+\delta^2) > 0.$$

Thus, the columns of Ψ consisting of the first p orthonormal eigenvectors of A are the stable minimum points of $E(W)$. ■

Theorem 5.9. *Let A5.1 and A5.3 hold. Then, all the equilibrium points of the ODE (5.56) are up to an arbitrary permutation of the eigenvectors of A weighted by 0, +1 or -1, i.e., any point $W = [d_{(1)}\phi_{(1)} d_{(2)}\phi_{(2)} \dots d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0, +1$ or -1 , is an equilibrium point of the ODE (5.56).*

Proof. We need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$2AW - WUT(W^T AW) - AWUT(W^T W) - \mu WUT(W^T W - I_p) = 0. \quad (5.57)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. We multiply (5.57) on the left by W^T , and define $G = W^T AW + (\mu/2)W^T W$, and $H = W^T W$. From (5.57), we obtain:

$$2G = H UT(G) + G UT(H). \quad (5.58)$$

Since A is positive definite by Assumption A5.1, and since W is assumed to be nonzero, both G and H have positive diagonal elements. From (5.58), we conclude that $H = I_p$, and G is diagonal. Thus, $W^T W = I_p$, and $W^T AW$ is diagonal. Then, the columns of W consist of the p orthonormal eigenvectors of A , i.e., $W = \Phi P D$, where $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $D = [D_1 | 0]^T \in \mathbb{R}^{n \times p}$ where $D_1 \in \mathbb{R}^{p \times p}$ is a diagonal matrix with elements $d_i = \pm 1$ for $i=1, \dots, p$. Combining all results, we conclude that $W = \Phi P D$, where $d_i = 0$ or ± 1 . ■

Theorem 5.10. *Let A5.1 and A5.3 hold. Then, the points $W^* = [\pm\phi_1 \pm\phi_2 \dots \pm\phi_p]$ are the strict global minimum points of the objective function (5.54) and stable equilibrium points of the ODE (5.56). In addition, the points $W = [d_{(1)}\phi_{(1)} d_{(2)}\phi_{(2)} \dots d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0$ or $\phi_{(i)} \neq \phi_i$ for $i \in \{1, 2, \dots, p\}$ are unstable equilibrium points of the ODE (5.56).*

Proof. From (5.54), the energy function $E(W)$ for the AL2 Deflation adaptive algorithm is:

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where}$$

$$E_i(\mathbf{w}_i) = -2\mathbf{w}_i^T A \mathbf{w}_i + \mathbf{w}_i^T A \mathbf{w}_i \mathbf{w}_i^T \mathbf{w}_i + 2 \sum_{j=1}^{i-1} \mathbf{w}_j^T A \mathbf{w}_i \mathbf{w}_j^T \mathbf{w}_i + \mu \left(\sum_{j=1}^{i-1} (\mathbf{w}_j^T \mathbf{w}_i)^2 + \frac{1}{2} (\mathbf{w}_i^T \mathbf{w}_i - 1)^2 \right).$$

From Theorem 5.9, $\mathbf{w}_i = d_{(i)} \phi_{(i)}$, $d_{(i)} = 0, \pm 1$, is the i^{th} column of W for $i=1, \dots, p$. The rest of the proof is same as Theorem 5.8 with ψ_j substituted by \mathbf{w}_j . We conclude that the columns of W consisting of the first p orthonormal eigenvectors of A as $W = [\pm \phi_1 \dots \pm \phi_p]$ are stable minimum points of $E(W)$. ■

Theorem 5.11. *Let A5.1 and A5.3 hold. Then, all the equilibrium points of the ODE (5.61) are up to an arbitrary permutation of the eigenvectors of A weighted by 0, +1 or -1, i.e., any point $W [d_{(1)} \phi_{(1)} d_{(2)} \phi_{(2)} \dots d_{(p)} \phi_{(p)}]$, where $d_{(j)} = 0, +1$ or -1 , is an equilibrium point of the ODE (5.61).*

Proof. In order to satisfy the first order conditions (see Section 2.10) for the existence of the equilibrium points of the joint objective functions $J(\mathbf{w}_i; A)$ (in (5.39)) for $i=1, \dots, p$, we need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$2AWC - WCW^T AW - AWCW^T W - \mu WC(W^T W - I_p) = 0. \quad (5.62)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. We multiply (5.62) on the left by W^T , and define $G = W^T A W + (\mu/2) W^T W$, and $H = W^T W$. From (5.62), we obtain:

$$2GC = HCG + GCH. \quad (5.63)$$

Since H , G , and $HCG + GCH$ are symmetric matrices, we conclude that $GC = CG$. Since C is diagonal with distinct diagonal elements, G is diagonal. Let $GC = CG = D \in \mathbb{R}^{p \times p}$ be a diagonal matrix. Then, from (5.63), we get:

$$2D = HD + DH. \quad (5.64)$$

Since W is assumed to be nonzero, the diagonal elements of $H = W^T W$ are positive. This implies that $H = I_p$. The rest of the proof is similar to Theorem 5.9 above. ■

Theorem 5.12. *Let A5.1 and A5.3 hold. Then, the points $W^* = [\pm \phi_1 \pm \phi_2 \dots \pm \phi_p]$ are stable equilibrium points of the ODE (5.61) and strict global minimum points of the objective function in (5.59). In*

addition, the points $W = [d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0$ or $\phi_{(i)} \neq \phi_i$ for $i \in \{1, 2, \dots, p\}$ are unstable equilibrium points of the ODE (5.61).

Proof. From (5.59), the energy function $E(W)$ for the PF Weighted adaptive algorithm is:

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where}$$

$$E_i(\mathbf{w}_i) = -2c_i \mathbf{w}_i^T A \mathbf{w}_i + c_i \mathbf{w}_i^T A \mathbf{w}_i \mathbf{w}_i^T \mathbf{w}_i + 2 \sum_{j=1, j \neq i}^p c_j \mathbf{w}_j^T A \mathbf{w}_i \mathbf{w}_j^T \mathbf{w}_i + \mu \left(\sum_{j=1, j \neq i}^p c_j (\mathbf{w}_j^T \mathbf{w}_i)^2 + \frac{c_i}{2} (\mathbf{w}_i^T \mathbf{w}_i - 1)^2 \right)$$

Here $\mathbf{w}_i = d_{(i)}\phi_{(i)}$, $d_{(i)} = 0$ or ± 1 , is the i^{th} column of W for $i=1, \dots, p$. The rest of the proof is similar to Theorem 5.10 above. ■

Theorem 7.1. For the ordinary differential equation (7.17), let A7.1 and A7.3 hold. Then $W = \Phi P D U$ are equilibrium points of (7.17), where $D = [D_1 \mid 0]^T \in \mathbb{R}^{n \times p}$ with $D_1 \in \mathbb{R}^{p \times p}$ is diagonal with elements $d_i = +1, -1$ or 0 , $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $U \in \mathbb{R}^{p \times p}$ is an arbitrary rotation matrix, i.e., $U^T U = U U^T = I_p$.

Proof. We need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$2AW - AWW^T BW - BWW^T AW = 0. \quad (7.18)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. Let $W = B^{-1/2} Q D V$ be the singular value decomposition of W , where $Q \in \mathbb{R}^{n \times n}$, and $V \in \mathbb{R}^{p \times p}$ are orthonormal, and $D \in \mathbb{R}^{n \times p}$ is diagonal. Replacing $Q D V$ for W in (7.28) and defining $M = Q^T B^{-1/2} A B^{-1/2} Q$, we get from (7.28):

$$2MD - MDD^T D - DD^T MD = 0. \quad (7.19)$$

Let $D = [D_1 \mid 0]^T$ where D_1 is a $p \times p$ diagonal matrix with diagonal elements d_i for $i=1, \dots, p$. Let

$M = \begin{bmatrix} M_1 & M_2 \\ M_2^T & M_3 \end{bmatrix}$ be a partition of M where $M_1 \in \mathbb{R}^{p \times p}$. From (7.19) we get:

$$(2M_1 - M_1 D_1^2 - D_1^2 M_1) D_1 = 0. \quad (7.20)$$

We conclude:

$$D_1 = 0 \text{ or } 2M_1 - M_1 D_1^2 - D_1^2 M_1 = 0. \quad (7.21)$$

From the second equation in (7.21), we obtain:

$$m_{ii}(d_i^2 - 1) = 0 \text{ for } i=j, \text{ and } m_{ij}(d_i^2 + d_j^2 - 2) = 0 \text{ for } i \neq j, \quad (7.22)$$

where $M_1 = [m_{ij}]$ for $i, j=1, \dots, p$. We satisfy both equations with $d_i = \pm 1$. Thus, $W = B^{-1/2} Q D V$, where $d_i = \pm 1$ for $i=1, \dots, p$. Clearly, $W^T B W = I_p$. Then, from (7.18), we obtain:

$$AW = BWW^T AW. \quad (7.23)$$

Let $W^T A W = U^T L U$ be the eigen-decomposition of $W^T A W$, where $U \in \mathbb{R}^{p \times p}$ is orthonormal, and $L \in \mathbb{R}^{p \times p}$ is diagonal. From (7.23), we get $A W U^T = B W U^T L$, where $W U^T$ is orthonormal with respect to B . Then, the columns of $W U^T$ consist of the p orthonormal eigenvectors of A with respect to B , i.e., $W U^T = \Phi P D$, where $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $D = [D_1 \mid 0]^T \in \mathbb{R}^{n \times p}$ where $D_1 \in \mathbb{R}^{p \times p}$ is a diagonal matrix with elements $d_i = \pm 1$ for $i=1, \dots, p$. Combining all results, we conclude that $W = \Phi P D U$, where U is orthonormal, and $d_i = 0$ or ± 1 . ■

Theorem 7.2. *Let A7.1 and A7.3 hold. Then $W = \Phi D U$, where $D = [D_1 \mid 0]^T \in \mathbb{R}^{n \times p}$, $D_1 = \text{diag}(d_1, \dots, d_p) \in \mathbb{R}^{p \times p}$, $d_i = \pm 1$ for $i=1, \dots, p$, and $U \in \mathbb{R}^{p \times p}$ is an arbitrary rotation matrix, are stable equilibrium points of the ODE (7.17) and strict global minimum points of the objective function (7.14). In addition, $W = \Phi P D U$, where $d_i = 0$ for $i \leq p$ or $P \neq I$, are unstable equilibrium points of the ODE (7.27).*

Proof. From (7.14), the energy function $E(W)$ for the XU Homogeneous adaptive algorithm is:

$$E(W) = -2\text{tr}(W^T A W) + \text{tr}(W^T A W W^T B W). \quad (7.24)$$

From Theorem 7.1, $W = \Phi P D U = \Psi U$, where $\Psi = \Phi P D$. Then from (7.14):

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where}$$

$$E_i(\mathbf{w}_i) = -2\Psi_i^T A \Psi_i + \Psi_i^T A \Psi_i \Psi_i^T B \Psi_i + \sum_{j=1, j \neq i}^p \Psi_j^T A \Psi_i \Psi_j^T B \Psi_i.$$

Here $\Psi_i = d_{(i)} \phi_{(i)}$ is the i^{th} column of Ψ for $i=1, \dots, p$, and $d_{(i)} = 0$ or ± 1 . We first prove that $d_{(i)} = 0$ is an unstable equilibrium point of $E(W)$. We perturb Ψ_i by $\delta \phi_i$. Then

$$E(\Psi_i = \delta_i) - E(\Psi_i = 0) = -2\delta^2 \lambda_i + \delta^4 \lambda_i < 0 \text{ for } 0 < \delta < \sqrt{2}.$$

We next prove that $\Psi_i = \pm \phi_r$ ($r > p$) is an unstable equilibrium point of $E(W)$. We perturb Ψ_i by $\delta \phi_p$, i.e., $\Psi_i = \pm(\phi_r + \delta \phi_p) / \sqrt{(1 + \delta^2)}$. Then

$$E\left(\pm(\phi_r + \delta \phi_p) / \sqrt{(1 + \delta^2)}\right) - E(\pm \phi_r) = -\delta^2(\lambda_p - \lambda_r) / (1 + \delta^2) < 0.$$

We next prove that $\Psi_i = \pm \phi_r$ ($r \leq p$) is a stable equilibrium point of $E(W)$. We perturb Ψ_i by $\delta \phi_s$, i.e., $\Psi_i = \pm(\phi_r + \delta \phi_s) / \sqrt{(1 + \delta^2)}$ for $s > p$. Then

$$E\left(\pm(\phi_r + \delta \phi_s) / \sqrt{(1 + \delta^2)}\right) - E(\pm \phi_r) = \delta^2(\lambda_r - \lambda_s) / (1 + \delta^2) > 0.$$

Thus, the columns of Ψ consisting of the first p orthonormal eigenvectors of A with respect to B are the stable minimum points of $E(W)$. ■

Theorem 7.3. Let A7.1 and A7.3 hold. Then, all the equilibrium points of the ODE (7.29) are up to an arbitrary permutation of the eigenvectors of A weighted by 0, +1 or -1, i.e., any point $W = [d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0, +1$ or -1 , is an equilibrium point of the ODE (7.29).

Proof. In order to satisfy the first order conditions (see Section 2.10) for the existence of the equilibrium points of the joint objective functions $J(w_i; A, B)$ (in (7.27)) for $i=1, \dots, p$, we need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$2AWC - BWCW^T AW - AWCW^T BW = 0. \quad (7.30)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. We multiply (7.40) on the left by W^T , and define $G=W^TAW$, and $H=W^TBW$. From (7.30), we obtain:

$$2GC = HCG + GCH. \quad (7.31)$$

Since H , G , and C are symmetric matrices, we conclude that GC is symmetric. Since C is diagonal with distinct diagonal elements, G is also diagonal. Let $GC=CG=D \in \mathbb{R}^{p \times p}$ be a diagonal matrix. Then, from (7.31), we get:

$$2D=HD+DH. \quad (7.32)$$

Since W is assumed to be nonzero, the diagonal elements of $H=W^TBW$ are positive. This implies that $H=I_p$. Thus, $W^TBW=I_p$, and W^TAW is diagonal. Then, the columns of W consist of the p orthonormal eigenvectors of A with respect to B , i.e., $W=\Phi PD$, where $P \in \mathbb{R}^{n \times n}$ is an arbitrary permutation matrix, and $D=[D_1 \ 0]^T \in \mathbb{R}^{n \times p}$ where $D_1 \in \mathbb{R}^{p \times p}$ is a diagonal matrix with elements $d_i = \pm 1$ for $i=1, \dots, p$. Combining all results, we conclude that $W = \Phi PD$, where $d_i = 0$ or ± 1 . ■

Theorem 7.4. Let A7.1 and A7.3 hold. Then, the points $W^* = [\pm\phi_1 \ \pm\phi_2 \ \dots \ \pm\phi_p]$ are stable equilibrium points of the ODE (7.29) and strict global minimum points of the objective function in (7.27). In addition, the points $W = [d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0$ or $\phi_{(i)} \neq \phi_i$ for $i \in \{1, 2, \dots, p\}$ are unstable equilibrium points of the ODE (7.29).

Proof. From (7.27), the energy function $E(W)$ for the XU Weighted adaptive algorithm is:

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where}$$

$$E_i(\mathbf{w}_i) = -2c_i \mathbf{w}_k^{iT} A_k \mathbf{w}_k^i + c_i \left(\mathbf{w}_k^{iT} A_k \mathbf{w}_k^i \right) \left(\mathbf{w}_k^{iT} B_k \mathbf{w}_k^i \right) + 2 \sum_{j=1, j \neq i}^p c_j \mathbf{w}_k^{iT} A_k \mathbf{w}_k^j \mathbf{w}_k^{jT} B_k \mathbf{w}_k^i.$$

Here $w_i = d_{(i)}\phi_{(i)}$, $d_{(i)} = 0$ or ± 1 , is the i^{th} column of W for $i=1, \dots, p$. The rest of the proof is similar to Theorem 7.2 above. ■

Theorem 7.5. Let A7.1 and A7.3 hold. Then, all the equilibrium points of the ODE (7.68) are up to an arbitrary permutation of the eigenvectors of A weighted by 0, +1 or -1, i.e., any point $W = [d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0, +1$ or -1 , is an equilibrium point of the ODE (7.68).

Proof. We need to find a $W \in \mathbb{R}^{n \times p}$ such that

$$(AW - BWUT(W^T AW))(W^T BW)^{-1} = 0. \quad (7.69)$$

The trivial solution is $W=0$. We next assume that $W \neq 0$. We multiply (7.69) on the left by W^T , and define $G=W^TAW$, and $H=W^TBW$. From (7.79), we obtain:

$$G = H UT(G). \quad (7.70)$$

Since G is symmetric, $HUT(G)$ is also symmetric. Since W is assumed to be nonzero, both G and H have positive diagonal elements. From (7.70), we conclude that $H=I_p$, and G is diagonal. Thus, $W^TBW=I_p$, and W^TAW is diagonal. The rest of the proof is similar to Theorem 7.3. We conclude that $W = \Phi PD$, where $d_i = 0$ or ± 1 . ■

Theorem 7.6. Let A7.1 and A7.3 hold. Then, the points $W^* = [\pm\phi_1 \ \pm\phi_2 \ \dots \ \pm\phi_p]$ are the strict global minimum points of the objective function (7.66) and stable equilibrium points of the ODE (7.68). In addition, the points $W = [d_{(1)}\phi_{(1)} \ d_{(2)}\phi_{(2)} \ \dots \ d_{(p)}\phi_{(p)}]$, where $d_{(i)} = 0$ or $\phi_{(i)} \neq \phi_i$ for $i \in \{1, 2, \dots, p\}$ are unstable equilibrium points of the ODE (7.68).

Proof. From (7.66), the energy function $E(W)$ for the RQ Deflation adaptive algorithm is:

$$E(W) = \sum_{i=1}^p E_i(\mathbf{w}_i), \text{ where}$$

$$E_i(\mathbf{w}_i) = -\frac{\mathbf{w}_k^{iT} A_k \mathbf{w}_k^i}{\mathbf{w}_k^{iT} B_k \mathbf{w}_k^i} + \alpha \left(\mathbf{w}_k^{iT} B_k \mathbf{w}_k^i - 1 \right) + 2 \sum_{j=1}^i \beta_j \mathbf{w}_k^{jT} B_k \mathbf{w}_k^i.$$

From Theorem 7.5, $\mathbf{w}_j = d_{(i)}\phi_{(i)}$, $d_{(i)} = 0, \pm 1$, is the i^{th} column of W for $i=1, \dots, p$. Rest of the proof is same as Theorem 7.4 with ψ_j substituted by \mathbf{w}_j . We conclude that the columns of W consisting of the first p orthonormal eigenvectors of A since $W = [\pm\phi_1 \ \dots \ \pm\phi_p]$ are stable minimum points of $E(W)$. ■