

# A Nonconvex Unconstrained Method for Eigenvalue Problems and A Nonsingular System for Eigenvector Estimation

Yunho Kim\*

November 28, 2016

## Abstract

We propose a nonconvex unconstrained minimization problem for eigenvalue problems. In this framework, given a symmetric matrix  $A$ , it turns out that any nonzero critical point is an eigenvector of  $A$  and any local minimizer is a global minimizer, an eigenvector of  $A$  corresponding to the smallest eigenvalue. Unlike the conventional way of estimating an eigenvector with a given eigenvalue, our proposed problem estimates the corresponding eigenvalue from the estimated eigenvector of  $A$ . We analyze two algorithms to solve the proposed problem, one of which guarantees convergence to a global minimizer. This makes our method applicable even to a singular matrix  $A$ , where conventional methods avoid finding an eigenvector corresponding to the eigenvalue 0. The other tries to find a critical point that may not be a global minimizer, however, is faster and gives rise to a nonsingular linear system, as a byproduct, whose unique solution is an eigenvector of  $A$ , which we believe is a new approach to eigenvector estimation. Our proposed method applies to nonsquare matrices, complex hermitian matrices and also to infinite dimensional cases such as self-adjoint elliptic operators, some of which will be presented as specific applications.

## 1 Introduction

Given an  $N \times N$  matrix  $A$ , the eigenvalue problem of our interest is to find an eigenvalue and its corresponding eigenvector of  $A$ , that is, to solve  $Ax = \lambda x$  for  $x$  and  $\lambda$ . This is one of the most fundamental problems in mathematics with applications to all other fields of science. Especially, one may be interested in estimating the largest and the smallest eigenvalues. If  $A$  is symmetric and positive definite, then finding the smallest eigenvalue of  $A$  is the same as finding the largest eigenvalue of  $A^{-1}$ . However, finding the smallest eigenvalue of  $A$  involves solving systems of linear equations of type  $Ax = b$ . For example, the inverse power method to estimate the smallest eigenvalue involves solving  $Ax = b$  as well as matrix multiplication, whereas the power method to estimate the largest eigenvalue involves only matrix multiplication. From these observations, we asked ourselves a question, “Is it possible to compute the smallest eigenvalue of  $A$  through only basic matrix operations such as multiplication and addition, without solving  $Ax = b$ ?”

Our interest of estimating the smallest eigenvalue and its corresponding eigenvector extends to the following infinite dimensional application, as well. On a compact manifold  $\mathcal{M}$ ,

---

\*Department of Mathematical Sciences, UNIST, South Korea. Email: yunhokim@unist.ac.kr

eigenvalues of the Laplacian  $-\Delta$  reveals important structures of  $\mathcal{M}$ , which makes understanding the eigenvalues of  $-\Delta$  on  $\mathcal{M}$  very important. This has interesting applications. For example, in image processing there are a few interesting works (e.g. [7], [3], [9]) to distinguish reconstructed objects from point cloud data by evaluating  $-\Delta$  on the surfaces of the objects. We can even consider general self-adjoint linear elliptic operators and find their eigenvalues and eigenfunctions. With these theoretical and numerical points of view in mind, our main discussion will be concentrated on finding the smallest eigenvalue and a corresponding eigenvector of a nonzero symmetric matrix, which leads us to begin with the following well-known constrained problem: given a symmetric and positive definite matrix  $A$ ,

$$\min_{x \in \mathbb{R}^N} \langle x, Ax \rangle \quad \text{subject to} \quad \|x\|^2 = 1. \quad (1)$$

There have been a large number of works to solve (1) by the name of inverse iteration methods. In particular, we would like to mention the work [4] by J.E. Dennis and R.A. Tapia, which surveys historical developments of inverse and shifted inverse iterations and of Rayleigh quotient iteration, and which approaches the listed methods from the viewpoint of the Newton's method. The unconstrained version of (1) analyzed in [4] is

$$\min_{x \in \mathbb{R}^N} \langle x, Ax \rangle + \frac{\gamma}{2} (\|x\|^2 - 1)^2. \quad (2)$$

The authors of [4] explained why the inverse and the shifted inverse Rayleigh quotient iterations are fast and effective by showing the equivalence between the inverse Rayleigh quotient iteration and the Newton's method for (2), and also between the shifted inverse Rayleigh quotient iteration and the Newton's method for the shifted version of (2), when the given matrix  $A$  is symmetric and invertible. We refer to [4], and references therein, any interested reader in the developments of the inverse and shifted inverse Rayleigh quotient iteration methods.

There are, however, a few disadvantageous features of the functional in (2) that we paid attention to. First of all, [4] considered only nonsingular matrices for (2) just as all other conventional methods do. Second of all, in the simplest case when  $A$  is symmetric and positive definite, if  $0 < \gamma < \lambda_1$ , then the zero vector is the only critical point of the functional in (2) and, even if  $\gamma \geq \lambda_1$ , the critical points of the functional in (2) are only the eigenvalues of  $A$  less than  $\gamma$ . Hence, our goal is of two folds: 1) extending existing theories to singular matrices, and 2) removing the additional limitations imposed by the parameter  $\gamma$  in (2). Noting that the functional in (2) contains the term

$$(\|x\|^2 - 1)^2 = (\|x\| + 1)^2 (\|x\| - 1)^2,$$

we believed that the factor  $(\|x\| + 1)^2$ , even though this is convex, is not desirable because the factor tries to reduce the norm  $\|x\|$  to 0 during minimization. Therefore, our analysis begins without this factor.

The rest of this manuscript is organized as follows. In Section 2, we define and analyze our proposed minimization problem in the case of a nonzero real symmetric matrix  $A$ . In the analysis, we prove that our proposed method does not possess the disadvantageous features mentioned above. The same analysis applies to complex hermitian matrices, as well. Moreover, we analyze two minimization algorithms, the Gradient Descent method and the Newton's method, to find an eigenvalue and its corresponding eigenvector, where the former ensures convergence to a global minimizer and the latter ensures fast convergence to critical points. In fact, as for the Newton's method, we show that the a generated sequence  $\{x_k\}$

converges if and only if the sequence of their norms  $\{\|x_n\|\}$  converges and that the limit of  $\{\|x_k\|\}$  determines the corresponding eigenvalue. What is more interesting is that we can provide a nonsingular system of linear equations whose unique solution is an eigenvector of  $A$  when an exact eigenvalue of  $A$  is known. The same is true for general diagonalizable matrices. This is unusual because an exact eigenvalue gives rise to a singular system of linear equations, which is why conventional methods assume nonsingularity of  $A$  and generate a sequence of approximate eigenvalues and justify that the sequence converges to an exact eigenvalue. In Section 3, we provide a few applications such as generalized eigenvalue problems and an infinite dimensional example for self-adjoint linear elliptic operators, which shows applicability of our method to various types of eigenvalue and eigenvector (or eigenfunction) estimation, not only in the finite dimensional setting, but also in the infinite dimensional setting.

## 2 Eigenvalue problem on a finite dimensional space

A general setting is the following. Let  $\mathbb{F}$  be either  $\mathbb{R}$  or  $\mathbb{C}$ . Given a matrix  $A \in \mathbf{M}_{M \times N}(\mathbb{F})$ , we define  $F_A : \mathbb{F}^N \rightarrow \mathbb{F}$  by

$$F_A(x) = \frac{1}{2} \langle x, Bx \rangle + \frac{\gamma}{2} \|x\|^2 - \gamma \|x\|, \quad (3)$$

where  $B$  is

$$B = \begin{cases} A, & \text{if } A \text{ is real symmetric or complex hermitian,} \\ A^*A, & \text{if } A \text{ is neither real symmetric nor complex hermitian,} \end{cases}$$

and  $\langle x, y \rangle = y^*x$  and  $\|x\| = \sqrt{\langle x, x \rangle}$ . As usual,  $x \in \mathbb{F}^N$ ,  $Ax \in \mathbb{F}^M$  are considered as  $N \times 1$  and  $M \times 1$  column vectors, respectively, and  $|A| = \sqrt{A^*A}$ . We use the convention that  $x^* = x^T$  and  $|A| = \sqrt{A^T A}$  if  $\mathbb{F} = \mathbb{R}$ .

Our discussion on the real symmetric case extends to the complex hermitian case without alteration. Moreover, the same discussion applies to the singular value decomposition of a nonsymmetric matrix  $A$  with  $B = A^*A$  in the definition (3). Hence, and we will only present the real symmetric case, where the minimization problem that we will consider is

$$\min_{x \in \mathbb{R}^N} F_A(x),$$

which is equivalent to

$$\min_{x \in \mathbb{R}^N} \frac{1}{2} \langle x, Ax \rangle + \frac{\gamma}{2} (\|x\| - 1)^2.$$

### 2.1 The Real Symmetric Case

To begin with, we will consider the case that  $A \in \mathbf{M}_N(\mathbb{R})$  is nonzero and symmetric, where the functional that we minimize in (3) has the form

$$F_A(x) = \frac{1}{2} \langle x, Ax \rangle + \frac{\gamma}{2} \|x\|^2 - \gamma \|x\| \quad (4)$$

**Lemma 1.** *Let  $\lambda_1 \in \mathbb{R}$  be the smallest eigenvalue of  $A$ . For  $\gamma > \max(0, -\lambda_1)$ , the set of nonzero critical points of  $F_A$  is*

$$\left\{ x \in \mathbb{R}^n : x \text{ is an eigenvector of } A \text{ corresponding to an eigenvalue } \lambda \text{ with norm } \frac{\gamma}{\gamma + \lambda} \right\},$$

and

$$\min_{x \in \mathbb{R}^n} F_A(x) = -\frac{\gamma^2}{2(\gamma + \lambda_1)}.$$

*Proof.* Let  $x_0 \neq 0$  be a critical point of  $F_A$ . For any unit vector  $\theta \in \mathbb{R}^N$ , we note that

$$\lim_{\epsilon \rightarrow 0} \frac{F_A(x_0 + \epsilon\theta) - F_A(x_0)}{\epsilon} = \langle Ax_0, \theta \rangle + \gamma \langle x_0, \theta \rangle - \frac{\gamma}{\|x_0\|} \langle x_0, \theta \rangle = 0.$$

This implies

$$\left[ A + \gamma \left( 1 - \frac{1}{\|x_0\|} \right) I \right] x_0 = 0 \quad \Leftrightarrow \quad Ax_0 = \gamma \left( \frac{1}{\|x_0\|} - 1 \right) x_0.$$

Hence,  $x_0 \neq 0$  is a critical point of  $F_A$  if and only if  $x_0$  is an eigenvector of  $A$  corresponding to the eigenvalue

$$\lambda_0 = \gamma \left( \frac{1}{\|x_0\|} - 1 \right).$$

Moreover,

$$\|x_0\| = \frac{\gamma}{\gamma + \lambda_0} \quad \text{and} \quad F_A(x_0) = -\frac{\gamma}{2} \|x_0\| = -\frac{\gamma^2}{2(\gamma + \lambda_0)}$$

and the set of nonzero critical points of  $F_A$  is

$$\left\{ x \in \mathbb{R}^N : x \text{ is an eigenvector of } A \text{ corresponding to an eigenvalue } \lambda \text{ with norm } \frac{\gamma}{\gamma + \lambda} \right\}.$$

It is easy to see that a global minimizer  $x_*$  of  $F_A$  exists and is an eigenvector of  $A$  corresponding to an eigenvalue  $\lambda_*$  with norm  $\frac{\gamma}{\gamma + \lambda_*}$ . Therefore,  $F_A(x_*) = \min_{x \in \mathbb{R}^N} F_A(x)$  implies that  $\lambda_* = \lambda_1$  and

$$\min_{x \in \mathbb{R}^N} F_A(x) = -\frac{\gamma^2}{2(\gamma + \lambda_1)} < 0.$$

□

Throughout the section, we will assume  $\gamma > \max(0, -\lambda_1)$ , where  $\lambda_1$  is the smallest eigenvalue of  $A$  if no condition on  $\gamma$  is stated.

**Theorem 1.** *Any local minimizer  $x_*$  of  $F_A$  is a global minimizer.*

*Proof.* First of all, 0 is not a local minimizer of  $F_A$  because for any nonzero  $x_0 \in \mathbb{R}$ , we have

$$\lim_{t \rightarrow 0^+} \frac{F_A(tx_0) - F_A(0)}{t\|x_0\|} = -\gamma < 0.$$

Suppose that  $x_* \neq 0$  is a local minimizer of  $F_A$ . Lemma 1 says that  $x_*$  is an eigenvector of  $A$  corresponding to an eigenvalue  $\lambda_*$  with norm  $\frac{\gamma}{\gamma + \lambda_*}$  and that

$$F_A(x_*) = -\frac{\gamma^2}{2(\gamma + \lambda_*)} \geq -\frac{\gamma^2}{2(\gamma + \lambda_1)} = \min_{x \in \mathbb{R}} F_A(x),$$

where  $\lambda_1$  is the smallest eigenvalue of  $A$ . We may diagonalize  $A$  such that  $A = Q\Lambda Q^T$ , where  $\Lambda$  is a diagonal matrix with nondecreasing diagonal entries  $\lambda_1 \leq \dots \leq \lambda_N$  and  $Q$  is an

orthogonal matrix having  $\frac{x_*}{\|x_*\|}$  as the  $j^{\text{th}}$  column for some  $j$  implying  $\lambda_* = \lambda_j$ . Then, with  $y = Q^T x$ , (3) becomes

$$F_A(x) = F_A(Qy) = F_\Lambda(y). \quad (5)$$

For  $k = 1, 2, \dots, N$ , we set  $\mathbf{e}_k$  to be the  $k^{\text{th}}$  column of the identity matrix  $I \in \mathbf{M}_N(\mathbb{R})$ . Then,  $x_*$  being a local minimizer of  $F_A$  is equivalent to  $\frac{\gamma}{\gamma + \lambda_j} \mathbf{e}_j$  being a local minimizer of  $F_\Lambda$ . Note that

$$F_\Lambda(y) = \frac{1}{2} \langle \Lambda y, y \rangle + \frac{\gamma}{2} \|y\|^2 - \gamma \|y\| = \left( \frac{1}{2} \sum_{k=1}^N (\lambda_k + \gamma) y_k^2 \right) - \gamma \|y\| \quad (6)$$

and that  $\frac{\gamma}{\gamma + \lambda_1} \mathbf{e}_1$  is a global minimizer of  $F_\Lambda$ . If we consider  $F_\Lambda$  on the subspace spanned by  $\{\mathbf{e}_1, \mathbf{e}_j\}$  by defining  $H : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$H(y_1, y_2) = F_\Lambda(y_1 \mathbf{e}_1 + y_2 \mathbf{e}_j) = \frac{1}{2} (\lambda_1 + \gamma) y_1^2 + \frac{1}{2} (\lambda_j + \gamma) y_2^2 - \gamma \sqrt{y_1^2 + y_2^2},$$

then  $\nabla^2 H$  at  $\left(0, \frac{\gamma}{\gamma + \lambda_j}\right)$  must be positive semidefinite, that is,

$$\det \begin{bmatrix} (\gamma + \lambda_1) - \frac{\gamma y_2^2}{\sqrt{y_1^2 + y_2^2}^3} & \frac{\gamma y_1 y_2}{\sqrt{y_1^2 + y_2^2}^3} \\ \frac{\gamma y_1 y_2}{\sqrt{y_1^2 + y_2^2}^3} & (\gamma + \lambda_j) - \frac{\gamma y_1^2}{\sqrt{y_1^2 + y_2^2}^3} \end{bmatrix} = (\lambda_1 - \lambda_j)(\gamma + \lambda_j) \geq 0$$

at  $\left(0, \frac{\gamma}{\gamma + \lambda_j}\right)$ . This implies  $\lambda_1 \geq \lambda_j = \lambda_*$ , i.e.,  $\lambda_* = \lambda_1$  and

$$F_A(x_*) = -\frac{\gamma^2}{2(\gamma + \lambda_*)} = \min_{x \in \mathbb{R}^N} F_A(x).$$

Therefore, any local minimizer  $x_*$  of  $F_A$  is a global minimizer of  $F_A$ .  $\square$

In addition, we may be able to find all the eigenvalues and their corresponding eigenvectors of  $A$ .

**Theorem 2.** *Let  $A$  be symmetric with eigenvalues  $\lambda_1 \leq \dots \leq \lambda_N$ . Let  $\mathbf{x}_1, \dots, \mathbf{x}_k$  be the first  $k$  eigenvectors of  $A$  corresponding to the eigenvalues  $\lambda_1 \leq \dots \leq \lambda_k$ . We consider the following problem*

$$\min_{x \in \mathbb{R}^N} F_A(x) \quad \text{subject to} \quad \langle x, \mathbf{x}_i \rangle = 0, \quad i = 1, 2, \dots, k. \quad (7)$$

*Then, any local minimizer of (7) is a global minimizer corresponding to the eigenvalue  $\lambda_{k+1}$  with norm  $\frac{\gamma}{\gamma + \lambda_{k+1}}$ .*

*Proof.* With a diagonalization  $QDQ^T$  of  $A$ , where  $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$  and the first  $k$  columns of  $Q$  are  $\mathbf{x}_1, \dots, \mathbf{x}_k$ , and  $y = Q^T x$ , we have that  $F_A(x) = F_D(y)$  and  $\|x\| = \|y\|$ , so (7) is equivalent to

$$\min\{F_D(y) : y = (0, \dots, 0, y_{k+1}, \dots, y_n) \in \mathbb{R}^N\}. \quad (8)$$

Let  $D_k$  be the last  $(n - k) \times (n - k)$  block of  $D$ , i.e.,  $D_k = \text{diag}(\lambda_{k+1}, \dots, \lambda_N)$ . Then, (8) is equivalent to

$$\min_{z \in \mathbb{R}^{n-k}} F_{D_k}(z).$$

Theorem 1 applies to  $F_{D_k}$  and we are done.  $\square$

### 2.1.1 Algorithm 1 : The Gradient Descent Method

We have seen in the previous section that our model does not possess a local minimizer that is not a global minimizer. To find a global minimizer of our model, we will consider and analyze the following gradient descent method with stepsize  $\alpha_k > 0$ :

$$x_{k+1} = x_k - \alpha_k \nabla F_A(x_k) = x_k - \alpha_k \left( Ax_k + \gamma x_k - \frac{\gamma}{\|x_k\|} x_k \right), \quad \text{with } x_0 \neq 0. \quad (9)$$

Let  $\{q_1, \dots, q_N\}$  be an orthonormal basis for  $\mathbb{R}^N$  consisting of unit eigenvectors of  $A$  corresponding to the eigenvalues  $\lambda_1 \leq \dots \leq \lambda_N$ , respectively. If  $x_k = \mu_{k,1}q_1 + \dots + \mu_{k,N}q_N$ , then

$$x_{k+1} = \mu_{k+1,1}q_1 + \dots + \mu_{k+1,N}q_N = \sum_{i=1}^N \mu_{k,i} \left[ 1 - \alpha_k \left( \lambda_i + \gamma - \frac{\gamma}{\|x_k\|} \right) \right] q_i$$

implies

$$\mu_{k+1,i} = \mu_{k,i} \left[ 1 - \alpha_k \left( \lambda_i + \gamma - \frac{\gamma}{\|x_k\|} \right) \right], \quad i = 1, \dots, N,$$

resulting in

$$x_{k+1} = \sum_{i=1}^N \mu_{0,i} \prod_{j=0}^k \left[ 1 - \alpha_j \left( \lambda_i + \gamma - \frac{\gamma}{\|x_j\|} \right) \right] q_i. \quad (10)$$

For simplicity, we assume a fixed stepsize  $0 < \alpha_k = \alpha < \frac{1}{\lambda_N + \gamma}$  for all  $k \in \mathbb{N}$ .

One very distinctive feature of our model is that (9) can find a global minimizer, i.e., an eigenvector of  $A$  corresponding to the smallest eigenvalue, even when  $A$  is singular. Theorem 10 below is given in a general form.

**Theorem 3.** *A sequence  $\{x_k\}$  generated by (9) with  $x_0 = \mu_{0,1}q_1 + \dots + \mu_{0,N}q_N \neq 0$  converges to a critical point  $x_*$  of  $F_A$  which is an eigenvector of  $A$  corresponding the eigenvalue  $\lambda_l$  with*

$$\|x_*\| = \frac{\gamma}{\gamma + \lambda_l},$$

where  $l = \min\{j \in \{1, \dots, N\} | \mu_{0,j} \neq 0\}$ . More precisely,  $x_*$  is

$$\left( \frac{\gamma}{\gamma + \lambda_l} \right) \left( \frac{1}{\sqrt{\sum_{m:\lambda_m=\lambda_l} \mu_{0,m}^2}} \right) \sum_{m:\lambda_m=\lambda_l} \mu_{0,m} q_m.$$

*Proof.* For any  $x \neq 0$ , we get

$$\left\langle \nabla F_A(x), \frac{x}{\|x\|} \right\rangle = \left( \frac{\langle Ax, x \rangle}{\|x\|^2} + \gamma \right) \|x\| - \gamma = (\lambda + \gamma) \|x\| - \gamma, \quad (11)$$

where  $\lambda_1 \leq \lambda \leq \lambda_N$ . Then, since  $\alpha < \frac{1}{\lambda_N + \gamma}$ , we obtain

$$\|x_{k+1}\|^2 \geq \left| \left\langle x_k - \alpha \nabla F_A(x_k), \frac{x_k}{\|x_k\|} \right\rangle \right|^2 \geq (\|x_k\| - \alpha(\lambda + \gamma)\|x_k\| + \alpha\gamma)^2 > (\alpha\gamma)^2,$$

which implies  $\|x_k\| > \alpha\gamma$  for all  $k \in \mathbb{N}$ . Note that for  $x \neq 0$ ,

$$\nabla^2 F_A(x) = A + \gamma I - \frac{\gamma}{\|x\|} \left[ I - \left( \frac{x}{\|x\|} \right) \left( \frac{x}{\|x\|} \right)^T \right].$$

Then, it is easy to see that  $\|\nabla^2 F_A(x)\| \leq \max(\lambda_N + \gamma, \frac{1}{\alpha}) \leq \frac{1}{\alpha}$  for  $\|x\| \geq \alpha\gamma$ . Moreover, the line segment connecting  $x_k$  and  $x_{k+1}$  for  $k \in \mathbb{N}$  lies entirely in  $\{x \in \mathbb{R}^N : \|x\| \geq \alpha\gamma\}$ . To see this, we take  $0 < t < 1$  and observe that for  $k \in \mathbb{N}$ ,

$$\|x_k - t\alpha\nabla F(x_k)\|^2 \geq (\|x_k\| - t\alpha(\lambda + \gamma)\|x_k\| + t\alpha\gamma)^2 > (\alpha\gamma)^2(1 - t\alpha(\lambda + \gamma) + t)^2 > (\alpha\gamma)^2.$$

Then, for each  $k \geq 1$ , we have that

$$\begin{aligned} F_A(x_{k+1}) &\leq F_A(x_k) + \langle \nabla F_A(x_k), x_{k+1} - x_k \rangle + \frac{1}{2\alpha} \|x_{k+1} - x_k\|^2 \\ &= F_A(x_k) - \frac{\alpha}{2} \|\nabla F_A(x_k)\|^2, \end{aligned}$$

which implies

$$\|\nabla F_A(x_k)\|^2 \leq \frac{2}{\alpha} (F_A(x_k) - F_A(x_{k+1})),$$

and that for any  $K \geq 1$ ,

$$\frac{1}{\alpha^2} \sum_{k=1}^K \|x_{k+1} - x_k\|^2 = \sum_{k=1}^K \|\nabla F_A(x_k)\|^2 \leq \frac{2}{\alpha} (F_A(x_1) - (\min_{x \in \mathbb{R}^N} F_A(x))). \quad (12)$$

Since  $F_A$  is coercive and  $F_A(x_k) \leq F_A(x_1) < \infty$  for all  $k \geq 1$ ,  $\{x_k\}$  must be a bounded sequence in

$$\{x \in \mathbb{R}^N : \|x\| \geq \alpha\gamma\}.$$

Choosing a subsequence,  $\{x_{k_n}\}$ , converging to  $x_*$ , (12) implies that  $\nabla F_A(x_*) = 0$ , i.e., according to Lemma 1,  $x_*$  is an eigenvector of  $A$  corresponding to an eigenvalue  $\lambda_i$  for some  $1 \leq i \leq N$  with norm

$$\|x_*\| = \frac{\gamma}{\gamma + \lambda_i}.$$

Knowing that  $\{F_A(x_k)\}$  is a decreasing and bounded sequence, we can easily derive that any subsequential limit  $\tilde{x}$  of  $\{x_k\}$  satisfies  $F_A(\tilde{x}) = F_A(x_*)$  and  $A\tilde{x} = \lambda_i\tilde{x}$  with norm

$$\|\tilde{x}\| = \frac{\gamma}{\gamma + \lambda_i}.$$

Hence, we can conclude that

$$\lim_{k \rightarrow \infty} \|x_k\| = \frac{\gamma}{\gamma + \lambda_i}.$$

Now, we set  $l = \min\{j \in \{1, \dots, N\} | \mu_{0,j} \neq 0\}$ . Then,  $\lambda_i \geq \lambda_l$ . Suppose  $\lambda_i > \lambda_l$ . We note that for all  $k \in \mathbb{N}$ ,

$$1 - \alpha \left( \lambda_l + \gamma - \frac{\gamma}{\|x_k\|} \right) > 1 - \frac{\lambda_l + \gamma}{\lambda_N + \gamma} + \frac{\alpha\gamma}{\|x_k\|} \geq \frac{\alpha\gamma}{\|x_k\|} > 0$$

and that as  $k \rightarrow \infty$ ,

$$1 - \alpha \left( \lambda_l + \gamma - \frac{\gamma}{\|x_k\|} \right) \rightarrow 1 - \alpha(\lambda_l - \lambda_i) > 1.$$

This implies from (10) that

$$\mu_{0,l} \prod_{j=0}^k \left[ 1 - \alpha \left( \lambda_l + \gamma - \frac{\gamma}{\|x_j\|} \right) \right] \rightarrow \infty \text{ as } k \rightarrow \infty.$$

This is a contradiction because  $\{x_k\}$  is a bounded sequence. Therefore,  $\lambda_i = \lambda_l$ .

Moreover, for  $\lambda_p > \lambda_l$ , we have

$$1 - \alpha\left(\lambda_p + \gamma - \frac{\gamma}{\|x_k\|}\right) \rightarrow 1 - \alpha(\lambda_p - \lambda_l) \in (0, 1) \text{ as } k \rightarrow \infty,$$

implying

$$\mu_{0,p} \Pi_{j=0}^k \left[1 - \alpha\left(\lambda_p + \gamma - \frac{\gamma}{\|x_j\|}\right)\right] \rightarrow 0 \text{ as } k \rightarrow \infty.$$

Hence, referring to (10), we can see that

$$x_{k+1} - \left(\sum_{m:\lambda_m=\lambda_l} \mu_{0,m} q_m\right) \Pi_{j=0}^k \left[1 - \alpha\left(\lambda_l + \gamma - \frac{\gamma}{\|x_j\|}\right)\right] \rightarrow 0 \text{ as } k \rightarrow \infty.$$

Moreover, convergence of the norm  $\|x_k\|$  implies that

$$\Pi_{j=0}^{\infty} \left[1 - \alpha\left(\lambda_l + \gamma - \frac{\gamma}{\|x_j\|}\right)\right] = \frac{\gamma}{\gamma + \lambda_l} \left(\frac{1}{\sqrt{\sum_{m:\lambda_m=\lambda_l} \mu_{0,m}^2}}\right).$$

and  $x_k$  converges to

$$\left(\sum_{m:\lambda_m=\lambda_l} \mu_{0,m} q_m\right) \Pi_{j=0}^{\infty} \left[1 - \alpha\left(\lambda_l + \gamma - \frac{\gamma}{\|x_j\|}\right)\right].$$

□

If  $\lambda_1$  has multiplicity  $p \geq 1$ , i.e.,  $\lambda_1 = \dots = \lambda_p < \lambda_{p+1}$ , and if  $|\mu_{0,1}| + \dots + |\mu_{0,p}| > 0$ , then the theorem says that the sequence generated by the gradient descent method converges to an eigenvector of  $A$  corresponding to the smallest eigenvalue  $\lambda_1$  with norm  $\frac{\gamma}{\gamma + \lambda_1}$ . In fact, if we begin the algorithm with a randomly chosen initial guess  $x_0$ , then it is easy to see that it converges to a global minimizer with probability 1. Therefore, regardless of the definiteness of  $A$  and of the multiplicity of the smallest eigenvalue of  $A$ , the gradient descent method (9) finds a global minimizer.

When compared with the inverse power method, the two have the same assumption for convergence. However, our approach only requires matrix multiplications, whereas the inverse power method needs to solve a linear system at every iteration.

### 2.1.2 Strongly convex $F_A$ around a global minimum

It is easy to see that  $F_A$  is strongly convex near a global minimum point  $x_*$  since  $\nabla^2 F_A(x_*)$  is positive definite. In this section, we would like to see how much  $F_A$  could be strongly convex, and how large such a strongly convex region near a global minimum point could be, which will give us some idea about when the Newton's method, if ever applied, will converge.

**Theorem 4.** *Suppose that  $A$  is symmetric with the smallest eigenvalue  $\lambda_1$  of multiplicity 1 and that  $\gamma > \max(0, -\lambda_1)$ . Let  $F_A$  attain the global minimum value at  $x_*$ . Let  $\eta = \min(2(\gamma + \lambda_1), \lambda_2 - \lambda_1)$ . Then,  $F_A$  is  $\frac{\eta}{4}$ -strongly convex in the region*

$$\left\{x : \|x - x_*\| < \min\left(\frac{\gamma(\lambda_2 - \lambda_1)}{4(\gamma + \lambda_2)(\gamma + \lambda_1)}, \frac{\gamma}{5\gamma + 4\lambda_2 + \lambda_1}\right)\right\},$$

where  $\lambda_2 > \lambda_1$  is the second smallest eigenvalue of  $A$ .



*Proof.* For simplicity, we may assume  $\lambda_1 \geq 0$  since the proof below will not change when  $\lambda_1 < 0$  with  $\gamma + \lambda_1 > 0$ . Let  $k \geq 2$ . For some  $0 < \epsilon < \frac{\|x_*\|}{k}$ , we suppose that  $x$  satisfies  $\|x - x_*\| < \epsilon$ . Then, we have

$$\frac{\gamma(k-1)}{(\gamma + \lambda_1)k} = \|x_*\| \left(1 - \frac{1}{k}\right) \leq \|x\| \leq \|x_*\| \left(1 + \frac{1}{k}\right) = \frac{\gamma(k+1)}{(\gamma + \lambda_1)k}$$

and

$$\left\langle \frac{x}{\|x\|}, \frac{x_*}{\|x_*\|} \right\rangle > \frac{1}{2} \left( \frac{\|x\|}{\|x_*\|} + \frac{\|x_*\|}{\|x\|} \right) - \frac{\epsilon^2}{2\|x\|\|x_*\|} \geq 1 - \frac{1}{2k(k-1)} \geq 1 - \frac{1}{k^2}.$$

Noting that

$$\nabla^2 F(x) = A + \gamma \left(1 - \frac{1}{\|x\|}\right) I + \frac{\gamma}{\|x\|^3} x x^T,$$

we will estimate  $\langle z, \nabla^2 F_A(x) z \rangle$  for  $z \in \mathbb{R}^N$  with  $\|z\| = 1$ . Let  $\{(x_i, \lambda_i) : i = 1, \dots, N\}$  be a set of eigenpairs of  $A$  such that  $\{x_1, \dots, x_n\}$  with  $x_1 = \frac{x_*}{\|x_*\|}$  is an orthonormal basis for  $\mathbb{R}^n$  and  $0 \leq \lambda_1 < \lambda_2 \leq \dots \leq \lambda_N$ . Let  $V = \{c x_* : c \in \mathbb{R}\}$  and  $W = \text{span}\{x_i : i = 2, \dots, N\}$ . Then, for any  $z \in \mathbb{R}^N$ , there exists a unique decomposition of  $z$  into  $v + w$ , where  $v \in V$  and  $w \in W$ . Hence, with  $\|z\| = 1$ , we have

$$\begin{aligned} \langle z, \nabla^2 F_A(x) z \rangle &= \langle v, A v \rangle + \langle w, A w \rangle + \gamma \left(1 - \frac{1}{\|x\|}\right) + \frac{\gamma}{\|x\|} \left( \left\langle \frac{x}{\|x\|}, v \right\rangle + \left\langle \frac{x}{\|x\|}, w \right\rangle \right)^2 \\ &= a^2 \lambda_1 + b^2 \lambda_2 + \gamma \left(1 - \frac{1}{\|x\|}\right) + \frac{\gamma}{\|x\|} \left( a \left\langle \frac{x}{\|x\|}, \frac{v}{\|v\|} \right\rangle + b \left\langle \frac{x}{\|x\|}, \frac{w}{\|w\|} \right\rangle \right)^2, \end{aligned}$$

where  $\|v\| = a$  and  $\|w\| = b$  and  $a^2 + b^2 = 1$  and  $\lambda = \langle w, A w \rangle \geq \lambda_2$ .

If  $a \geq b$ , then  $a \geq \frac{1}{\sqrt{2}} \geq b$  and we have that for  $k \geq 2$ ,

$$\left( a \left\langle \frac{x}{\|x\|}, \frac{v}{\|v\|} \right\rangle + b \left\langle \frac{x}{\|x\|}, \frac{w}{\|w\|} \right\rangle \right)^2 \geq \left( a \left(1 - \frac{1}{k^2}\right) - b \sqrt{\frac{2}{k^2} - \frac{1}{k^4}} \right)^2,$$

which implies that for  $k \in \mathbb{N}$  with  $k \geq \frac{2(\gamma + \lambda_2)}{\lambda_2 - \lambda_1} > 1$ ,

$$\begin{aligned} \langle z, \nabla^2 F_A(x) z \rangle &\geq a^2(\lambda_1 + \gamma) + b^2 \left( \lambda_2 + \gamma - \frac{(\gamma + \lambda_1)k}{k-1} \right) + \frac{\gamma}{\|x\|} \left[ \frac{2(b^2 - a^2)}{k^2} - 2ab \sqrt{\frac{2}{k^2} - \frac{1}{k^4}} \right] \\ &\geq a^2(\lambda_1 + \gamma) + b^2 \frac{(\lambda_2 - \lambda_1)(\gamma + \lambda_2)}{2\gamma + \lambda_2 + \lambda_1} - \frac{(\gamma + \lambda_2)}{k-1} \left[ \frac{2(a^2 - b^2)}{k} + 2ab \sqrt{2 - \frac{1}{k^2}} \right] \\ &> a^2(\lambda_1 + \gamma) + b^2 \frac{(\lambda_2 - \lambda_1)}{2} - \frac{(\gamma + \lambda_2)}{k-1} \left[ \frac{2}{k} + \sqrt{2} \right] a^2 \\ &= a^2 \frac{(\lambda_1 + \gamma)}{2} + b^2 \frac{(\lambda_2 - \lambda_1)}{2} + \left( \frac{(\lambda_1 + \gamma)}{2} - \frac{(\gamma + \lambda_2)}{k-1} \left[ \frac{2}{k} + \sqrt{2} \right] \right) a^2. \end{aligned}$$

In addition, if  $k \geq \max\left(\frac{2(\gamma + \lambda_2)}{\lambda_2 - \lambda_1}, \frac{4(\gamma + \lambda_2)}{\gamma + \lambda_1} + 1\right)$ , then

$$\langle z, \nabla^2 F_A(x) z \rangle > a^2 \frac{(\lambda_1 + \gamma)}{2} + b^2 \frac{(\lambda_2 - \lambda_1)}{2} \geq \min\left(\frac{\gamma + \lambda_2}{4}, \frac{\gamma + \lambda_1}{2}\right).$$

On the other hand, if  $a < b$ , then  $a < \frac{1}{\sqrt{2}} < b$  and  $a^2 \lambda_1 + b^2 \lambda_2 > \frac{\lambda_1 + \lambda_2}{2}$  and we have that, with  $k > \frac{4\gamma + \lambda_2 + 3\lambda_1}{\lambda_2 - \lambda_1}$ ,

$$\langle z, \nabla^2 F_A(x) z \rangle \geq a^2 \lambda_1 + b^2 \lambda_2 + \gamma - \frac{(\gamma + \lambda_1)k}{k-1} > \frac{\lambda_1 + \lambda_2}{2} + \gamma - \frac{(\gamma + \lambda_1)k}{k-1} > \frac{\lambda_2 - \lambda_1}{4}.$$

Therefore, if  $k \geq \max(\frac{4(\gamma+\lambda_2)}{\lambda_2-\lambda_1}, 5 + \frac{4(\lambda_2-\lambda_1)}{\gamma+\lambda_1})$ , then  $\nabla^2 F_A(x) - \frac{\eta}{4}I$  is positive definite for  $\|x - x_*\| < \epsilon$ , where  $\eta = \min(2(\gamma + \lambda_1), \lambda_2 - \lambda_1)$ , i.e.,  $\mathcal{F}_A$  is  $\frac{\eta}{4}$ -strongly convex in

$$\{x \in \mathbb{R} : \|x - x_*\| < \epsilon\}.$$

□

### 2.1.3 Algorithm 2 : The Newton's Method

We have analyzed the gradient descent method to minimize the model (3) and confirmed that the method finds a global minimizer of (3) if an initial guess  $x_0$  is not orthogonal to the eigenspace corresponding to the smallest eigenvalue. However, it converges slowly. After realizing strongly convex regions near global minimizers in the previous sections, we now discuss the Newton's method applied to (3) to enhance the convergence rate. As a byproduct of the discussion, we will be able to provide a nonsingular linear system whose unique solution is an eigenvector of  $A$  when the corresponding eigenvalue of multiplicity 1 is known.

Since the functional (3) is continuously twice differentiable at  $x \neq 0$ , if we apply the Newton's method, we will get

$$x_{k+1} = x_k - (\nabla^2 F(x_k))^{-1} \nabla F(x_k), \quad k = 0, 1, 2, \dots, \quad (13)$$

with an initial  $x_0 \neq 0$  unless  $\nabla^2 F(x_k)$  is singular. Noting that

$$\begin{aligned} (\nabla^2 F(x))^{-1} \nabla F(x) &= \left[ A + \gamma \left( 1 - \frac{1}{\|x\|} \right) I + \frac{\gamma}{\|x\|^3} x x^T \right]^{-1} \left[ A + \gamma \left( 1 - \frac{1}{\|x\|} \right) I \right] x \\ &= x - \left[ A + \gamma \left( 1 - \frac{1}{\|x\|} \right) I + \frac{\gamma}{\|x\|^3} x x^T \right]^{-1} \left[ \frac{\gamma}{\|x\|} x \right] \\ &= x - \left[ \frac{1}{\gamma} A + \left( 1 - \frac{1}{\|x\|} \right) I + \frac{1}{\|x\|} \left( \frac{x}{\|x\|} \right) \left( \frac{x}{\|x\|} \right)^T \right]^{-1} \left[ \frac{x}{\|x\|} \right], \end{aligned}$$

we can observe that (13) becomes

$$x_{k+1} = \left[ \frac{1}{\gamma} A + \left( 1 - \frac{1}{\|x_k\|} \right) I + \frac{1}{\|x_k\|} \left( \frac{x_k}{\|x_k\|} \right) \left( \frac{x_k}{\|x_k\|} \right)^T \right]^{-1} \frac{x_k}{\|x_k\|}, \quad k = 0, 1, 2, \dots$$

Hence, we propose the following scheme:

- Initialize  $x_0 \neq 0$ ,
- For  $k = 0, 1, 2, \dots$ , compute  $y_k$  and  $x_{k+1}$  by

$$y_k = \frac{x_k}{\|x_k\|},$$

and

$$\left[ \frac{1}{\gamma} A + \left( 1 - \frac{1}{\|x_k\|} \right) I + \frac{1}{\|x_k\|} \left( \frac{x_k}{\|x_k\|} \right) \left( \frac{x_k}{\|x_k\|} \right)^T \right] x_{k+1} = y_k. \quad (14)$$

We wrote (14) in the given form to enhance its similarity to the inverse power method. Under the assumptions of Theorem 4, choosing an initial  $x_0$  in the strongly convex region near a global minimum point  $x_*$ , the above algorithm presents, at least, the quadratic convergence. Moreover, unlike the Inverse Rayleigh Quotient Iteration, the system of linear equations in

(14) that we need to solve near  $x_*$  is far from being singular because  $\nabla^2 F(x)$  is  $\frac{\eta}{4}$ -strongly convex in the region where we apply the above algorithm.

We will now show that the Newton's method discussed above converges to a critical point under a certain condition. In fact, the following theorem says that convergence of the norm  $\|x_k\|$  is equivalent to convergence of  $x_k$ .

**Theorem 5.** *Let  $A$  be an  $N \times N$  nonzero real symmetric matrix with eigenvalues  $\lambda_1 \leq \dots \leq \lambda_N$ . Let  $\gamma > \max(0, -\lambda_1)$ , and  $x_0 \neq 0$  with  $\|x_0\| \neq \frac{\gamma}{\gamma + \lambda_j}$  for any  $1 \leq j \leq N$ .*

1. *Suppose that a sequence  $\{x_k\}$  can be generated by (14) and that  $\|x_k\|$  converges to  $\eta > 0$ . If  $\|x_k\| \neq \frac{\gamma}{\gamma + \lambda_j}$  for any  $1 \leq j \leq N$ ,  $k \in \mathbb{N} \cup \{0\}$ , then there exists  $1 \leq i_0 \leq N$  such that  $\eta = \frac{\gamma}{\gamma + \lambda_{i_0}}$  and  $x_k$  converges to an eigenvector  $x_*$  of  $A$  corresponding to the eigenvalue  $\lambda_{i_0}$  with  $\|x_*\| = \eta$ .*
2. *On the other hand, we assume that  $\lambda_i$  for some  $1 \leq i \leq N$  has multiplicity 1 and  $q_i$  is a unit eigenvector of  $A$  corresponding to  $\lambda_i$  and  $\|x_{k_0}\| = \frac{\gamma}{\gamma + \lambda_i}$  for some  $k_0 \in \mathbb{N}$ . If  $x_{k_0}$  is not a critical point of  $F_A$  with  $|q_i^T x_{k_0}| > 0$ , then  $x_{k_0+1}$  is an eigenvector of  $A$  corresponding to the eigenvalue  $\lambda_i$ . If  $|q_i^T x_{k_0}| \neq \frac{\gamma(\gamma + \lambda_j)}{(\gamma + \lambda_i)^2}$  for  $j < i$ , then  $x_{k_0+2}$  is a critical point of  $F_A$ , i.e., an eigenvector of  $A$  corresponding to the eigenvalue  $\lambda_i$  with norm  $\|x_{k_0+2}\| = \frac{\gamma}{\gamma + \lambda_i}$ . However, if  $|q_i^T x_{k_0}| = \frac{\gamma(\gamma + \lambda_j)}{(\gamma + \lambda_i)^2}$  for some  $j < i$ , then the system becomes singular and we may not compute  $x_{k_0+2}$  uniquely. In any case, the algorithm terminates in  $k_0 + 2$  iterations.*

*Proof.* It suffices to consider the case that  $A$  is a diagonal matrix with diagonal entries  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$ . Then, (14) becomes

$$\frac{1}{\gamma} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{bmatrix} \begin{bmatrix} x_{k+1,1} \\ \vdots \\ x_{k+1,N} \end{bmatrix} + \left(1 - \frac{1}{\|x_k\|}\right) \begin{bmatrix} x_{k+1,1} \\ \vdots \\ x_{k+1,N} \end{bmatrix} + \frac{1}{\|x_k\|} (y_k^T x_{k+1}) y_k = y_k, \quad (15)$$

where  $x_{k+1} = [x_{k+1,1} \ \dots \ x_{k+1,N}]^T$  and  $y_k = \frac{x_k}{\|x_k\|}$ . Let  $\{x_k\}$  be a sequence generated by (15) with  $\|x_k\| \neq \frac{\gamma}{\gamma + \lambda_j}$  for any  $1 \leq j \leq N$  and for all  $k \in \mathbb{N} \cup \{0\}$ .

Firstly, we consider that  $\|x_k\|$  converges to  $\eta > 0$ . Suppose that  $\eta \neq \frac{\gamma}{\gamma + \lambda_j}$  for any  $1 \leq j \leq N$ . Note that for  $1 \leq j \leq N$ , we have

$$\begin{aligned} \left(1 + \frac{\lambda_j}{\gamma} - \frac{1}{\|x_k\|}\right) x_{k+1,j} &= y_{k,j} \left(1 - \frac{\|x_{k+1}\|}{\|x_k\|} (y_k^T y_{k+1})\right) \\ \Leftrightarrow \frac{x_{k+1,j}}{\|x_{k+1}\|} &= \left[ \frac{\frac{\|x_k\|}{\|x_{k+1}\|} - (y_k^T y_{k+1})}{\|x_k\| \left(1 + \frac{\lambda_j}{\gamma} - 1\right)} \right] \frac{x_{k,j}}{\|x_k\|}. \end{aligned} \quad (16)$$

Since  $\eta > 0$  and  $\left(1 + \frac{\lambda_j}{\gamma} - \frac{1}{\|x_k\|}\right) \neq 0$  for any  $1 \leq j \leq N$  and  $k \in \mathbb{N} \cup \{0\}$ , we know that  $\|x_k\| \neq 0$  for all  $k \geq 0$ . This implies that

$$1 - \frac{\|x_{k+1}\|}{\|x_k\|} (y_k^T y_{k+1}) \neq 0, \quad k \geq 0.$$

If we set  $\mathcal{J}_0 := \{j \in \{1, 2, \dots, N\} : x_{0,j} \neq 0\}$ , then we can see that for  $k \geq 0$ ,

$$x_{k,j} \neq 0 \text{ if and only if } j \in \mathcal{J}_0.$$

We will now show that  $\limsup_k y_k^T y_{k+1} = 1$ . Suppose that  $\limsup_k y_k^T y_{k+1} = \epsilon < 1$ . Given  $0 < \delta < 1 - \epsilon$ , we may choose  $l_1 \in \mathbb{N}$  so that  $k \geq l_1$  implies

$$\|x_k\| - \eta < \frac{\delta}{2} \quad \text{and} \quad \left| \frac{\|x_k\|}{\|x_{k+1}\|} - 1 \right| < \frac{\delta}{2} \quad \text{and} \quad y_k^T y_{k+1} < \epsilon + \frac{\delta}{2}. \quad (17)$$

We also choose  $J \in \mathcal{J}_0$  satisfying

$$\left| \eta \left( 1 + \frac{\lambda_J}{\gamma} \right) - 1 \right| \leq \min_{j \in \mathcal{J}_0} \left| \eta \left( 1 + \frac{\lambda_j}{\gamma} \right) - 1 \right|.$$

From (16), we see that for  $k \geq l_1$ ,

$$\frac{|x_{k+1,J}|}{\|x_{k+1}\|} \geq \frac{1 - \delta - \epsilon}{\left| \eta \left( 1 + \frac{\lambda_J}{\gamma} \right) - 1 \right| + \frac{\delta}{2} \left( 1 + \frac{\lambda_J}{\gamma} \right)} \frac{|x_{k,J}|}{\|x_k\|}. \quad (18)$$

If  $\left| \eta \left( 1 + \frac{\lambda_J}{\gamma} \right) - 1 \right| < 1 - \epsilon$ , then with  $\delta > 0$  satisfying

$$\left| \eta \left( 1 + \frac{\lambda_J}{\gamma} \right) - 1 \right| + \frac{\delta}{2} \left( 1 + \frac{\lambda_J}{\gamma} \right) < 1 - \delta - \epsilon,$$

we can see from (18) that  $\lim_{k \rightarrow \infty} \frac{|x_{k+1,J}|}{\|x_{k+1}\|} = \infty$ , which is impossible. Hence,

$$\min_{j \in \mathcal{J}_0} \left| \eta \left( 1 + \frac{\lambda_j}{\gamma} \right) - 1 \right| \geq 1 - \epsilon$$

i.e.,

$$\eta \leq \frac{\epsilon \gamma}{\gamma + \lambda^*} \quad \text{or} \quad \eta \geq \frac{(2 - \epsilon) \gamma}{\gamma + \lambda^*},$$

where  $\lambda_* = \min_{j \in \mathcal{J}_0} \lambda_j$  and  $\lambda^* = \max_{j \in \mathcal{J}_0} \lambda_j$ .

Suppose that  $\eta \leq \frac{\epsilon \gamma}{\gamma + \lambda^*}$ . For any  $0 < \delta < \frac{(1 - \epsilon) \gamma}{2\gamma + \lambda^*} < 1 - \epsilon$ , we can see from (17) that  $k \geq l_1$  implies that for each  $j \in \mathcal{J}_0$ , we have

$$\|x_k\| \left( 1 + \frac{\lambda_j}{\gamma} \right) - 1 < \left( \eta + \frac{\delta}{2} \right) \left( 1 + \frac{\lambda_j}{\gamma} \right) - 1 \leq \frac{\epsilon(\gamma + \lambda_j)}{\gamma + \lambda^*} + \frac{\delta}{2} \left( 1 + \frac{\lambda_j}{\gamma} \right) - 1 < 0. \quad (19)$$

In addition, from (17) we know that for  $k \geq l_1$ ,

$$\frac{\|x_k\|}{\|x_{k+1}\|} - (y_k^T y_{k+1}) > 1 - \frac{\delta}{2} - (y_k^T y_{k+1}) > 0,$$

since  $\epsilon + \delta < 1$  implies

$$y_k^T y_{k+1} < \epsilon + \frac{\delta}{2} < 1 - \frac{\delta}{2}.$$

This results in, together with (19),

$$\begin{aligned} y_k^T y_{k+1} &= \left( \frac{\|x_k\|}{\|x_{k+1}\|} - (y_k^T y_{k+1}) \right) \left( \sum_{j=1}^N \frac{y_{k,j}^2}{\|x_k\| \left( 1 + \frac{\lambda_j}{\gamma} \right) - 1} \right) \\ &\leq \left( 1 - \frac{\delta}{2} - (y_k^T y_{k+1}) \right) \left( \sum_{j \in \mathcal{J}_0} \frac{y_{k,j}^2}{\|x_k\| \left( 1 + \frac{\lambda_j}{\gamma} \right) - 1} \right) < 0. \end{aligned}$$

This implies that  $\epsilon \leq 0$ , i.e.,  $\eta \leq 0$ . which is a contradiction. Hence, we must have  $\eta \geq \frac{(2-\epsilon)\gamma}{\gamma+\lambda_*}$ .

With  $0 < \delta < \frac{(1-\epsilon)\gamma}{2\gamma+\lambda_*} < 1 - \epsilon$ , we can also see that for  $k \geq l_1$ , and for each  $j \in \mathcal{J}_0$ ,

$$\begin{aligned} \|x_k\| \left(1 + \frac{\lambda_j}{\gamma}\right) - 1 &> \left(\eta - \frac{\delta}{2}\right) \left(1 + \frac{\lambda_j}{\gamma}\right) - 1 \\ &\geq \frac{(2-\epsilon)(\gamma + \lambda_j)}{\gamma + \lambda_*} - \frac{\delta}{2} \left(1 + \frac{\lambda_j}{\gamma}\right) - 1 \\ &> \frac{(2-\epsilon)(\gamma + \lambda_j)}{\gamma + \lambda_*} - \frac{(1-\epsilon)(\gamma + \lambda_j)}{2(2\gamma + \lambda_*)} - 1 > 0. \end{aligned}$$

Together with (17), we can conclude that for  $k \geq l_1$ , since  $\epsilon + \delta < 1$ ,

$$\begin{aligned} y_k^T y_{k+1} &= \left(\frac{\|x_k\|}{\|x_{k+1}\|} - (y_k^T y_{k+1})\right) \left(\sum_{j=1}^N \frac{y_{k,j}^2}{\|x_k\| \left(1 + \frac{\lambda_j}{\gamma}\right) - 1}\right) \\ &\geq \left(1 - \frac{\delta}{2} - (y_k^T y_{k+1})\right) \left(\sum_{j \in \mathcal{J}_0} \frac{y_{k,j}^2}{\|x_k\| \left(1 + \frac{\lambda_j}{\gamma}\right) - 1}\right) > 0. \end{aligned}$$

Hence, we have

$$0 \leq \epsilon < 1. \quad (20)$$

If we extract a subsequence  $y_{k_n}$  such that  $y_{k_n}^T y_{k_n+1} \rightarrow \epsilon$  as  $n \rightarrow \infty$ , then using the form of (14) and knowing that

$$\liminf_{n \rightarrow \infty} \frac{y_{k_n+1}^T}{\|x_{k_n+1}\|} \frac{1}{\gamma} A x_{k_n+1} = \liminf_{n \rightarrow \infty} \frac{1}{\gamma} y_{k_n+1}^T A y_{k_n+1} \geq \frac{\lambda_*}{\gamma}$$

we can see that

$$\begin{aligned} &\frac{y_{k_n+1}^T}{\|x_{k_n+1}\|} \left[ \frac{1}{\gamma} A + \left(1 - \frac{1}{\|x_{k_n}\|}\right) I + \frac{1}{\|x_{k_n}\|} \left(\frac{x_{k_n}}{\|x_{k_n}\|}\right) \left(\frac{x_{k_n}}{\|x_{k_n}\|}\right)^T \right] x_{k_n+1} = \frac{y_{k_n+1}^T y_{k_n}}{\|x_{k_n+1}\|} \\ \Leftrightarrow &\frac{1}{\gamma} y_{k_n+1}^T A y_{k_n+1} + \left(1 - \frac{1}{\|x_{k_n}\|}\right) + \frac{1}{\|x_{k_n}\|} (y_{k_n}^T y_{k_n+1})^2 = \frac{y_{k_n+1}^T y_{k_n}}{\|x_{k_n+1}\|} \\ \Rightarrow &\frac{\lambda_*}{\gamma} + \left(1 - \frac{1}{\eta}\right) + \frac{\epsilon^2}{\eta} \leq \frac{\epsilon}{\eta} \\ \Leftrightarrow &\eta \leq \frac{(1 + \epsilon - \epsilon^2)\gamma}{\gamma + \lambda_*}. \end{aligned}$$

Hence, we have

$$\frac{(2-\epsilon)\gamma}{\gamma + \lambda_*} \leq \eta \leq \frac{(1 + \epsilon - \epsilon^2)\gamma}{\gamma + \lambda_*}. \quad (21)$$

However, (20) implies  $1 + \epsilon - \epsilon^2 < 2 - \epsilon$ , which contradicts (21). Therefore, we conclude that  $\epsilon = 1$ , i.e.,  $\limsup_k y_k^T y_{k+1} = 1$ .

By considering a subsequence  $y_{k_n}$  such that  $\lim_{n \rightarrow \infty} y_{k_n}^T y_{k_n+1} = 1$ , we can also see from (16) that

$$\begin{aligned} \infty &= \lim_{n \rightarrow \infty} \frac{1}{\left(\frac{\|x_{k_n}\|}{\|x_{k_n+1}\|} - (y_{k_n}^T y_{k_n+1})\right)^2} = \lim_{n \rightarrow \infty} \sum_{j=1}^N \frac{y_{k_n,j}^2}{\left(\|x_{k_n}\| \left(1 + \frac{\lambda_j}{\gamma}\right) - 1\right)^2} \\ &\leq \max_{1 \leq j \leq N} \frac{1}{\left(\eta \left(1 + \frac{\lambda_j}{\gamma}\right) - 1\right)^2} < \infty, \end{aligned}$$

which is a contradiction. Therefore, we conclude that  $\eta = \frac{\gamma}{\gamma + \lambda_{i_0}}$  for some  $1 \leq i_0 \leq N$ .

We will now show that  $x_k$  converges to an eigenvector  $x_*$  of  $A$  corresponding to  $\lambda_{i_0}$  with norm  $\|x_*\| = \eta$ . Firstly, we show that there exists  $j \in \mathcal{J}_0$  such that  $\lambda_j = \lambda_{i_0}$ . It is easy to see that  $\limsup_{k \rightarrow \infty} y_k^T y_{k+1} = 1$ , whose proof is almost exactly the same as above and is easier using  $\eta = \frac{\gamma}{\gamma + \lambda_{i_0}}$ . So we will omit it. By choosing a subsequence  $y_{k_n}$  such that  $y_{k_n}^T y_{k_n+1} \rightarrow 1$  as  $n \rightarrow \infty$ , we can also see that

$$\infty = \lim_{n \rightarrow \infty} \frac{1}{\left(\frac{\|x_{k_n}\|}{\|x_{k_n+1}\|} - (y_{k_n}^T y_{k_n+1})\right)^2} = \lim_{n \rightarrow \infty} \sum_{j=1}^N \frac{y_{k_n, j}^2}{(\|x_{k_n}\| (1 + \frac{\lambda_j}{\gamma}) - 1)^2}.$$

Since  $\|x_{k_n}\| \rightarrow \eta$ , there must be  $j \in \mathcal{J}_0$  with  $\lambda_j = \lambda_{i_0}$ . That is,  $\{j \in \mathcal{J}_0 : \lambda_j = \lambda_{i_0}\} \neq \emptyset$ . Hence, without loss of generality we will say that  $i_0 \in \mathcal{J}_0$ .

Let  $k_0 \in \mathbb{N}$  be such that  $k \geq k_0$  implies

$$\left|1 + \frac{\lambda_{i_0}}{\gamma} - \frac{1}{\|x_k\|}\right| < \min_{\lambda_j \neq \lambda_{i_0}} \frac{|\lambda_j - \lambda_{i_0}|}{3\gamma}.$$

Note that  $k \geq k_0$ , and for  $\lambda_j \neq \lambda_{i_0}$ ,

$$\left|\frac{1 + \frac{\lambda_{i_0}}{\gamma} - \frac{1}{\|x_k\|}}{1 + \frac{\lambda_j}{\gamma} - \frac{1}{\|x_k\|}}\right| < \frac{1}{2}.$$

Moreover, for  $1 \leq j \leq N$ , we can represent

$$\frac{x_{K+1, j}}{\|x_{K+1}\|} = \left(\prod_{k=0}^K \left[\frac{1}{\|x_{k+1}\|} - \frac{1}{\|x_k\|} (y_k^T y_{k+1})\right]\right) \frac{x_{0, j}}{\|x_0\|}, \quad K \geq 0. \quad (22)$$

Since  $i_0 \in \mathcal{J}_0$ , we have that for  $\lambda_j \neq \lambda_{i_0}$ , as  $K \rightarrow \infty$ ,

$$\begin{aligned} \left|\frac{x_{K+1, j}}{x_{K+1, i_0}}\right| &= \left(\prod_{k=k_0}^K \left|\frac{1 + \frac{\lambda_{i_0}}{\gamma} - \frac{1}{\|x_k\|}}{1 + \frac{\lambda_j}{\gamma} - \frac{1}{\|x_k\|}}\right|\right) \left(\prod_{k=0}^{k_0-1} \left|\frac{1 + \frac{\lambda_{i_0}}{\gamma} - \frac{1}{\|x_k\|}}{1 + \frac{\lambda_j}{\gamma} - \frac{1}{\|x_k\|}}\right|\right) \left|\frac{x_{0, j}}{x_{0, i_0}}\right| \\ &\leq \frac{1}{2^{K-k_0+1}} \left|\frac{x_{0, j}}{x_{0, i_0}}\right| \rightarrow 0. \end{aligned}$$

Therefore, as  $K \rightarrow \infty$ ,

$$(x_{K+1, i_0})^2 \left( \sum_{\substack{j \in \mathcal{J}_0, \\ \lambda_j = \lambda_{i_0}}} \left(\frac{x_{K+1, j}}{x_{K+1, i_0}}\right)^2 + \sum_{\lambda_j \neq \lambda_{i_0}} \left(\frac{x_{K+1, j}}{x_{K+1, i_0}}\right)^2 \right) = \|x_{K+1}\|^2 \rightarrow \eta^2 = \left(\frac{\gamma}{\gamma + \lambda_{i_0}}\right)^2,$$

which implies that as  $K \rightarrow \infty$ , we have

$$(x_{K+1, i_0})^2 \sum_{\substack{j \in \mathcal{J}_0, \\ \lambda_j = \lambda_{i_0}}} \left(\frac{x_{k_0, j}}{x_{k_0, i_0}}\right)^2 \rightarrow \left(\frac{\gamma}{\gamma + \lambda_{i_0}}\right)^2.$$

Let  $m = \left(\sum_{\substack{j \in \mathcal{J}_0, \\ \lambda_j = \lambda_{i_0}}} \left(\frac{x_{k_0, j}}{x_{k_0, i_0}}\right)^2\right)^{\frac{1}{2}} \geq 1$ . Then, we can see that as  $K \rightarrow \infty$ ,

$$(x_{K+1, i_0})^2 \rightarrow \frac{1}{m^2} \left(\frac{\gamma}{\gamma + \lambda_{i_0}}\right)^2. \quad (23)$$

This implies that  $y_{K+1,i_0}^2 \rightarrow \frac{1}{m^2}$  as  $K \rightarrow \infty$ . Noting that for all  $k$ ,

$$\begin{aligned} 1 &\geq |y_k^T y_{k+1}| = \left| \frac{\|x_k\|}{\|x_{k+1}\|} - (y_k^T y_{k+1}) \right| \left| \sum_{j=1}^N \frac{y_{k,j}^2}{\|x_k\|(1 + \frac{\lambda_j}{\gamma}) - 1} \right| \\ &\geq \left| \frac{\|x_k\|}{\|x_{k+1}\|} - (y_k^T y_{k+1}) \right| \left( \left| \sum_{\lambda_j = \lambda_{i_0}} \frac{y_{k,j}^2}{\|x_k\|(1 + \frac{\lambda_{i_0}}{\gamma}) - 1} \right| - \left| \sum_{\lambda_j \neq \lambda_{i_0}} \frac{y_{k,j}^2}{\|x_k\|(1 + \frac{\lambda_j}{\gamma}) - 1} \right| \right), \end{aligned}$$

and

$$\lim_{k \rightarrow \infty} \left( \left| \sum_{\lambda_j = \lambda_{i_0}} \frac{y_{k,j}^2}{\|x_k\|(1 + \frac{\lambda_{i_0}}{\gamma}) - 1} \right| - \left| \sum_{\lambda_j \neq \lambda_{i_0}} \frac{y_{k,j}^2}{\|x_k\|(1 + \frac{\lambda_j}{\gamma}) - 1} \right| \right) = \infty,$$

we can see that

$$\lim_{k \rightarrow \infty} \left| \frac{\|x_k\|}{\|x_{k+1}\|} - (y_k^T y_{k+1}) \right| = 0,$$

i.e.,  $y_k^T y_{k+1}$  converges to 1 as  $k \rightarrow \infty$ . Since (22) implies that for  $j \in \mathcal{J}_0$  with  $\lambda_j = \lambda_{i_0}$ ,

$$x_{k,j} = x_{k,i_0} \begin{pmatrix} x_{0,j} \\ x_{0,i_0} \end{pmatrix},$$

we have that as  $k \rightarrow \infty$ ,

$$\begin{aligned} y_k^T y_{k+1} &= \frac{1}{\|x_k\| \|x_{k+1}\|} \sum_{j=1}^N x_{k,j} x_{k+1,j} \\ &= \frac{1}{\|x_k\| \|x_{k+1}\|} \left( \sum_{\substack{j \in \mathcal{J}_0, \\ \lambda_j = \lambda_{i_0}}} x_{k,i_0} x_{k+1,i_0} \left( \frac{x_{0,j}}{x_{0,i_0}} \right)^2 + \sum_{\substack{j \in \mathcal{J}_0, \\ \lambda_j \neq \lambda_{i_0}}} x_{k,j} x_{k+1,j} \right) \\ &\rightarrow \left( \frac{\gamma + \lambda_{i_0}}{\gamma} \right)^2 m^2 \lim_{k \rightarrow \infty} (x_{k,i_0} x_{k+1,i_0}) = 1. \end{aligned}$$

Together with (23), we can easily see that  $x_{k,i_0}$  converges either to  $\frac{1}{m} \frac{\gamma}{\gamma + \lambda_{i_0}}$  or  $-\frac{1}{m} \frac{\gamma}{\gamma + \lambda_{i_0}}$ .

Therefore,  $x_k$  converges to  $x_*$  with

$$x_{*,j} = \begin{cases} \begin{pmatrix} x_{0,j} \\ x_{0,i_0} \end{pmatrix} x_{*,i_0}, & \text{for } j \in \mathcal{J}_0 \text{ and } \lambda_j = \lambda_{i_0}, \\ 0, & \text{otherwise,} \end{cases}$$

where  $x_{*,i_0} = \pm \frac{1}{m} \frac{\gamma}{\gamma + \lambda_{i_0}}$ . Note that  $x_*$  is an eigenvector of  $A$  corresponding to  $\lambda_{i_0}$  with norm  $\|x_*\| = \eta = \frac{\gamma}{\gamma + \lambda_{i_0}}$ .

Secondly, we consider that  $\lambda_i$  for some  $1 \leq i \leq N$  has multiplicity 1 and  $\|x_{k_0}\| = \frac{\gamma}{\gamma + \lambda_i}$  for some  $k_0 \in \mathbb{N}$ . If  $x_{k_0}$  is not a critical point of  $F_A$  and  $x_{k_0,i} \neq 0$ , then we have that  $|x_{k_0,i}| < \frac{\gamma}{\gamma + \lambda_i}$  and  $\frac{1}{\|x_{k_0}\|} (y_k^T x_{k_0+1}) = 1$  in (15), which leads to

$$\frac{1}{\gamma} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_N \end{bmatrix} \begin{bmatrix} x_{k_0+1,1} \\ \vdots \\ x_{k_0+1,N} \end{bmatrix} = \frac{\lambda_i}{\gamma} \begin{bmatrix} x_{k_0+1,1} \\ \vdots \\ x_{k_0+1,N} \end{bmatrix}. \quad (24)$$

Hence, there exists a unique solution  $x_{k_0+1}$ , which is  $x_{k_0+1} = \alpha \mathbf{e}_i$ , where  $\mathbf{e}_i$  is the standard basis element in  $\mathbb{R}^N$  with  $\mathbf{e}_{i,j} = \delta_{ij}$  and  $\alpha = \frac{\gamma^2}{x_{k_0,i}(\gamma + \lambda_i)^2}$ . Since  $|x_{k_0,i}| < \frac{\gamma}{\gamma + \lambda_i}$ , we have that

$|\alpha| > \frac{\gamma}{\gamma + \lambda_i}$  and  $x_{k_0+1}$  is an eigenvector of  $A$  corresponding to the eigenvalue  $\lambda_i$ , which is not a critical point of  $F_A$ .

Since  $x_{k_0+2}$  satisfies (15) and  $y_{k_0+1}$  is either  $\mathbf{e}_i$  or  $-\mathbf{e}_i$ , if we multiply (15) by  $\mathbf{e}_j^T$  from the left, then we have

$$\begin{cases} \left(\frac{\lambda_j}{\gamma} + 1 - \frac{1}{|\alpha|}\right)x_{k_0+2,j} + \frac{\delta_{ij}}{|\alpha|}x_{k_0+2,i} = \delta_{ij}, & \text{if } y_{k_0+1} = \mathbf{e}_i, \\ \left(\frac{\lambda_j}{\gamma} + 1 - \frac{1}{|\alpha|}\right)x_{k_0+2,j} + \frac{\delta_{ij}}{|\alpha|}x_{k_0+2,i} = -\delta_{ij}, & \text{if } y_{k_0+1} = -\mathbf{e}_i. \end{cases}$$

Note that  $|\alpha| = \frac{\gamma}{\gamma + \lambda_j}$  is equivalent to  $|x_{k_0,i}| = \frac{\gamma(\gamma + \lambda_j)}{(\gamma + \lambda_i)^2}$ . Since  $|x_{k_0,i}| < \frac{\gamma}{\gamma + \lambda_i}$ , we know that it is possible to have  $|\alpha| = \frac{\gamma}{\gamma + \lambda_j}$  only if  $j < i$ .

Hence, if  $|x_{k_0,i}| \neq \frac{\gamma(\gamma + \lambda_j)}{(\gamma + \lambda_i)^2}$  for  $j < i$ , then  $|\alpha| \neq \frac{\gamma}{\gamma + \lambda_j}$  for  $j \neq i$ , and we have

$$x_{k_0+2,j} = \begin{cases} \pm \frac{\gamma}{\gamma + \lambda_i}, & \text{if } j = i, \\ 0, & \text{if } j \neq i. \end{cases}$$

depending on  $y_{k_0+1} = \pm \mathbf{e}_i$ . That is,  $x_{k_0+2}$  is a critical point of  $A$ , that is, an eigenvector of  $A$  corresponding to the eigenvalue  $\lambda_i$  with norm  $\|x_{k_0+2}\| = \frac{\gamma}{\gamma + \lambda_i}$  and the algorithm terminates.

On the other hand, if  $|x_{k_0,i}| = \frac{\gamma(\gamma + \lambda_j)}{(\gamma + \lambda_i)^2}$  for some  $j < i$ , then  $|\alpha| = \frac{\gamma}{\gamma + \lambda_j}$  and the system becomes singular and the algorithm terminates.  $\square$

#### 2.1.4 A Nonsingular Linear System for Eigenvector Estimation

Inspired by the Newton's method applied to our proposed model (3) with the convergence results, we realize that we can propose a nonsingular linear system for eigenvector estimation.

It is well known that an eigenvector of  $A$  is a solution to a singular linear system, which makes it difficult to find a corresponding eigenvector, and iterative methods are designed to approximate eigenvectors. However, our nonsingular system guarantees to find a corresponding eigenvector by solving the system once when knowing that the eigenvalue has multiplicity 1. Moreover, even if we have an estimated eigenvalue, we can guaranteed that the error in the eigenvector estimation is comparable with the error from the estimated eigenvalue. For eigenvalues with multiplicity greater than 1, the same form of a nonsingular system provides an estimated eigenvector within arbitrarily small error.

**One Step Eigenvector Estimation.** Given an  $N \times N$  symmetric matrix  $A$ , and an eigenvalue  $\tilde{\lambda}$  of  $A$ , and  $\gamma > 0$  with  $\gamma \neq -\tilde{\lambda}$ , we choose  $x_0$  uniformly at random from  $S^{N-1}$  and solve

$$(A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})x_0x_0^T)x = \gamma x_0. \quad (25)$$

The equation (25), which is what we propose to solve for an eigenvector of  $A$  corresponding to  $\tilde{\lambda}$ , is inspired by (14) and the proof of the second part of Theorem 5. By dividing (25) by  $\gamma$ , we have the same form as in (14). The difference is the condition on  $\gamma$ . In solving (25), any  $\gamma > 0$  with  $\gamma \neq -\tilde{\lambda}$  works. The following proposition says that (25) is, indeed, a nonsingular linear system with probability 1, whose unique solution is a corresponding eigenvector.

**Proposition 1.** *Suppose that  $\tilde{\lambda}$  has multiplicity 1. With probability 1, the equation (25) has a unique solution  $\tilde{x}$  that is an eigenvector of  $A$  corresponding to the eigenvalue  $\tilde{\lambda}$ .*

*Proof.* Let  $q$  be a unit eigenvector of  $A$  corresponding to  $\tilde{\lambda}$ . Note that if we choose  $x_0 \in S^{N-1}$  uniformly at random, then we have  $q^T x_0 \neq 0$  with probability 1. Moreover, if  $q^T x_0 \neq 0$ , then  $(A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0x_0^T))z = 0$  implies  $x_0^T z = 0$ . Hence, we have  $Az = \tilde{\lambda}z$ . Since  $\tilde{\lambda}$  has



multiplicity 1,  $z = aq$  for some  $a \in \mathbb{R}$ . In addition, since  $0 = z^T x_0 = aq^T x_0$  and  $q^T x_0 \neq 0$ , we must have  $a = 0$ . That is,  $z = 0$ . Hence,  $A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T)$  is nonsingular and there exists a unique solution  $\tilde{x}$  to (25). By multiplying (25) by  $q^T$ , we have  $(\gamma + \tilde{\lambda})x_0^T \tilde{x} = \gamma$ , which implies that  $\tilde{x}$  also satisfies

$$(A - \tilde{\lambda}I)\tilde{x} = 0.$$

□

For practical reasons, when only approximate eigenvalues are available, a good estimate of an eigenvalue guarantees a good estimate of a corresponding eigenvalue using a nonsingular linear system in the same form as (25).

**Proposition 2.** *Let  $A$  be an  $N \times N$  symmetric matrix. Let  $\tilde{\lambda}$  be an eigenvalue of  $A$  with multiplicity 1. Let  $\gamma > 0$  be such that  $\gamma \neq -\tilde{\lambda}$ . Let  $x_0 \in S^{N-1}$  be such that  $0 < |q^T x_0| < 1$ , where  $q$  is a unit eigenvector of  $A$  corresponding to  $\tilde{\lambda}$ . Then, for  $\lambda$  close enough to  $\tilde{\lambda}$ , but not equal,  $A - \lambda I + (\gamma + \lambda)(x_0 x_0^T)$  is nonsingular and there exist  $0 < c \leq C$  such that*

$$c|\lambda - \tilde{\lambda}| < \|x - \tilde{x}\| < C|\lambda - \tilde{\lambda}|,$$

where  $x$  and  $\tilde{x}$  are the solutions to (25), respectively, i.e.,

$$(A - \lambda I + (\gamma + \lambda)(x_0 x_0^T))x = \gamma x_0 \quad \text{and} \quad (A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T))\tilde{x} = \gamma x_0.$$

*Proof.* Since  $q^T x_0 \neq 0$  implies, by Proposition 1, that  $A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T)$  is nonsingular, there exists  $\epsilon > 0$  such that if  $|\lambda - \tilde{\lambda}| < \epsilon$ ,

$$A - \lambda I + (\gamma + \lambda)(x_0 x_0^T)$$

is nonsingular.

Let  $\alpha, \beta > 0$  be the smallest and the largest singular values of  $A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T)$ . Note that

$$\alpha\|\tilde{x}\| \leq \|(A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T))\tilde{x}\| = \gamma \leq \beta\|\tilde{x}\| \Rightarrow \alpha \leq \frac{\gamma}{\|\tilde{x}\|} \leq \beta,$$

and

$$\alpha\|x\| \leq \|(A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T))x\| = \|\gamma x_0 + (\lambda - \tilde{\lambda})(I - x_0 x_0^T)x\| \leq \gamma + |\lambda - \tilde{\lambda}|\|x\|.$$

So, if  $|\lambda - \tilde{\lambda}| < \min(\frac{\alpha}{2}, \epsilon)$ , then  $\|x\| < \frac{2\gamma}{\alpha}$ . In addition, noting that

$$(A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T))(x - \tilde{x}) = (\lambda - \tilde{\lambda})(I - x_0 x_0^T)x, \quad (26)$$

we obtain that for  $\|x\| < \frac{2\gamma}{\alpha}$ ,

$$\alpha\|x - \tilde{x}\| \leq \|(A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T))(x - \tilde{x})\| \leq |\lambda - \tilde{\lambda}|\|x\| < \frac{2\gamma}{\alpha}|\lambda - \tilde{\lambda}|.$$

Hence, if we set

$$C = \frac{2\gamma}{\alpha^2} > 0,$$

and if  $|\lambda - \tilde{\lambda}| < \min(\epsilon, \frac{\alpha}{2})$ , then we have

$$\|x - \tilde{x}\| < C|\lambda - \tilde{\lambda}|$$

Secondly, since (26) is the same as

$$(\lambda - \tilde{\lambda})(I - x_0 x_0^T)(\tilde{x} - x) + (A - \tilde{\lambda}I + (\gamma + \tilde{\lambda}(x_0 x_0^T)))(x - \tilde{x}) = (\lambda - \tilde{\lambda})(I - x_0 x_0^T)\tilde{x},$$

we have

$$(|\lambda - \tilde{\lambda}| + \beta)\|x - \tilde{x}\| \geq |\lambda - \tilde{\lambda}|\|\tilde{x} - (x_0^T \tilde{x})x_0\|.$$

Therefore, if  $|\lambda - \tilde{\lambda}| < \min(\epsilon, \frac{\alpha}{2})$ , then we have

$$\left(\frac{\alpha}{2} + \beta\right)\|x - \tilde{x}\| > |\lambda - \tilde{\lambda}|\|\tilde{x} - (x_0^T \tilde{x})x_0\|$$

implying

$$c|\lambda - \tilde{\lambda}| < \|x - \tilde{x}\|$$

with

$$c = \frac{2\|\tilde{x} - (x_0^T \tilde{x})x_0\|}{(\alpha + 2\beta)} > 0$$

since  $|q^T x_0| < 1$  implies  $\tilde{x} \neq (x_0^T \tilde{x})x_0$ .  $\square$

When the multiplicity of an eigenvalue  $\tilde{\lambda}$  whose eigenvector is being sought is greater than 1, we can see that  $A - \tilde{\lambda}I + (\gamma + \tilde{\lambda})(x_0 x_0^T)$  is singular and may not apply the previous propositions. However, we can still estimate a corresponding eigenvector with an estimated eigenvalue.

**Proposition 3.** *Let  $A$  be an  $N \times N$  symmetric matrix with eigenvalues  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$  and linearly independent corresponding unit eigenvectors  $q_1, \dots, q_N$ , respectively. Let  $\tilde{\lambda}$  be an eigenvalue  $\lambda_{k_0}$ , for some  $1 \leq k_0 \leq N$ , of  $A$  with multiplicity  $m > 1$  so that  $\tilde{\lambda} = \lambda_{k_0} = \dots = \lambda_{k_0+m-1}$ . Let  $\gamma > 0$  be such that  $\gamma \neq -\tilde{\lambda}$ . Let  $x_0 \in S^{N-1}$  be such that*

$$\zeta_0 := \sqrt{\sum_{j=0}^{m-1} |q_{k_0+j}^T x_0|^2} > 0.$$

For  $\lambda$  close enough to  $\tilde{\lambda}$ , but not equal, there exists a unique solution  $x_\lambda$  to

$$(A - \lambda I + (\gamma + \lambda)(x_0 x_0^T))x_\lambda = \gamma x_0.$$

Moreover, there exists  $\eta > 0$  such that

$$\|(A - \tilde{\lambda}I)x_\lambda\| \leq \eta|\tilde{\lambda} - \lambda|$$

and  $\left\{\frac{x_\lambda}{\|x_\lambda\|}\right\}$  has exactly two limit points  $\tilde{x}$  and  $-\tilde{x}$ , where  $\tilde{x}$  satisfies

$$q_i^T \tilde{x} = \begin{cases} \frac{q_i^T x_0}{\zeta_0}, & \text{for } k_0 \leq i \leq k_0 + m - 1, \\ 0, & \text{otherwise.} \end{cases}$$

*Proof.* It suffices to consider a diagonal matrix  $A$  with diagonal entries  $\lambda_1 \leq \dots \leq \lambda_N$ . Let  $\tilde{\lambda}$  be of multiplicity  $m > 1$  and  $\tilde{\lambda} = \lambda_{k_0} = \dots = \lambda_{k_0+m-1}$  for some  $1 \leq k_0 \leq N - m + 1$ . We assume that

$$\zeta_0 := \sqrt{\sum_{j=0}^{m-1} |x_{0, k_0+j}|^2} > 0.$$

Note that if  $\lambda \neq \lambda_j$  for all  $1 \leq j \leq N$ ,  $A - \lambda I$  is nonsingular and  $x_\lambda$  must satisfy

$$(A - \lambda I)x_\lambda = \begin{bmatrix} (\lambda_1 - \lambda)x_{\lambda,1} \\ \vdots \\ (\lambda_N - \lambda)x_{\lambda,N} \end{bmatrix} = (\gamma - (\gamma + \lambda)(x_0^T x_\lambda)) \begin{bmatrix} x_{0,1} \\ \vdots \\ v_{0,N} \end{bmatrix},$$

i.e.,  $(A - \lambda I)x_\lambda = \alpha_\lambda x_0$  for some  $\alpha_\lambda \in \mathbb{R}$ . Since  $x_\lambda = 0$  cannot be a solution, we know that  $\alpha_\lambda$  must be nonzero and

$$x_\lambda = \alpha_\lambda \begin{bmatrix} \frac{1}{\lambda_1 - \lambda} x_{0,1} \\ \vdots \\ \frac{1}{\lambda_N - \lambda} x_{0,N} \end{bmatrix}.$$

Since  $\alpha_\lambda = \gamma - (\gamma + \lambda)(x_0^T x_\lambda)$ , we have

$$\alpha_\lambda = \gamma - \alpha_\lambda \left( \frac{\gamma + \lambda}{\lambda_1 - \lambda} x_{0,1}^2 + \cdots + \frac{\gamma + \lambda}{\lambda_N - \lambda} x_{0,N}^2 \right)$$

implying

$$\alpha_\lambda = \frac{\gamma}{1 + \frac{\gamma + \lambda}{\lambda_1 - \lambda} x_{0,1}^2 + \cdots + \frac{\gamma + \lambda}{\lambda_N - \lambda} x_{0,N}^2} = \frac{\gamma}{\sum_{j=1}^N \frac{(\gamma + \lambda_j)}{\lambda_j - \lambda} x_{0,j}^2}.$$

For  $\lambda$  close enough to  $\tilde{\lambda}$ , not equal to, it is easy to see that  $\alpha_\lambda \neq 0$  exists and is unique, i.e.  $x_\lambda$  exists and is unique. In fact, since

$$\begin{aligned} & \left| \frac{\gamma + \tilde{\lambda}}{\tilde{\lambda} - \lambda} \left| \left( \sum_{j=0}^{m-1} x_{0,k_0+j}^2 \right) - \left| \sum_{j \notin \{k_0, \dots, k_0+m-1\}} \frac{(\gamma + \lambda_j)}{\lambda_j - \lambda} x_{0,j}^2 \right| \right. \right| \\ & \leq \left| \sum_{j=1}^N \frac{(\gamma + \lambda_j)}{\lambda_j - \lambda} x_{0,j}^2 \right| \\ & \leq \left| \frac{\gamma + \tilde{\lambda}}{\tilde{\lambda} - \lambda} \left| \left( \sum_{j=0}^{m-1} x_{0,k_0+j}^2 \right) + \left| \sum_{j \notin \{k_0, \dots, k_0+m-1\}} \frac{(\gamma + \lambda_j)}{\lambda_j - \lambda} x_{0,j}^2 \right| \right. \right|, \end{aligned}$$

we can see that there exists  $0 < \epsilon$  such that  $0 < |\lambda - \tilde{\lambda}| < \epsilon$  implies

$$\frac{1}{2} \left| \frac{\gamma + \tilde{\lambda}}{\tilde{\lambda} - \lambda} \right| \left| \left( \sum_{j=0}^{m-1} x_{0,k_0+j}^2 \right) \right| \leq \left| \sum_{j=1}^N \frac{(\gamma + \lambda_j)}{\lambda_j - \lambda} x_{0,j}^2 \right| \leq \frac{3}{2} \left| \frac{\gamma + \tilde{\lambda}}{\tilde{\lambda} - \lambda} \right| \left| \left( \sum_{j=0}^{m-1} x_{0,k_0+j}^2 \right) \right|,$$

that is,

$$\frac{2\gamma}{3(\gamma + \tilde{\lambda})\zeta_0^2} |\lambda - \tilde{\lambda}| \leq |\alpha_\lambda| \leq \frac{2\gamma}{(\gamma + \tilde{\lambda})\zeta_0^2} |\lambda - \tilde{\lambda}|.$$

Let  $c = \min(|\lambda_{k_0-1} - \tilde{\lambda}|, |\lambda_{k_0+m} - \tilde{\lambda}|)$ . We may assume that  $\epsilon < \min(c, 1)$ . Therefore, for  $0 < |\lambda - \tilde{\lambda}| < \epsilon$ ,  $x_\lambda$  exists and is unique and

$$A - \lambda I + (\gamma + \lambda)(x_0 x_0^T)$$

is nonsingular, as well.

Since  $\|x_\lambda\| = |\alpha_\lambda| \sqrt{\sum_{j=1}^N \left(\frac{1}{\lambda_j - \lambda}\right)^2 x_{0,j}^2} \geq \frac{|\alpha_\lambda|}{|\lambda - \tilde{\lambda}|} \zeta_0$  and  $0 < \epsilon < c$ , we have that

$$\frac{2\gamma}{3(\gamma + \tilde{\lambda})\zeta_0} \leq \|x_\lambda\| \leq \frac{2\gamma}{(\gamma + \tilde{\lambda})\zeta_0^2} \sqrt{\frac{\epsilon^2}{c^2} + \zeta_0^2} \leq \frac{2\gamma\sqrt{1 + \zeta_0^2}}{(\gamma + \tilde{\lambda})\zeta_0^2} \leq \frac{2\sqrt{2}\gamma}{(\gamma + \tilde{\lambda})\zeta_0^2}$$

for  $0 < |\lambda - \tilde{\lambda}| < \epsilon$ . Together with

$$\|(A - \lambda I)x_\lambda\| = \|\alpha_\lambda x_0\| = |\alpha_\lambda| \leq \frac{2\gamma}{(\gamma + \tilde{\lambda})\zeta_0^2} |\lambda - \tilde{\lambda}|,$$

we have that  $0 < |\tilde{\lambda} - \lambda| < \epsilon$  implies

$$\|(A - \tilde{\lambda} I)x_\lambda\| \leq \|(A - \lambda I)x_\lambda\| + \|x_\lambda\| |\tilde{\lambda} - \lambda| \leq \eta |\tilde{\lambda} - \lambda|,$$

where  $\eta = \frac{2(1+\sqrt{2})\gamma}{(\gamma+\tilde{\lambda})\zeta_0^2}$ .

Moreover,

$$y_\lambda := \frac{x_\lambda}{\|x_\lambda\|} = \begin{bmatrix} \frac{\frac{1}{\lambda_1 - \lambda} x_{0,1}}{\sqrt{\sum_{j=1}^N (\frac{1}{\lambda_j - \lambda})^2 x_{0,j}^2}} \\ \vdots \\ \frac{\frac{1}{\lambda_N - \lambda} x_{0,N}}{\sqrt{\sum_{j=1}^N (\frac{1}{\lambda_j - \lambda})^2 x_{0,j}^2}} \end{bmatrix}$$

and if  $j \notin \{k_0, k_0 + 1, \dots, k_0 + m - 1\}$ , then as  $\lambda \rightarrow \tilde{\lambda} = \lambda_{k_0}$ , we have

$$\begin{aligned} y_{\lambda,j} &= \frac{\frac{1}{\lambda_j - \lambda} x_{0,j}}{\sqrt{\sum_{j=1}^N (\frac{1}{\lambda_j - \lambda})^2 x_{0,j}^2}} \\ &= \frac{\frac{|\lambda_{k_0} - \lambda|}{\lambda_j - \lambda} x_{0,j}}{\sqrt{\sum_{j \notin \{k_0, \dots, k_0 + m - 1\}} (\frac{\lambda_{k_0} - \lambda}{\lambda_j - \lambda})^2 x_{0,j}^2 + x_{0,k_0}^2 + \dots + x_{0,k_0 + m - 1}^2}} \\ &\rightarrow 0. \end{aligned}$$

If  $j \in \{k_0, k_0 + 1, \dots, k_0 + m - 1\}$ , then it is easy to see that

$$y_{\lambda,j} \rightarrow \begin{cases} \frac{x_{0,j}}{\zeta_0}, & \text{as } \lambda \uparrow \lambda_{k_0}, \\ -\frac{x_{0,j}}{\zeta_0}, & \text{as } \lambda \downarrow \lambda_{k_0}. \end{cases}$$

Hence,

$$y_\lambda \rightarrow \begin{cases} \tilde{x}, & \text{as } \lambda \uparrow \lambda_{k_0}, \\ -\tilde{x}, & \text{as } \lambda \downarrow \lambda_{k_0}, \end{cases}$$

where  $\tilde{x} = \begin{cases} \frac{x_{0,i}}{\zeta_0}, & \text{for } k_0 \leq i < k_0 + m, \\ 0, & \text{otherwise.} \end{cases}$  □

Surprisingly, the propositions above say that the error in estimating an eigenvector  $\tilde{x}$  is no worse than the error in estimating an eigenvalue  $\tilde{\lambda}$ . That is, if we have a very accurate estimation of an eigenvalue, then we can estimate a corresponding eigenvector within the same accuracy using (25), which is a new observation. For example, there is a recently proposed method to estimate eigenvalues using contour integrals in [1] and similar previous works, therein, where the authors mentioned that if more accurate estimation is needed, then part of their algorithm could be run again with another randomly set of vectors. However, simply solving our one step estimation method (25) with the estimated eigenvalues provides better eigenvalue estimation since an  $O(\epsilon)$  error in eigenvector estimation gives rise to an  $O(\epsilon^2)$  error in eigenvalue estimation (Rayleigh quotient effect), and this error in eigenvalue guarantees the same amount of error  $O(\epsilon^2)$  in eigenvector using the simple one step (25).

### 2.1.5 One Step Estimation For Nonsymmetric Diagonalizable Matrices

We can also observe that the above propositions provide us with a method either to find or to approximate eigenvectors of a nonsymmetric diagonalizable matrix  $A$  within a prescribed error. Since finding an eigenvector of  $A$  corresponding to an eigenvalue  $\lambda$  is equivalent to finding a nonzero vector in the null space of  $A - \lambda I$ , we will provide the following corollary.

**Corollary 1.** *Let  $A$  be an  $N \times N$  nonzero real diagonalizable matrix. Suppose that  $A$  has a nontrivial null space  $\mathcal{N}(A)$ . Let  $\gamma > 0$ .*

1. *If  $\mathcal{N}(A)$  has dimension 1, then choosing a unit vector  $x_0$  uniformly at random, we have that with probability 1,*

$$(A + \gamma x_0 x_0^T)x = \gamma x_0$$

*has a unique solution  $\tilde{x}$  that spans  $\mathcal{N}(A)$ .*

2. *If  $\mathcal{N}(A)$  has dimension greater than 1, then for small enough  $\lambda > 0$ , choosing a unit vector  $x_0$  uniformly at random, we have that with probability 1,*

$$(A - \lambda I + (\gamma + \lambda)x_0 x_0^T)x = \gamma x_0 \tag{27}$$

*has a unique solution  $x_\lambda$  satisfying  $\|Ax_\lambda\| \leq \eta\lambda$  with  $\eta$  depending on  $x_0$ .*

*Proof.* The first part can be proven in the same way as we proved Proposition 1 with a choice of  $x_0$  satisfying  $q^T x_0 \neq 0$  where  $q$  is a unit vector spanning  $\mathcal{N}(A)$ . Since the second part can be proven with a minor modification in the proof of Proposition 3, we will only provide a sketch of this proof.

Note that  $A - \lambda I$  for small enough  $\lambda > 0$  is invertible and that if we choose a unit vector  $x_0$  uniformly at random, then  $x_0$  is not orthogonal to  $\mathcal{N}(A)$  with probability 1. In fact, since  $\mathbb{R}^N = \text{eigenspace}(0) \oplus \bigoplus_{\lambda_j \neq 0} \text{eigenspace}(\lambda_j)$ , we can represent  $x_0$  as  $q_0 + r_0$ , where  $0 \neq q_0 \in \text{eigenspace}(0) = \mathcal{N}(A)$  and  $r_0 \in \bigoplus_{\lambda_j \neq 0} \text{eigenspace}(\lambda_j)$ .

With such an  $x_0$ , we need to see if we can solve  $(A - \lambda I)x = (\gamma x_0 - (\gamma + \lambda)x_0 x_0^T x) = \alpha_\lambda x_0$ , with

$$\alpha_\lambda = \gamma - (\gamma + \lambda)x_0^T x.$$

Noting that if there exists such a solution  $x_\lambda$ , we have

$$x_\lambda = \alpha_\lambda (A - \lambda I)^{-1} x_0$$

and  $\alpha_\lambda$  satisfies

$$\alpha_\lambda = \gamma - \alpha_\lambda (\gamma + \lambda) x_0^T (A - \lambda I)^{-1} x_0,$$

i.e.,

$$\alpha_\lambda (1 + (\gamma + \lambda) x_0^T (A - \lambda I)^{-1} x_0) = \gamma \Leftrightarrow \alpha_\lambda \left( 1 + \left( \frac{\gamma}{\lambda} + 1 \right) x_0^T \left( \frac{1}{\lambda} A - I \right)^{-1} x_0 \right) = \gamma.$$

It is not difficult to see that

$$\left\| \left( \frac{1}{\lambda} A - I \right)^{-1} r_0 \right\| \rightarrow 0 \quad \text{as} \quad \lambda \downarrow 0.$$

Hence, there exists  $\lambda_0 < \min(\|A\|, \frac{\gamma |x_0^T q_0|}{2(2 - |x_0^T q_0|)})$  such that  $0 < \lambda < \lambda_0$  implies

- $\left\| \left( \frac{1}{\lambda} A - I \right)^{-1} r_0 \right\| < \frac{|x_0^T q_0|}{2}$ ,
- $\frac{\gamma}{\lambda} + 1 > \frac{2}{|x_0^T q_0|}$ .

This implies that for  $0 < \lambda < \lambda_0$ , we have

$$\frac{1}{2}|x_0^T q_0| < \left| x_0^T \left( \frac{1}{\lambda} A - I \right)^{-1} x_0 \right| = \left| -x_0^T q_0 + x_0^T \left( \frac{1}{\lambda} A - I \right)^{-1} r_0 \right| < \frac{3}{2}|x_0^T q_0|.$$

It is now easy to see that for  $0 < \lambda < \lambda_0$ ,  $\alpha_\lambda = \frac{\gamma}{(1+(\gamma+\lambda)x_0^T(A-\lambda I)^{-1}x_0)}$  exists uniquely and satisfies

$$|\alpha_\lambda| < \frac{\gamma}{\frac{1}{2}\left(\frac{\gamma}{\lambda} + 1\right)|x_0^T q_0| - 1} = \left( \frac{2\gamma}{\gamma|x_0^T q_0| - \lambda(2 - |x_0^T q_0|)} \right) \lambda < \frac{4\lambda}{|x_0^T q_0|}$$

and

$$|\alpha_\lambda| > \frac{\gamma}{\frac{3}{2}\left(\frac{\gamma}{\lambda} + 1\right)|x_0^T q_0| + 1} = \left( \frac{2\gamma}{3\gamma|x_0^T q_0| + \lambda(2 + 3|x_0^T q_0|)} \right) \lambda > \frac{4\lambda}{11|x_0^T q_0|}$$

That is, for  $0 < \lambda < \lambda_0$ , there exists a unique solution

$$x_\lambda = \alpha_\lambda(A - \lambda I)^{-1}x_0 = \frac{\gamma(A - \lambda I)^{-1}x_0}{(1 + (\gamma + \lambda)x_0^T(A - \lambda I)^{-1}x_0)}$$

with

$$\frac{4}{11|x_0^T q_0|} \|\lambda(A - \lambda I)^{-1}x_0\| < \|x_\lambda\| < \frac{4}{|x_0^T q_0|} \|\lambda(A - \lambda I)^{-1}x_0\|$$

Note that for  $0 < \lambda < \lambda_0$ , we have

$$\frac{\|q_0\|}{2} \leq \|\lambda(A - \lambda I)^{-1}x_0\| = \left\| \left( \frac{1}{\lambda} A - I \right)^{-1} x_0 \right\| = \left\| -q_0 + \left( \frac{1}{\lambda} A - I \right)^{-1} r_0 \right\| \leq \frac{3\|q_0\|}{2},$$

implying that for  $0 < \lambda < \lambda_0$ ,

$$\frac{2\|q_0\|}{11|x_0^T q_0|} < \|x_\lambda\| < \frac{6\|q_0\|}{|x_0^T q_0|}.$$

Finally, we note that  $0 < \lambda < \lambda_0$  implies

$$\begin{aligned} \|Ax_\lambda\| &= \left\| \frac{\gamma A(A - \lambda I)^{-1}x_0}{(1 + (\gamma + \lambda)x_0^T(A - \lambda I)^{-1}x_0)} \right\| \\ &= \left\| \alpha_\lambda x_0 + \lambda \alpha_\lambda (A - \lambda I)^{-1}x_0 \right\| \\ &\leq |\alpha_\lambda| (1 + \|\lambda(A - \lambda I)^{-1}x_0\|) \\ &= \eta \lambda, \end{aligned}$$

where  $\eta = \frac{4+6\|q_0\|}{|x_0^T q_0|}$ . □

### 3 Applications

So far, we have analyzed our proposed model (3) when the given matrix  $A$  is real symmetric. As was mentioned earlier, the same analysis can be carried over to complex hermitian matrices. When  $A$  is real and nonsymmetric, our proposed model (3) with  $B = A^T A$  will present a singular value decomposition of  $A$ . Since  $B$  is symmetric, the previous discussion applies without alteration. For this reason, we will not give a detailed presentation for singular value estimation, but rather provide other applications such as generalized eigenvalue problems and eigenfunctions of self-adjoint elliptic operators.

### 3.1 Generalized Eigenvalue Problem : $Ax = \lambda Bx$

Given that  $A, B \in \mathbf{M}_N(\mathbb{R})$  are symmetric matrices, and  $B$  is positive definite with the smallest eigenvalue  $\mu_{(B,1)} > 0$ , and  $A$  has the smallest eigenvalue  $\mu_{(A,1)}$ , we would like to consider the generalized eigenvalue problem:

$$Ax = \lambda Bx. \quad (28)$$

We will denote by  $\mu_{(A,j)}, \mu_{(B,j)}$  the  $j^{\text{th}}$  smallest eigenvalues of  $A$  and  $B$ , respectively. In addition, it is enough to consider that  $\mu_{(B,1)} = 1$  and that either  $-\mu_{(A,N)} \leq \mu_{(A,1)} < 0$  or  $\mu_{(A,1)} \geq 0$  holds because (28) is equivalent to

$$\pm Ax = (\pm \lambda \cdot \mu_{(B,1)}) \left( \frac{1}{\mu_{(B,1)}} B \right) x.$$

With  $\gamma > \max(0, -\mu_{(A,1)})$ , we know that  $A + \gamma B$  is positive definite and we can solve (28) by minimizing

$$F_{A,B}(x) = \frac{1}{2} \langle Ax, x \rangle + \frac{\gamma}{2} \langle Bx, x \rangle - \gamma \sqrt{\langle Bx, x \rangle}, \quad (29)$$

since we have

$$\nabla F_{A,B}(x) = Ax + \gamma \left( 1 - \frac{1}{\sqrt{\langle Bx, x \rangle}} \right) Bx = 0$$

if and only if

$$Ax = \gamma \left( \frac{1}{\sqrt{\langle Bx, x \rangle}} - 1 \right) Bx.$$

Therefore, if  $x_* \neq 0$  is a critical point of  $F_{A,B}$ , then  $(x_*, \lambda_*)$  is a solution pair of (28), i.e.,  $Ax_* = \lambda_* Bx_*$  with

$$\lambda_* = \gamma \left( \frac{1}{\sqrt{\langle Bx_*, x_* \rangle}} - 1 \right).$$

Moreover,  $(y_*, \mu_*)$  with  $y_* = \sqrt{B}x_*$  and  $\mu_* = \gamma \left( \frac{1}{\sqrt{\langle Bx_*, x_* \rangle}} - 1 \right)$  is an eigenpair of the matrix  $C = \sqrt{B}^{-1} A \sqrt{B}^{-1}$ .

In this section, we set  $r_1, \dots, r_N$  to be linearly independent unit eigenvectors of  $C$  corresponding to the eigenvalues  $\lambda_1 \leq \dots \leq \lambda_N$ . We also let  $q_j = \sqrt{B}^{-1} r_j$  for  $1 \leq j \leq N$ . Then, it is easy to see that  $\{q_1, \dots, q_N\}$  is an orthonormal basis for  $\mathbb{R}^N$  with respect to the inner product  $\langle x, y \rangle_B := x^T B y$ .

We can also apply the Gradient Descent method to (29) to solve (28) just as we did in Theorem 3.

**Proposition 4.** *Let  $A, B$  be as mentioned above. If we choose*

$$0 < \alpha < \frac{1}{(\mu_{(A,N)} + \gamma)(\mu_{(B,N)})^3},$$

*then with a randomly chosen  $x_0 \neq 0$ ,*

$$x_{k+1} = x_k - \alpha \nabla F_{A,B}(x_k), \quad k = 0, 1, 2, \dots, \quad (30)$$

*produce a sequence  $\{x_k\}$  converging to  $x_*$  with  $\langle Bx_*, x_* \rangle = \frac{\gamma}{\gamma + \lambda_*}$ , where  $(x_*, \lambda_*)$  is a solution pair of (28) with  $\lambda_* = \lambda_1 = \min\{\lambda \in \mathbb{R} : A - \lambda B \text{ is singular}\}$ .*

*Proof.* This proof mimics that of Theorem 3 with minor differences in detail. Let  $\{x_k\}$  be the sequence generated by (30). First of all, we note that for any  $x_k \in \mathbb{R}^N$ ,

$$\left| \frac{\langle Ax_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \right| \leq \mu_{(A,N)} \mu_{(B,N)} \quad \text{and} \quad 1 = \mu_{(B,1)} \leq \frac{\langle Bx_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \leq \mu_{(B,N)},$$

which implies that if  $0 < \alpha < \frac{1}{(\mu_{(A,N)} + \gamma)(\mu_{(B,N)})^3}$ , then with  $0 < \epsilon := 1 - \frac{1}{(\mu_{(B,N)})^2}$ , we have

$$1 - \alpha \left[ \frac{\langle Ax_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} + \gamma \frac{\langle Bx_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \right] \geq 1 - \alpha(\mu_{(A,N)} + \gamma)\mu_{(B,N)} > \epsilon.$$

We can see that for  $k = 1, 2, \dots$ ,

$$\begin{aligned} \langle Bx_{k+1}, x_{k+1} \rangle &\geq \frac{\langle Bx_{k+1}, x_k \rangle^2}{\langle Bx_k, x_k \rangle} = \left( \frac{\langle x_k - \alpha \nabla F_A(x_k), Bx_k \rangle}{\sqrt{\langle Bx_k, x_k \rangle}} \right)^2 \\ &= \left( \sqrt{\langle Bx_k, x_k \rangle} - \alpha \left[ \frac{\langle Ax_k, Bx_k \rangle}{\sqrt{\langle Bx_k, x_k \rangle}} + \gamma \frac{\langle Bx_k, Bx_k \rangle}{\sqrt{\langle Bx_k, x_k \rangle}} \right] + \alpha \gamma \frac{\langle Bx_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \right)^2 \\ &= \left[ \sqrt{\langle Bx_k, x_k \rangle} \left( 1 - \alpha \left[ \frac{\langle Ax_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} + \gamma \frac{\langle Bx_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \right] \right) + \alpha \gamma \frac{\langle Bx_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \right]^2 \\ &\geq (\sqrt{\langle Bx_k, x_k \rangle} \epsilon + \alpha \gamma)^2 = \langle Bx_k, x_k \rangle \epsilon^2 + 2\sqrt{\langle Bx_k, x_k \rangle} \alpha \gamma \epsilon + (\alpha \gamma)^2. \end{aligned} \quad (31)$$

This proves that  $\langle Bx_k, x_k \rangle \geq (\alpha \gamma)^2$  for  $k \geq 1$ , which enables us to improve the inequality in (31) as follows: for  $k \geq 1$ ,

$$\begin{aligned} \langle Bx_{k+1}, x_{k+1} \rangle &\geq (\sqrt{\langle Bx_k, x_k \rangle} \epsilon + \alpha \gamma)^2 = \langle Bx_k, x_k \rangle \epsilon^2 + 2\sqrt{\langle Bx_k, x_k \rangle} \alpha \gamma \epsilon + (\alpha \gamma)^2 \\ &\geq (\alpha \gamma)^2 (1 + 2\epsilon) + \langle Bx_k, x_k \rangle \epsilon^2, \end{aligned}$$

resulting in

$$\langle Bx_{k+1}, x_{k+1} \rangle \geq (\alpha \gamma)^2 (1 + 2\epsilon) \sum_{l=0}^{k-1} \epsilon^{2l} + (\alpha \gamma)^2 \epsilon^{2k}$$

Hence,

$$\liminf_{k \rightarrow \infty} \langle Bx_k, x_k \rangle \geq (\alpha \gamma)^2 \frac{1 + 2\epsilon}{1 - \epsilon^2} > \frac{(\alpha \gamma)^2}{1 - \epsilon} = (\alpha \gamma \mu_{(B,N)})^2$$

and there exists  $K \in \mathbb{N}$  such that  $k \geq K$  implies

$$\langle Bx_k, x_k \rangle > (\alpha \gamma \mu_{(B,N)})^2.$$

Moreover, we can estimate  $\langle B(x_k + t(x_{k+1} - x_k)), (x_k + t(x_{k+1} - x_k)) \rangle$  for  $t \in (0, 1)$  as follows: for  $k > 1, 2, \dots$ , and  $0 < t < 1$ ,

$$\begin{aligned} &\langle B(x_k + t(x_{k+1} - x_k)), (x_k + t(x_{k+1} - x_k)) \rangle = \|\sqrt{B}(x_k - t\alpha \nabla F_A(x_k))\|^2 \\ &\geq \left\langle \sqrt{B}(x_k - t\alpha \nabla F_{A,B}(x_k)), \frac{\sqrt{B}x_k}{\|\sqrt{B}x_k\|} \right\rangle^2 = \left( \frac{\langle x_k - t\alpha \nabla F_A(x_k), Bx_k \rangle}{\sqrt{\langle Bx_k, x_k \rangle}} \right)^2 \\ &= \left[ \sqrt{\langle Bx_k, x_k \rangle} \left( 1 - t\alpha \left[ \frac{\langle Ax_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} + \gamma \frac{\langle Bx_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \right] \right) + t\alpha \gamma \frac{\langle Bx_k, Bx_k \rangle}{\langle Bx_k, x_k \rangle} \right]^2 \\ &\geq (\sqrt{\langle Bx_k, x_k \rangle} (1 - t\epsilon_1) + t\alpha \gamma)^2 = (\sqrt{\langle Bx_k, x_k \rangle} + t(\alpha \gamma - \sqrt{\langle Bx_k, x_k \rangle} \epsilon_1))^2 \\ &\geq \begin{cases} \langle Bx_k, x_k \rangle, & \text{if } \alpha \gamma \geq \sqrt{\langle Bx_k, x_k \rangle} \epsilon_1, \\ (\sqrt{\langle Bx_k, x_k \rangle} (1 - \epsilon_1) + \alpha \gamma)^2, & \text{if } \alpha \gamma < \sqrt{\langle Bx_k, x_k \rangle} \epsilon_1. \end{cases} \end{aligned}$$



where  $\epsilon_1 = 1 - \epsilon$ . Noting that  $\alpha\gamma < \sqrt{\langle Bx_k, x_k \rangle} \epsilon_1$  implies

$$\sqrt{\langle Bx_k, x_k \rangle} (1 - \epsilon_1) + \alpha\gamma > \frac{\alpha\gamma(1 - \epsilon_1)}{\epsilon_1} + \alpha\gamma = \frac{\alpha\gamma}{1 - \epsilon},$$

we have that for  $k \geq K$ ,

$$\min_{t \in [0,1]} \langle B(x_k + t(x_{k+1} - x_k)), (x_k + t(x_{k+1} - x_k)) \rangle > (\alpha\gamma\mu_{(B,N)})^2.$$

Therefore, we have

$$\{x_k + t(x_{k+1} - x_k) : t \in [0, 1], k \geq K\} \subset \{x \in \mathbb{R}^N : \langle Bx, x \rangle \geq (\alpha\gamma\mu_{(B,N)})^2\}.$$

Since  $\sqrt{B} \left( \frac{\gamma}{\sqrt{\langle Bx, x \rangle}} \left[ I - \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right) \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right)^T \right] \right) \sqrt{B}$  is positive semidefinite and

$$\nabla^2 F_{A,B}(x) = A + \gamma B - \sqrt{B} \left( \frac{\gamma}{\sqrt{\langle Bx, x \rangle}} \left[ I - \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right) \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right)^T \right] \right) \sqrt{B},$$

we can see that

$$\|\nabla^2 F_{A,B}(x)\| \leq \max \left\{ \|A + \gamma B\|, \left\| \sqrt{B} \left( \frac{\gamma}{\sqrt{\langle Bx, x \rangle}} \left[ I - \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right) \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right)^T \right] \right) \sqrt{B} \right\| \right\}.$$

We also note that

$$\|A + \gamma B\| \leq \mu_{(A,N)} + \gamma\mu_{(B,N)} < \frac{1}{\alpha},$$

and that for  $x \in \{z \in \mathbb{R}^N : \langle Bz, z \rangle \geq (\alpha\gamma\mu_{(B,N)})^2\}$

$$\left\| \sqrt{B} \left( \frac{\gamma}{\sqrt{\langle Bx, x \rangle}} \left[ I - \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right) \left( \frac{\sqrt{B}x}{\|\sqrt{B}x\|} \right)^T \right] \right) \sqrt{B} \right\| \leq \frac{\mu_{(B,N)}}{\alpha\mu_{(B,N)}} = \frac{1}{\alpha}.$$

Therefore, for each  $k \geq K$ ,

$$\sup_{t \in [0,1]} \|\nabla^2 F_{A,B}(x_k + t(x_{k+1} - x_k))\| \leq \frac{1}{\alpha}$$

and we finally have that for  $k \geq K$ ,

$$\begin{aligned} F_{A,B}(x_{k+1}) &\leq F_{A,B}(x_k) + \langle \nabla F_{A,B}(x_k), x_{k+1} - x_k \rangle + \frac{1}{2\alpha} \|x_{k+1} - x_k\|^2 \\ &= F_{A,B}(x_k) - \frac{\alpha}{2} \|\nabla F_{A,B}(x_k)\|^2, \end{aligned}$$

which implies

$$\sum_{k=K}^{\infty} \|\nabla F_{A,B}(x_k)\|^2 \leq \frac{2}{\alpha} (F_{A,B}(x_K) - \min_{x \in \mathbb{R}^N} (F_{A,B}(x))).$$

We omit the rest of the proof because we can now follow the proof of Theorem 3 for the convergence of the sequence  $\{x_k\}_{k>K}$  to  $x_*$  satisfying

$$\sqrt{\langle Bx_*, x_* \rangle} = \|\sqrt{B}x_*\| = \frac{\gamma}{\gamma + \lambda_*},$$

where  $(x_*, \lambda_*)$  is a solution pair of (28).

Let  $r_1, \dots, r_N$  be linearly independent unit eigenvectors of  $C = \sqrt{B}^{-1} A \sqrt{B}^{-1}$  corresponding to the eigenvalues  $\lambda_1 \leq \dots \leq \lambda_N$ . If we define  $q_j$ ,  $1 \leq j \leq N$  by  $q_j = \sqrt{B}^{-1} r_j$  and represent  $x_k$ ,  $k \geq 1$ , as

$$x_k = \mu_{k,1} q_1 + \dots + \mu_{k,N} q_N,$$

then (30) becomes

$$\sum_{j=1}^N \mu_{k+1,j} q_j = \sum_{j=1}^N \mu_{k,j} q_j - \alpha \left[ \sum_{j=1}^N \mu_{k,j} A q_j + \gamma \left( 1 - \frac{1}{\|\sqrt{B} x_k\|} \right) \sum_{j=1}^N \mu_{k,j} B q_j \right].$$

Noting that  $\{q_1, \dots, q_N\}$  is an orthonormal basis for  $\mathbb{R}^N$  with respect to the inner product  $\langle x, x \rangle_B := x^T B x$ , by taking the inner product  $\langle q_j, \sum_{j=1}^N \mu_{k+1,j} q_j \rangle$  for each  $j = 1, 2, \dots, N$ , we have

$$\mu_{k+1,j} = \mu_{k,j} \left( 1 - \alpha \left( \lambda_j + \gamma - \frac{\gamma}{\|\sqrt{B} x_k\|} \right) \right) = \mu_{0,j} \prod_{m=0}^k \left( 1 - \alpha \left( \lambda_j + \gamma - \frac{\gamma}{\|\sqrt{B} x_m\|} \right) \right)$$

since  $\langle q_j, A q_l \rangle = \langle r_j, C r_l \rangle = \lambda_j \delta_{jl}$ . The rest of the proof is the same as that of Theorem 3. Hence, we have that

$$\lambda_* = \lambda_1 = \min\{\lambda \in \mathbb{R} : A - \lambda B \text{ is singular}\}.$$

□

If it is easy to compute  $B^{-1}$ , then we can do a little better than Proposition 4. Due to the equivalence between  $Ax = \lambda Bx$  and  $B^{-1}Ax = \lambda x$ , all we need to do is to find eigenvectors and their corresponding eigenvalues of  $E = B^{-1}A$ . Since  $E$  may not be symmetric, we cannot use Theorem 3. However, realizing that  $F_{A,B}$  can be rewritten using a different inner product and proposing a gradient descent method with respect to the new inner product, we can find an eigenvector of  $E$  corresponding to the smallest eigenvalue. We will also see that this method allows for a bigger stepsize.

**Proposition 5.** *With  $F_{A,B}$  defined as in (29) and  $\gamma > \max(0, -\mu_{(A,1)})$ , a sequence  $\{x_k\}$  generated by*

$$x_{k+1} = x_k - \alpha B^{-1} \nabla F_{A,B}(x_k), \quad k = 0, 1, 2, \dots, \quad (32)$$

*with a randomly chosen  $x_0 \neq 0$  and  $0 < \alpha < \frac{1}{\mu_{(A,N)} + \gamma}$ , converges to  $x_*$ , where  $(x_*, \lambda_*)$  is a solution pair of (28) satisfying*

$$\sqrt{\langle Bx_*, x_* \rangle} = \|\sqrt{B}x_*\| = \frac{\gamma}{\gamma + \lambda_*}$$

and

$$\lambda_* = \min\{\lambda : A - \lambda B \text{ is singular}\}.$$

*In fact, (32) is the gradient descent of  $F_{A,B}$  with respect to the inner product  $\langle x, y \rangle_B := x^T B y$ .*

*Proof.* First of all, as we mentioned, (28) is equivalent to

$$B^{-1}Ax = \lambda x,$$

and to

$$Cy = \lambda y, \quad \text{with } C = \sqrt{B}^{-1} A \sqrt{B}^{-1}, \quad y = \sqrt{B}x.$$

Therefore,  $(x, \lambda)$  is a solution pair of (28) if and only if  $x$  is an eigenvector of  $B^{-1}A$  corresponding to the eigenvalue  $\lambda$  if and only if  $\sqrt{B}x$  is an eigenvector of  $C$  corresponding to the eigenvalue  $\lambda$ . Note that

$$\lambda_1 = \min_{\|y\|=1} \langle y, Cy \rangle = \frac{\langle x, Ax \rangle}{\|x\|^2} \|x\|^2 \geq \mu_{(A,1)} \|x\|^2 \geq \begin{cases} \mu_{(A,1)}, & \text{if } \mu_{(A,1)} < 0, \\ 0, & \text{if } \mu_{(A,1)} \geq 0, \end{cases}$$

where  $x = \sqrt{B}^{-1}y$  and  $\|x\| = \|\sqrt{B}^{-1}y\| \leq \frac{1}{\sqrt{\mu_{(B,1)}}} = 1$ . Since  $\gamma > \max(0, -\mu_{(A,1)})$ , if we set  $G_C(y)$  to be

$$G(y) = \frac{1}{2} \langle Cy, y \rangle + \frac{\gamma}{2} \|y\|^2 - \gamma \|y\|,$$

and apply Theorem 3 to  $G$ , then since  $\lambda_N$ , the largest eigenvalue of  $C$ , is at most  $\mu_{(A,N)} > 0$ , a sequence  $\{y_k\}$  generated by

$$y_{k+1} = y_k - \alpha \nabla G(y_k), \quad k = 0, 1, 2, \dots,$$

with a randomly chosen  $y_0 \neq 0$ , and any  $0 < \alpha < \frac{1}{\lambda_N + \gamma} \leq \frac{1}{\mu_{(A,N)} + \gamma}$ , converges to  $y_*$ , an eigenvector of  $C$  corresponding to the smallest eigenvalue  $\lambda_* = \lambda_1$  with norm  $\|y_*\| = \frac{\gamma}{\gamma + \lambda_1}$ .

On the other hand, with  $y_k = \sqrt{B}x_k$ ,  $k = 0, 1, \dots$ , we have

$$\begin{aligned} \nabla G_C(y_k) &= Cy_k + \gamma \left(1 - \frac{1}{\|y_k\|}\right) y_k = (\sqrt{B})^{-1} Ax_k + \gamma \left(1 - \frac{1}{\sqrt{\langle Bx_k, x_k \rangle}}\right) \sqrt{B}x_k \\ &= (\sqrt{B})^{-1} \nabla F_{A,B}(x_k), \end{aligned}$$

which implies that

$$y_{k+1} = y_k - \alpha \nabla G_C(y_k)$$

is equivalent to (32), i.e.,  $x_{k+1} = x_k - \alpha B^{-1} \nabla F_{A,B}(x_k)$ .

Hence, the sequence generated by (32), with a randomly chosen  $x_0 \neq 0$ , converges to  $x_*$ , where  $(x_*, \lambda_*)$  is a solution pair of (28) satisfying  $\sqrt{\langle Bx_*, x_* \rangle} = \|y_*\| = \frac{\gamma}{\gamma + \lambda_*}$  and

$$\lambda_* = \lambda_1 = \min\{\lambda : A - \lambda B \text{ is singular}\}.$$

In addition, with respect to the inner product  $\langle x, y \rangle_B := x^T B y$ , we note that

$$B^{-1} \nabla F_{A,B}(x) = B^{-1} Ax + \gamma \left(1 - \frac{1}{\sqrt{\langle Bx, x \rangle}}\right) x = B^{-1} Ax + \gamma \left(1 - \frac{1}{\sqrt{\langle x, x \rangle_B}}\right) x,$$

which is the gradient of  $F_{A,B}$  with respect to the inner product  $\langle \cdot, \cdot \rangle_B$  because

$$F_{A,B}(x) = \frac{1}{2} \langle x, B^{-1} Ax \rangle_B + \frac{\gamma}{2} \|x\|_B^2 - \gamma \|x\|_B,$$

where  $\|x\|_B^2 = \langle x, x \rangle_B$ . Note that  $B^{-1}A$  is self-adjoint with respect to  $\langle \cdot, \cdot \rangle_B$ . Therefore, (32) is the gradient descent of  $F_{A,B}$  with respect to  $\langle \cdot, \cdot \rangle_B$ .  $\square$

In general, with a nonsymmetric square matrix  $E$ , minimizing the functional  $F_E$  in any of the forms in (3) does not guarantee to find an eigenvector of  $E$ . However, Proposition 4 and Proposition 5 make it possible in certain cases when  $E$  decomposes into  $E = B^{-1}A$  with a symmetric matrix  $A$  and a symmetric positive definite matrix  $B$ . If we know such matrices  $A, B$ , then Proposition 4 applies to find a global minimizer of  $F_E$ . Moreover, if it is easy to compute  $B^{-1}$ , then Proposition 5 applies to find a global minimizer of  $F_E$  with a better stepsize. We can also find subsequent eigenvectors in the same way as presented in Theorem 2.

**Proposition 6.** Let  $(\mathbf{x}_1, \lambda_1), \dots, (\mathbf{x}_m, \lambda_m)$  be solution pairs of (28) where  $\lambda_1 \leq \dots \leq \lambda_m$  are the first  $m \geq 1$  smallest ones in  $\{\lambda \in \mathbb{R} : A - \lambda B \text{ is singular}\}$ .

We consider the following problem

$$\min_{x \in \mathbb{R}^N} F_{A,B}(x) \quad \text{subject to} \quad \langle x, \mathbf{x}_i \rangle_B = 0, \quad i = 1, 2, \dots, m. \quad (33)$$

Then, any local minimizer  $x_*$  of (7) is a global minimizer with  $\sqrt{\langle Bx_*, x_* \rangle} = \frac{\gamma}{\gamma + \lambda_*}$  and

$$\lambda_* = \min(\{\lambda \in \mathbb{R} : A - \lambda B \text{ is singular}\} \setminus \{\lambda_1, \dots, \lambda_m\}).$$

*Proof.* We omit the proof because the proof is not so much different from that of Theorem 2. □

As was done with the real symmetric case, we can also consider a similar nonsingular system whose unique solution solves (28). Please note that the following proposition is, in fact, the same as Corollary 1.

**Proposition 7.** Let  $\tilde{\lambda}$  be a generalized eigenvalue of the pencil  $A - \lambda B$  of multiplicity 1. Let  $\gamma > 0$ . Let  $x_0$  be chosen uniformly at random from  $S^{N-1}$ . Then, with probability 1, there is a unique solution  $\tilde{x}$  to

$$(A - \tilde{\lambda}B + \gamma x_0 x_0^T)x = \gamma x_0$$

so that  $(\tilde{x}, \tilde{\lambda})$  is a solution to (28).

When a generalized eigenvalue has multiplicity more than 1, we can also apply Corollary 1 to find approximate eigenvalues.

## 3.2 Eigenvalue problems on infinite dimensional spaces

It is interesting to see that the same framework as (3) applies to eigenvalue problems on infinite dimensional spaces such as the Sturm-Liouville eigenvalue problem, the eigenvalue problem of self-adjoint elliptic operators, etc. We will present one such application, that is, finding an eigenfunction of a self-adjoint uniformly elliptic operator corresponding to the smallest eigenvalue.

### 3.2.1 Symmetric Uniformly Elliptic Operators

Let  $\Omega$  be a bounded open subset of  $\mathbb{R}^d$  with Lipschitz boundary  $\partial\Omega$ . We will denote by  $L$  a symmetric uniformly elliptic operator defined by

$$Lu = - \sum_{i,j=1}^d \partial_i(a_{i,j} \partial_j u) + cu, \quad (34)$$

where  $a_{i,j}, c \in L^\infty(\Omega)$ ,  $i, j = 1, \dots, d$ , are such that  $a_{i,j}(x) = a_{j,i}(x)$  a.e., and  $c(x) \geq 0$  a.e., and there exists  $0 < \alpha \leq \beta < \infty$  such that for a.e.  $x \in \Omega$  and for  $\xi \in \mathbb{R}^d$ ,

$$\alpha |\xi|^2 \leq \sum_{i,j=1}^d a_{i,j}(x) \xi_i \xi_j \leq \beta |\xi|^2.$$

The problem that we are interested in is to find  $\varphi \in H_0^1(\Omega)$  that solves

$$\begin{cases} L\varphi = \lambda\varphi, & \text{in } \Omega, \\ \varphi = 0, & \text{on } \partial\Omega. \end{cases} \quad (35)$$

It is known that the eigenvalues  $\lambda_1, \lambda_2, \dots$  of  $L$  are nonnegative and we may have them in an increasing order, that is,

$$0 < \lambda_1 < \lambda_2 \leq \dots \leq \lambda_3 \leq \dots.$$

Therefore, a natural question to ask is to find the smallest eigenvalue  $\lambda_1$  of  $L$  and its corresponding eigenfunction. It will be apparent that the same analysis applies with the Neumann boundary condition. We would like to stress that a fundamental problem that our method can solve is to find eigenvalues and eigenfunctions of the Laplace operator  $-\Delta$  on any domain  $\Omega$  with Lipschitz boundary.

**Definition 1.**  $\varphi \in H_0^1(\Omega)$  is a weak solution of (35) if for any  $\psi \in C_0^\infty(\Omega)$ ,

$$\int_{\Omega} \sum_{i,j=1}^d a_{i,j}(x) \partial_j \varphi(x) \partial_i \psi(x) dx + \int_{\Omega} (c(x) - \lambda) \varphi(x) \psi(x) dx = 0.$$

Then, it is natural from the discussion in the previous sections that we want to define a functional  $F_L$  with  $\gamma > 0$  by

$$\begin{aligned} F_L(u) &= \frac{1}{2} \int_{\Omega} \sum_{i,j=1}^d a_{i,j}(x) \partial_j u(x) \partial_i u(x) dx + \frac{1}{2} \int_{\Omega} c(x) |u(x)|^2 dx \\ &\quad + \frac{\gamma}{2} \int_{\Omega} |u(x)|^2 dx - \gamma \left( \int_{\Omega} |u(x)|^2 dx \right)^{\frac{1}{2}} \end{aligned} \quad (36)$$

and solve the following minimization problem

$$\min_{u \in H_0^1(\Omega)} F_L(u) \quad (37)$$

and investigation the relation between (35) and (37). The existence of a minimizer of the problem (37) is obvious by the standard method using the compact embedding theorem by Rellich-Kondrachov. Moreover, we can observe the same characteristics of (36) as those of (3). For theorems and lemmas that follow, we will omit their proofs because they are essentially the same as what we presented in the finite dimensional case.

**Lemma 2.** *The set of nonzero critical points of  $F_L$  is*

$$\left\{ \varphi \in H_0^1(\Omega) : L\varphi = \lambda\varphi \text{ in } \Omega \text{ and } \|\varphi\|_{L^2(\Omega)} = \frac{\gamma}{\gamma + \lambda} \right\}.$$

and

$$\min_{u \in H_0^1(\Omega)} F_L(u) = -\frac{\gamma^2}{2(\gamma + \lambda_*)},$$

where  $\lambda_* \geq 0$  is the smallest eigenvalue of  $L$ .

**Theorem 6.** *Any local minimizer of (36) is a global minimizer.*

### 3.2.2 Corresponding parabolic PDEs

Corresponding to the uniformly elliptic PDEs of the form (34), we will consider the following parabolic PDE: for  $T \in (0, \infty)$ , we solve

$$\begin{cases} \frac{\partial u}{\partial t} &= -Lu - \gamma \left( 1 - \frac{1}{\|u(t)\|_2} \right) u \text{ in } \Omega_T := \Omega \times (0, T], \\ u &= 0 \text{ on } \partial\Omega_T := \partial\Omega \times [0, T], \\ u(0) &= u_0 \neq 0 \text{ in } H_0^1(\Omega), \end{cases} \quad (38)$$

where  $\|u(t)\|_2 = (\int_{\Omega} |u(x,t)|^2 dx)^{\frac{1}{2}}$ . Due to the condition  $c \geq 0$  in  $\Omega$ ,  $L$  is positive definite. In the context, we will use  $\langle \cdot, \cdot \rangle$  for both inner products in  $L^2(\Omega)$  and in  $\mathbb{R}^N$  for  $N \in \mathbb{N}$ . Note that the partial differential equation in (38) is the formal gradient flow of a functional  $F_L$  in (36), which is a difference of convex functionals, i.e.,

$$\frac{\partial u}{\partial t} = -\nabla F_L(u).$$

**Theorem 7.** *The problem (38) has a unique weak solution  $u$  for any  $T \in (0, \infty)$ . Moreover, if  $\psi_1$  is an eigenfunction of  $L$  corresponding to the smallest eigenvalue  $\lambda_1 > 0$  with  $\|\psi_1\|_2 = 1$ , and if  $\langle u_0, \psi_1 \rangle \neq 0$ , then the solution  $u$  satisfies*

$$u(t) \rightarrow \frac{\gamma}{\gamma + \lambda_1} \mathbf{v} \text{ in } L^2(\Omega) \text{ as } t \rightarrow \infty \text{ for some } \mathbf{v} \in \{\pm\psi_1\}.$$

*Proof.* The proof is inspired by the Galerkin's method. Note that we can find an orthonormal basis  $\{\psi_k\}_{k \in \mathbb{N}}$  for  $L^2(\Omega)$  such that  $\psi_k$  is an eigenfunction of  $L$  corresponding to the  $k^{\text{th}}$  smallest eigenvalue  $\lambda_k$  with  $\psi_k \in H_0^1(\Omega)$ . Then, we let  $V_N$ ,  $N \in \mathbb{N}$ , be the subspace of  $L^2(\Omega)$  spanned by  $\{\psi_1, \dots, \psi_N\}$  and let  $P_N$  be the projection of  $L^2(\Omega)$  onto  $V_N$ .

Suppose that  $u_0 \in H_0^1(\Omega)$  is given and satisfies  $\langle u_0, \psi_1 \rangle \neq 0$ . We then consider

$$\begin{cases} \frac{\partial u}{\partial t} &= -Lu - \gamma \left(1 - \frac{1}{\|u(t)\|_2}\right) u \text{ in } \Omega_T, \\ u &= 0 \text{ on } \partial\Omega_T, \\ u(0) &= P_N(u_0). \end{cases} \quad (39)$$

If  $u_N \in L^2([0, T]; H_0^1(\Omega))$  with  $\frac{d}{dt} u_N \in L^2([0, T]; H^{-1}(\Omega))$  is a solution of (39), then for  $k = 1, \dots, N$ , we have

$$\frac{d}{dt} \langle u_N, \psi_k \rangle = -(\lambda_k + \gamma) \langle u_N, \psi_k \rangle + \frac{\gamma}{\|u_N(t)\|_2} \langle u_N, \psi_k \rangle$$

and

$$\frac{1}{2} \frac{d}{dt} \langle u_N, \psi_k \rangle^2 = -(\lambda_k + \gamma) \langle u_N, \psi_k \rangle^2 + \frac{\gamma}{\|u_N(t)\|_2} \langle u_N, \psi_k \rangle^2$$

with  $\langle u_N(0), \psi_k \rangle = \langle P_N(u_0), \psi_k \rangle$ . Considering  $\phi_k(t) = \langle u_N(t), \psi_k \rangle$ ,  $k = 1, \dots, N$ , we can see that  $\phi_1, \dots, \phi_N$  solve (54). Since  $\Phi_N = [\phi_1 \ \dots \ \phi_N]^T$  satisfies that for  $t \in [0, \infty)$ ,

$$\|u_N(t)\|_2 \geq \|\Phi_N(t)\| > \omega = \frac{1}{2} \min \left( \frac{\gamma}{\gamma + \lambda_N}, \|\Phi_N(0)\| \right). \quad (40)$$

Therefore, the existence of a solution to (39) can be obtained by a linear system of ODEs (54) given in the Appendix below, with the initial condition

$$\Phi_N(0) = [\langle P_N(u_0), \psi_1 \rangle \ \dots \ \langle P_N(u_0), \psi_N \rangle]^T,$$

and setting

$$u_N(x, t) = \phi_1(t)\psi_1(x) + \dots + \phi_N(t)\psi_N(x).$$

Then,  $u_N$  is a weak solution of (39) such that

$$u_N(t) \rightarrow \frac{\gamma}{\gamma + \lambda_1} \mathbf{v} \text{ in } L^2(\Omega) \text{ as } t \rightarrow \infty,$$

where  $\phi_1, \dots, \phi_N$  satisfy (54) and  $\mathbf{v} \in L^2(\Omega)$  is either  $\psi_1$  or  $-\psi_1$ .

Suppose now that  $u_{N,1}, u_{N,2}$  are two such solutions of (39). Let  $v_N = u_{N,1} - u_{N,2}$ . Then, we have

$$\frac{\partial v_N}{\partial t} = -Lv_N - \gamma(f(u_{N,1}) - f(u_{N,2})), \quad (41)$$

where  $f(u(x, t)) = u(x, t) - \frac{u(x, t)}{\|u(t)\|_2}$ . Since  $f$  on  $L^2(\Omega)$  is Lipschitz with a Lipschitz constant  $\mu = (\frac{2}{\omega} - 1)$  on  $\{u : \|u\|_2 \geq \omega\}$  and the two solutions  $u_{N,1}, u_{N,2}$  satisfy (40) for all  $t \geq 0$ , by taking the inner product on  $L^2(\Omega)$  with  $v_N$  on both sides of (41), we have

$$\frac{1}{2} \frac{d}{dt} \langle v_N(t), v_N(t) \rangle \leq \langle -Lv_N(t), v_N(t) \rangle + \gamma \mu \langle v_N(t), v_N(t) \rangle,$$

which implies

$$\frac{1}{2} \frac{d}{dt} \langle v_N, v_N \rangle - \gamma \mu \langle v_N, v_N \rangle \leq -\langle Lv_N, v_N \rangle \leq 0. \quad (42)$$

Therefore,  $\|v_N(0)\|_2 = 0$  implies  $u_{N,1} = u_{N,2}$  for a.e.  $(x, t) \in \Omega \times [0, T]$  for any  $T \in (0, \infty)$ . In fact, we know that

$$u_N \in C^1([0, \infty); H_0^1(\Omega))$$

as well as

$$u_N \in L^2([0, \infty); H_0^1(\Omega)), \quad \frac{d}{dt} u_N \in L^2([0, \infty); H^{-1}(\Omega)).$$

We now solve (38). Firstly, we fix  $k_0 \in \mathbb{N}$  so that  $\frac{\|u_0\|_2}{\sqrt{2}} < \|P_{k_0}(u_0)\|_2$  with  $\lambda_{k_0} < \lambda_{k_0+1}$  and obtain the solution  $u_N$  of (39) with  $N > k_0$ . We may write  $u_N$  as

$$u_N(t) = \phi_{1,N}(t)\psi_1 + \cdots + \phi_{N,N}(t)\psi_N.$$

Let  $\varphi_{1,N}(t) = \phi_{1,N}^2(t) + \cdots + \phi_{k_0,N}^2(t)$  and  $\varphi_{2,N}(t) = \phi_{k_0+1,N}^2(t) + \cdots + \phi_{N,N}^2(t)$ . Then, we have

$$\frac{1}{2} \frac{d\varphi_{1,N}}{dt} \geq -(\lambda_{k_0} + \gamma)\varphi_{1,N} + \frac{\gamma}{\|\Phi_N(t)\|} \varphi_{1,N}, \quad (43)$$

$$\frac{1}{2} \frac{d\varphi_{2,N}}{dt} \leq -(\lambda_{k_0+1} + \gamma)\varphi_{2,N} + \frac{\gamma}{\|\Phi_N(t)\|} \varphi_{2,N}. \quad (44)$$

This implies that

$$\varphi_{1,N}(t) \geq \varphi_{1,N}(0) e^{-2(\lambda_{k_0} + \gamma)t + \int_0^t \frac{2\gamma}{\|\Phi_N(s)\|} ds} \geq \varphi_{2,N}(0) e^{-2(\lambda_{k_0+1} + \gamma)t + \int_0^t \frac{2\gamma}{\|\Phi_N(s)\|} ds} \geq \varphi_{2,N}(t)$$

i.e.,

$$\frac{\varphi_{1,N}(t)}{\varphi_{2,N}(t)} \geq \frac{\varphi_{1,N}(0)}{\varphi_{2,N}(0)} e^{2(\lambda_{k_0+1} - \lambda_{k_0})t} > e^{2(\lambda_{k_0+1} - \lambda_{k_0})t} \quad (45)$$

due to  $\varphi_{1,N}(0) > \frac{\|u_0\|_2^2}{2} > \varphi_{2,N}(0)$ . Let

$$\begin{cases} M_1 &= 2 \max\left(\frac{\gamma}{\gamma + \lambda_1}, \|u_0\|_2\right), \\ M_2 &= \frac{1}{\sqrt{2}} \min\left(\frac{\gamma}{\lambda_{k_0} + \gamma}, \|u_0\|_2\right). \end{cases}$$

Note that

$$\frac{1}{2} \frac{d\varphi_{1,N}}{dt} \leq -(\lambda_1 + \gamma)\varphi_{1,N} + \frac{\gamma}{\|\Phi_N(t)\|} \varphi_{1,N}$$

implies  $\varphi_{1,N}(t) \leq \|\Phi_N(t)\|^2 < M_1^2$ . And (45) implies

$$\varphi_{2,N}(t) < M_1^2 e^{-2(\lambda_{k_0+1}-\lambda_{k_0})t}. \quad (46)$$

Since we already saw that  $\phi_{j,N}(t)$ ,  $j = 1, 2, \dots, N$ , exist for all  $t \in (0, \infty)$ , if we suppose

$$t_0 = \inf\{t \in (0, \infty) | \varphi_{1,N}(t) = M_2^2\} < \infty,$$

then from (43) and (45), we have

$$\|\Phi_N(t_0)\|^2 = \varphi_{1,N}(t_0) + \varphi_{2,N}(t_0) < M_2^2(1 + e^{-2(\lambda_{k_0+1}-\lambda_{k_0})t_0}) < 2M_2^2 \leq \left(\frac{\gamma}{\lambda_{k_0} + \gamma}\right)^2.$$

Therefore, there exists  $\delta > 0$  such that for  $t \in (t_0 - \delta, t_0)$ ,

$$\|\Phi_N(t)\| < \frac{\gamma}{\lambda_{k_0} + \gamma},$$

which implies that

$$\frac{1}{2} \frac{d\varphi_{1,N}}{dt}(t) \geq -(\lambda_{k_0} + \gamma)\varphi_{1,N}(t) + \frac{\gamma}{\|\Phi_N\|}\varphi_{1,N}(t) > 0 \quad \text{on } (t_0 - \delta, t_0).$$

This is a contradiction since  $\varphi_{1,N}(t) > M_2^2$  for  $t \in [0, t_0)$ . Therefore, we end up with

$$M_2^2 < \varphi_{1,N}(t) < M_1^2 \quad \text{for } t \geq 0. \quad (47)$$

Note that (47) implies

$$\inf_{t \geq 0} \|\Phi_N(t)\| \geq M_2$$

That is, the solution  $u_N$  for  $N > k_0$  satisfies

$$\inf_{t \geq 0} \|u_N(t)\|_2 \geq M_2. \quad (48)$$

We fix  $T \in (0, \infty)$  and consider a sequence of solutions  $\{u_N(t)\}_{N > k_0}$ , where  $u_N$  is the solution of (39) with  $u_N(0) = P_N(u_0)$ . Using  $v = u_N - u_M$  with  $N, M > k_0$  in place of  $v_N$  in (42), and noting that we may choose  $\omega = M_2$  for  $\mu = (\frac{2}{\omega} - 1)$ , we obtain

$$\|u_N(t) - u_M(t)\|_2 \leq e^{\gamma\mu t} \|u_N(0) - u_M(0)\|_2 \leq e^{\gamma\mu T} \|u_N(0) - u_M(0)\|_2 \quad \text{for } t \in [0, T], \quad (49)$$

which implies that for any  $m \in \mathbb{N}$ ,  $\{\phi_{m,N}(t)\}_{N > k_0}$  converges uniformly to  $\phi_{m,*}(t)$  in  $[0, T]$  and

$$\|\Phi_N(t)\| \rightarrow \|\Phi_*(t)\| \quad \text{uniformly on } [0, T] \text{ as } N \rightarrow \infty,$$

where  $\phi_{m,N}(t) = \langle u_N(t), \psi_m \rangle$  and  $\|\Phi_*(t)\| = \sqrt{\sum_{m=1}^{\infty} \phi_{m,*}^2(t)}$ . Hence, for all  $m$ , we have

$$\frac{d}{dt} \phi_{m,*} = -(\lambda_m + \gamma)\phi_{m,*} + \frac{\gamma}{\|\Phi_*\|} \phi_{m,*}.$$

In addition, for  $k_0 < m \leq N$ , if we consider

$$\omega_{m,N}(t) = \phi_{m,N}^2(t) + \dots + \phi_{N,N}^2(t),$$

then we can obtain, by the same argument for (46),

$$\omega_{m,N}(t) < M_1^2 e^{-2(\lambda_m - \lambda_{k_0})t}, \quad t \geq 0 \quad (50)$$



implying  $|\phi_{m,N}(t)| < M_1 e^{-(\lambda_m - \lambda_{k_0})t}$  on  $[0, T]$  and, eventually, on  $(0, \infty)$ . By a slight modification on (46), we can see that even for each  $2 \leq m \leq k_0$ , there exists  $\zeta_m > 0$  such that  $|\phi_{m,N}(t)| < \zeta_m M_1 e^{-(\lambda_m - \lambda_1)t}$  on  $[0, \infty)$ , i.e., with  $\eta_{k_0} = \max_{k=2, \dots, k_0}(\zeta_k)$ ,

$$|\phi_{m,N}(t)| < \begin{cases} \eta_{k_0} M_1 e^{-(\lambda_m - \lambda_1)t}, & 2 \leq m \leq k_0, \\ M_1 e^{-(\lambda_m - \lambda_{k_0})t}, & m > k_0. \end{cases} \quad (51)$$

This implies that

$$\left\{ \sum_{m=1}^N \lambda_m \phi_{m,N}^2 \right\}_{N \in \mathbb{N}} \text{ converges uniformly to } \sum_{m=1}^{\infty} \lambda_m \phi_{m,*}^2 \text{ in } [0, T] \text{ as } N \rightarrow \infty.$$

Hence, by defining

$$u_*(x, t) = \sum_{m=1}^{\infty} \phi_{m,*}(t) \psi_m(x),$$

we can easily see that  $u_* \in L^2([0, T]; H_0^1(\Omega))$  and

$$\frac{d}{dt} u_* = \sum_{m=1}^{\infty} \frac{d\phi_{m,*}}{dt}(t) \psi_m(x) \in L^2([0, T]; H^{-1}(\Omega))$$

and that for a.e.  $t \in [0, T]$ , and for each  $v \in H_0^1(\Omega)$ ,

$$\begin{aligned} & \int_{\Omega} \frac{\partial}{\partial t} u_*(x, t) v(x) dx + \int_{\Omega} \sum_{i,j=1}^d a_{i,j}(x) \partial_j u_*(x) \partial_i v(x) dx + \int_{\Omega} c(x) u_*(x) v(x) dx \\ & \gamma \int_{\Omega} u_*(x) v(x) dx - \frac{\gamma}{\|u_*(t)\|_2} \int_{\Omega} u_*(x) v(x) dx = 0 \end{aligned}$$

with  $u_*(0) = u_0$  in  $L^2(\Omega)$ . Therefore,  $u_*$  is a weak solution of (38) for any  $T \in (0, \infty)$ . Next, using the functional in (36), we define  $f(t)$  for  $t \in (0, \infty)$  by

$$\begin{aligned} f(t) = F_L(u_*) &= \frac{1}{2} \int_{\Omega} \sum_{i,j=1}^d a_{i,j}(x) \partial_j u_*(x, t) \partial_i u_*(x, t) dx + \frac{1}{2} \int_{\Omega} c(x) |u_*(x, t)|^2 dx \\ &+ \frac{\gamma}{2} \int_{\Omega} |u_*(x, t)|^2 dx - \gamma \left( \int_{\Omega} |u_*(x, t)|^2 dx \right)^{\frac{1}{2}}, \end{aligned}$$

Note that

$$f(t) = \frac{1}{2} \sum_{m=1}^{\infty} \lambda_m \phi_{m,*}^2(t) + \frac{\gamma}{2} \sum_{m=1}^{\infty} \phi_{m,*}^2(t) - \gamma \|\Phi_*(t)\|.$$

Since  $\phi_{m,*}^2(t)$  decays exponentially to 0 in  $(0, \infty)$  as  $m \rightarrow \infty$ , we can see that

$$\begin{aligned} \frac{df}{dt}(t) &= \sum_{m=1}^{\infty} \left( \lambda_m \phi_{m,*}(t) + \gamma \phi_{m,*}(t) - \gamma \frac{\phi_{m,*}(t)}{\|\Phi_*(t)\|} \right) \frac{d\phi_{m,*}}{dt}(t) \\ &= - \sum_{m=1}^{\infty} \left| \frac{d\phi_{m,*}}{dt}(t) \right|^2. \end{aligned}$$

Therefore,  $f$  is non-increasing and bounded below, i.e.,  $a = \lim_{t \rightarrow \infty} f(t)$  exists. From the fact that for all  $N \in \mathbb{N}$ ,

$$\|\Phi_N(t)\| \rightarrow \frac{\gamma}{\gamma + \lambda_1} \text{ as } t \rightarrow \infty,$$

together with (51), we can see that

$$\|\Phi_*(t)\| \rightarrow \frac{\gamma}{\gamma + \lambda_1} \text{ as } t \rightarrow \infty.$$

which eventually implies that for each  $m \in \mathbb{N}$ ,  $|\frac{d\phi_{m,*}}{dt}(t)| \rightarrow 0$  as  $t \rightarrow \infty$ , i.e.,

$$\lim_{t \rightarrow \infty} \left( (\lambda_m + \gamma) - \frac{\gamma}{\|\Phi_*(t)\|} \right) \phi_{m,*}(t) = 0. \quad (52)$$

Moreover, since we have  $|\phi_{m,*}(t)| \leq \max_{k=1, \dots, k_0} \{M_1 \zeta_k e^{-(\lambda_2 - \lambda_1)t}\}$  for all  $m \geq 2$  with  $\zeta_1 = 1$  and

$$\inf_{t \geq 0} \|\Phi_*(t)\| \geq M_2,$$

we finally obtain, together with (52),

$$\lim_{t \rightarrow \infty} \|\Phi_*(t)\| = \frac{\gamma}{\lambda_1 + \gamma},$$

which implies that  $\lim_{t \rightarrow \infty} |\phi_{1,*}(t)| = \frac{\gamma}{\lambda_1 + \gamma}$ , and

$$u_*(t) \rightarrow \frac{\gamma}{\lambda_1 + \gamma} \mathbf{v} \text{ in } L^2(\Omega) \text{ as } t \rightarrow \infty$$

for some  $\mathbf{v} \in \{\pm\psi_1\}$ . □

## 4 Appendix: Corresponding ODEs

For a given symmetric positive semidefinite matrix  $A$  in  $\mathbf{M}_N(\mathbb{R})$ , as we saw in (53) and in (6), with an orthogonal matrix  $Q$  and a diagonal matrix  $\Lambda_N = \text{diag}(\lambda_1, \dots, \lambda_N)$ ,  $0 \leq \lambda_1 = \dots = \lambda_p < \lambda_{p+1} \leq \dots \leq \lambda_N$  for some  $1 \leq p < N$ , we have  $A = Q\Lambda_N Q^T$  and

$$F_A(\mathbf{x}) = F_{\Lambda_N}(\mathbf{y}), \text{ with } \mathbf{y} = Q^T \mathbf{x}. \quad (53)$$

The gradient descent flow associated with the functional  $F_{\Lambda_N}$  is

$$\frac{d}{dt} \Phi_N(t) = -\nabla F_{\Lambda_N}(\Phi_N(t)),$$

and we are interested in the existence of a solution of (54). So we solve on  $(0, \infty)$

$$\begin{cases} \frac{d}{dt} \Phi_N(t) &= -\nabla F_{\Lambda_N}(\Phi_N(t)) = -(\Lambda_N + \gamma I) \Phi_N(t) + \frac{\gamma}{\|\Phi_N(t)\|} \Phi_N(t), \\ \Phi_N(0) &\in \mathbb{R}^N, \phi_j(0) \neq 0 \text{ for some } 1 \leq j \leq p, \end{cases} \quad (54)$$

where  $\Phi_N = \Phi_N(t) = [\phi_1(t) \cdots \phi_N(t)]^T$  and  $\|\Phi_N(t)\| = \sqrt{\sum_{k=1}^N \phi_k^2(t)}$ .

**Theorem 8.** *There exist a unique solution  $\Phi_N \in C^1([0, \infty))$  of (54) and  $\mathbf{v} \in \mathbb{R}^N$  such that*

$$\Phi_N(t) \rightarrow \mathbf{v} \text{ as } t \rightarrow \infty,$$

where  $\mathbf{v}$  is an eigenvector of  $\Lambda_N$  corresponding to the eigenvalue  $\lambda_1$  with  $\|\mathbf{v}\| = \frac{\gamma}{\gamma + \lambda_1}$  and  $\langle \mathbf{v}, \mathbf{e}_j \rangle \neq 0$ . Note that  $\{\mathbf{e}_1, \dots, \mathbf{e}_N\}$  is the standard basis for  $\mathbb{R}^N$ .

*Proof.* The existence and uniqueness of a solution  $\Phi_N$  on  $[0, \infty)$  is easily guaranteed by the theory of ODEs once we establish a lower bound  $\omega > 0$  for  $\|\Phi_N(t)\|$ ,  $t \in [0, \infty)$ . Due to  $\phi_j(0) \neq 0$ , i.e.,  $\|\Phi_N(0)\| \neq 0$ , a solution exists and is unique on  $[0, \epsilon]$  with  $\epsilon \ll 1$ . Then, by setting

$$T_{max} = \sup\{T \in (0, \infty) : \text{a solution } \Phi_N \text{ exists and is unique on } [0, T]\},$$

we know that  $T_{max} \geq \epsilon$ . Note that for  $t \in (0, T_{max})$  and for  $1 \leq k \leq N$ , we have

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \phi_k^2(t) &= -(\lambda_k + \gamma) \phi_k^2(t) + \frac{\gamma}{\|\Phi_N(t)\|} \phi_k^2(t) \\ &> -(\lambda_N + \gamma) \phi_k^2(t), \end{aligned}$$

resulting in  $\phi_k(t) \neq 0$  for  $t \geq 0$  if and only if  $\phi_k(0) \neq 0$ . With  $\omega = \frac{1}{2} \min(\frac{\gamma}{\gamma + \lambda_N}, \|\Phi_N(0)\|)$ , we suppose  $\{t \in (0, T_{max}) : \|\Phi_N(t)\| \leq \omega\} \neq \emptyset$ . Then,  $T_\omega = \inf\{t \in (0, \epsilon) : \|\Phi_N(t)\| \leq \omega\}$  exists in  $(0, T_{max})$ . Since there exists  $\delta > 0$  such that  $\omega < \|\Phi_N(t)\| < \frac{3}{2}\omega$  on  $(T_\omega - \delta, T_\omega)$ , we have

$$\frac{1}{2} \frac{d}{dt} \|\Phi_N(t)\|^2 \geq \frac{1}{3} (\lambda_N + \gamma) \|\Phi_N(t)\|^2 > 0 \text{ on } (T_\omega - \delta, T_\omega),$$

that is,  $\|\Phi_N(t)\|$  is increasing in  $(T_\omega - \delta, T_\omega)$  and  $\lim_{t \rightarrow T_\omega} \|\Phi_N(t)\| > \omega = \|\Phi_N(T_\omega)\|$ . This is a contradiction. Therefore, we have

$$\inf_{t \in (0, T_{max})} \|\Phi_N(t)\| > \omega,$$

which eventually implies that  $T_{max} = \infty$ .

In addition, we note that

$$\frac{d}{dt} F_{\Lambda_N}(\Phi_N) = \left\langle \frac{d}{dt} \Phi_N, \Lambda_N \Phi_N + \gamma \Phi_N - \frac{\gamma}{\|\Phi_N\|} \Phi_N \right\rangle = - \left\| \frac{d\Phi_N}{dt} \right\|^2$$

and that since  $F_{\Lambda_N}(\Phi_N)$  is bounded below, there exists  $y \in \mathbb{R}$  such that

$$\lim_{t \rightarrow \infty} F_{\Lambda_N}(\Phi_N(t)) = y \tag{55}$$

implying that

$$\lim_{t \rightarrow \infty} \left\| \frac{d\Phi_N}{dt} \right\|^2 = 0 \Leftrightarrow \lim_{t \rightarrow \infty} \left\| \Lambda_N \Phi_N + \gamma \left(1 - \frac{1}{\|\Phi_N\|}\right) \Phi_N \right\|^2 = 0.$$

This implies that for each  $k = 1, 2, \dots, N$ ,

$$\lim_{t \rightarrow \infty} \left( \lambda_k + \gamma - \frac{\gamma}{\|\Phi_N(t)\|} \right)^2 \phi_k^2(t) = 0. \tag{56}$$

Since

$$0 = \lim_{t \rightarrow \infty} \left\| \Lambda_N \Phi_N + \gamma \left(1 - \frac{1}{\|\Phi_N\|}\right) \Phi_N \right\| \geq \lim_{t \rightarrow \infty} \left| \|(\Lambda_N + \gamma I) \Phi_N\| - \gamma \right|,$$

we have

$$\lim_{t \rightarrow \infty} \sum_{k=1}^N (\lambda_k + \gamma)^2 \phi_k^2(t) = \gamma^2, \quad (57)$$

which implies that there exist  $k_1$  and  $\{t_n\}$  such that  $t_n \rightarrow \infty$  as  $n \rightarrow \infty$  and

$$\lim_{n \rightarrow \infty} \phi_{k_1}^2(t_n) > 0.$$

Unless  $\lim_{n \rightarrow \infty} \phi_{k_1}^2(t_n) = \frac{\gamma^2}{(\gamma + \lambda_{k_1})^2}$ , there exists  $k_2 \neq k_1$  such that  $\lim_{n \rightarrow \infty} \phi_{k_2}^2(t_n) > 0$  by taking a subsequence of  $\{t_n\}$  if necessary. Then, (56) implies

$$\lambda_{k_1} = \lambda_{k_2} \quad \text{and} \quad \lim_{n \rightarrow \infty} \|\Phi_N(t_n)\| = \frac{\gamma}{\gamma + \lambda_{k_1}}.$$

We may repeat this process finitely many times to have

$$\lim_{n \rightarrow \infty} \sum_{\{k: \lambda_k = \lambda_{k_1}\}} \phi_k^2(t_n) = \frac{\gamma^2}{(\gamma + \lambda_{k_1})^2}, \quad \lim_{n \rightarrow \infty} \sum_{\{k: \lambda_k \neq \lambda_{k_1}\}} \phi_k^2(t_n) = 0.$$

Suppose that there exist  $l$  and  $\{s_n\}$  such that  $\lambda_l \neq \lambda_{k_1}$  and  $s_n \rightarrow \infty$  as  $n \rightarrow \infty$  and

$$\lim_{n \rightarrow \infty} \phi_l^2(s_n) > 0.$$

Then, as was done above, with a subsequence of  $\{s_n\}$  if necessary, we have

$$\lim_{n \rightarrow \infty} \|\Phi_N(s_n)\| = \frac{\gamma}{\gamma + \lambda_l}, \quad \lim_{n \rightarrow \infty} \sum_{\{k: \lambda_k = \lambda_l\}} \phi_k^2(s_n) = \frac{\gamma^2}{(\gamma + \lambda_l)^2}, \quad \lim_{n \rightarrow \infty} \sum_{\{k: \lambda_k \neq \lambda_l\}} \phi_k^2(s_n) = 0$$

and (55) implies that

$$\lim_{n \rightarrow \infty} F_{\Lambda_N}(\Phi_N(t_n)) = -\frac{\gamma^2}{2(\gamma + \lambda_{k_1})} = y = \lim_{n \rightarrow \infty} F_{\Lambda_N}(\Phi_N(s_n)) = -\frac{\gamma^2}{2(\gamma + \lambda_l)},$$

which is a contradiction. Therefore, we conclude that

$$\lim_{t \rightarrow \infty} \|\Phi_N(t)\| = \frac{\gamma}{\gamma + \lambda_{k_1}}, \quad \lim_{t \rightarrow \infty} \sum_{\{k: \lambda_k = \lambda_{k_1}\}} \phi_k^2(t) = \frac{\gamma^2}{(\gamma + \lambda_{k_1})^2}, \quad \lim_{t \rightarrow \infty} \sum_{\{k: \lambda_k \neq \lambda_{k_1}\}} \phi_k^2(t) = 0.$$

We will now claim that  $\lambda_{k_1} = \lambda_1$ . Suppose that  $\lambda_{k_1} > \lambda_1$ , i.e.,  $k_1 > p$ . Then, since there exists  $T > 0$  such that  $\|\Phi_N(t)\| < \frac{\gamma(2\gamma + \lambda_1 + \lambda_{k_1})}{2(\gamma + \lambda_1)(\gamma + \lambda_{k_1})}$  for  $t > T$ , we have

$$\frac{1}{2} \frac{d}{dt} \sum_{k=1}^p \phi_k^2(t) = \left( -(\lambda_1 + \gamma) + \frac{\gamma}{\|\Phi_N(t)\|} \right) \sum_{k=1}^p \phi_k^2(t) \geq \left( \frac{\lambda_{k_1} - \lambda_1}{2\gamma + \lambda_1 + \lambda_{k_1}} \right) \sum_{k=1}^p \phi_k^2(t)$$

for  $t > T$ , which results in  $\|\Phi_N(t)\| \rightarrow \infty$  as  $t \rightarrow \infty$ . This is a contradiction. Therefore, we have  $\lambda_{k_1} = \lambda_1$  and

$$\lim_{t \rightarrow \infty} \|\Phi_N(t)\|^2 = \lim_{t \rightarrow \infty} \sum_{k=1}^p \phi_k^2(t) = \frac{\gamma^2}{(\gamma + \lambda_1)^2}, \quad \text{and} \quad \lim_{t \rightarrow \infty} \sum_{k=p+1}^N \phi_k^2(t) = 0. \quad (58)$$

Lastly, we will claim that there exists  $\mathbf{v} \in \mathbb{R}^N$  such that  $\|\mathbf{v}\| = \frac{\gamma}{\gamma + \lambda_1}$ , and  $\langle \mathbf{v}, \mathbf{e}_j \rangle \neq 0$ , and

$$\Phi_N(t) \rightarrow \mathbf{v} \text{ as } t \rightarrow \infty.$$

Note that for  $k = 1, \dots, p$ , if  $\phi_k(0) \neq 0$ , then

$$\frac{1}{2\phi_k^2(t)} \frac{d}{dt} \phi_k^2(t) = -(\lambda_1 + \gamma) + \frac{\gamma}{\|\Phi_N(t)\|} \text{ for } t \geq 0,$$

implying that for  $T \in (0, \infty)$ ,

$$\ln(\phi_k^2(T)) - \ln(\phi_k^2(0)) = 2 \int_0^T \left( -(\lambda_1 + \gamma) + \frac{\gamma}{\|\Phi_N(t)\|} \right) dt := 2\Psi(T).$$

Together with (58), we have

$$\sum_{k=1}^p \phi_k^2(T) = e^{2\Psi(T)} \sum_{k=1}^p \phi_k^2(0) \rightarrow \left( \frac{\gamma}{\gamma + \lambda_1} \right)^2 \text{ as } T \rightarrow \infty.$$

Note that

$$\Psi_\infty = \lim_{t \rightarrow \infty} \Psi(t) = \int_0^\infty \left( -(\lambda_1 + \gamma) + \frac{\gamma}{\|\Phi_N(t)\|} \right) dt = \ln \left( \frac{\gamma}{\gamma + \lambda_1} \right) - \frac{1}{2} \ln \left( \sum_{k=1}^p \phi_k^2(0) \right).$$

Hence, for  $1 \leq k \leq p$ ,

$$\phi_k^2(t) \rightarrow v_k^2 \text{ as } t \rightarrow \infty,$$

where  $v_k = \phi_k(0)e^{\Psi_\infty}$ . Noting that  $\phi_k(t)\phi_k(0) > 0$  for all  $t \geq 0$  unless  $\phi_k(0) = 0$ , we can easily see that

$$\phi_k(t) \rightarrow v_k \text{ as } t \rightarrow \infty.$$

If we set

$$\mathbf{v} = [v_1 \ \cdots \ v_p \ 0 \ \cdots \ 0]^T,$$

then  $\mathbf{v}$  is an eigenvector of  $\Lambda_N$  corresponding to  $\lambda_1$  and  $\langle \mathbf{v}, \mathbf{e}_j \rangle = v_j \neq 0$  since  $\phi_j(0) \neq 0$  for some  $1 \leq j \leq p$ , and

$$\Phi_N(t) \rightarrow \mathbf{v} \text{ as } t \rightarrow \infty$$

and  $\|\mathbf{v}\| = \frac{\gamma}{\gamma + \lambda_1}$ . □

*Remark 1.* As was proved in Theorem 8, even though a symmetric positive semidefinite matrix  $A \in M_N(\mathbb{R})$  has the smallest eigenvalue of multiplicity greater than 1, the solution of (54) converges to a corresponding eigenvector just as the solution of (9) does. So the discrete and continuous settings present the same behavior. Moreover, due to the equivalence (53), the solution of the gradient descent flow associated to  $F_A$  converges to  $Q\mathbf{v}$ , which is an eigenvector of  $A$  corresponding to the smallest eigenvalue  $\lambda_1$ .

## 5 Conclusion

In this paper, we have proposed and analyzed to minimize the functional (3) using two well-known methods: The Gradient Descent and The Newton's methods. The two algorithms have shown their advantages in various situations. More interestingly, these methods revealed what conventional methods haven't observed.

More precisely, applying the Gradient Descent method to minimize (3) guarantees that we can find a global minimizer of our proposed functional and allows us to find eigenvectors in the increasing order of their corresponding eigenvalues without matrix inversion and produces many applications such as generalized eigenvalue problems and finding eigenfunctions of self-adjoint operators.

On the other hand, besides its faster convergence, the Newton's method proposes a new nonsingular linear system that we can solve once to find an exact eigenvector when an exact eigenvalue is given. Even with an approximate eigenvalue, the nonsingular linear system guarantees that the unique solution is close to an exact eigenvector as much the estimated eigenvalue is close to an exact eigenvalue as possible, which reveals that the error in eigenvalue estimation and the error in eigenvector estimation are comparable.

Moreover, it turns out that the same framework can extend to nonlinear operators as well. One such example is the  $p$ -Laplacian operator  $-\Delta_p$ ,  $p > 2$ . As before, we want to find  $\varphi$  such that

$$\begin{cases} -\Delta_p \varphi = \lambda \varphi, & \text{in } \Omega, \\ \varphi = 0, & \text{on } \partial\Omega, \end{cases} \quad (59)$$

where  $-\Delta_p \varphi = \operatorname{div}(|\nabla \varphi|^{p-2} \nabla \varphi)$  by minimizing the functional  $F_p$  defined by

$$F_p(u) = \frac{1}{2} \int_{\Omega} |\nabla u(x)|^p dx + \frac{\gamma}{2} \int_{\Omega} |u(x)|^2 dx - \gamma \left( \int_{\Omega} |u(x)|^2 dx \right)^{\frac{1}{2}} \quad (60)$$

Note that

$$\min_{u \in W_0^{1,p}(\Omega)} F_p(u)$$

exists and a minimizer satisfies (59).

We will provide detailed analyses and numerical computations in a subsequent paper on applications of our proposed framework, including those mentioned in this paper, to prove efficiency and usefulness of our proposed method by revealing what have not been known through conventional methods in finding eigenvalues and eigenvectors.

## References

- [1] A.P. Austin and L.N. Trefethen, *Computing eigenvalues of real symmetric matrices with rational filters in real arithmetic*, SIAM J. Sci. Comput., 37(3), pp. A1365-A1387, 2015
- [2] C.A. Beattie, M. Embree and D.C. Sorensen, *Convergence of polynomial restart Krylov methods for eigenvalue computations*, SIAM Review, 47(3), pp. 492-515, 2005
- [3] M. Belkin, J. Sun and Y. Wang, *Constructing Laplace operator from point clouds in  $\mathbb{R}^d$* , In Proceedings of the Twentieth Annual ACM-SIAM Symposium on Discrete Algorithms, pp. 1031-1040, Philadelphia, PA, USA, 2009
- [4] J.E. Dennis and R.A. Tapia, *Inverse, shifted inverse, and Rayleigh quotient iteration as Newton's method*, <http://www.caam.rice.edu/~rat/cv/RQI.pdf>
- [5] X.B. Gao, G.H. Golub and L.Z. Liao, *Continuous methods for symmetric generalized eigenvalue problems*, Linear Alg. Appl., 428, pp. 676-696, 2008
- [6] J.M. Gedicke, *On the numerical analysis of eigenvalue problems*, <http://edoc.hu-berlin.de/dissertationen/gedicke-joscha-micha-2013-06-10/PDF/gedicke.pdf>, Dissertation, 2013

- [7] R. Kolluri, J.R. Shewchuk and J.F. O'Brien, *Spectral surface reconstruction from noisy point clouds*, Proceedings of the 2004 Eurographics/ACM SIGGRAPH Symposium on Geometry Processing, SGP '04, pp. 11-21, 2004
- [8] V. Kuleshov, *Fast algorithms for sparse principal component analysis based on Rayleigh quotient iteration*, JMLR W& CP 28(3), pp. 1418-1425, 2013
- [9] R. Lai, J. Liang and H. Zhao, *A local mesh method for solving PDEs on point clouds*, Inverse Probl. Imaging, 7(3), pp. 737-755
- [10] E. Mengi, E.A. Yildirim and M. Kilic, *Numerical optimization of eigenvalues of hermitian matrix functions*, SIAM J. Matrix Anal. Appl., 35(2), pp. 699-724, 2014
- [11] Y. Notay, *Convergence analysis of inexact Rayleigh quotient iteration*, SIAM J. Matrix Anal. Appl., 24(3), pp. 627-644, 2003
- [12] G.L.G. Sleijpen and H.A. Van der Vorst, *A Jacobi-Davidson iteration method for linear eigenvalue problems*, SIAM J. Matrix Anal. Appl., 17(2), pp. 401-425, 1996
- [13] G. Still, *Computable Bounds for Eigenvalues and Eigenfunctions of Elliptic Differential Operators*, Numer. Math., 54, pp. 201-223, 1988