

# Numerical Solution of Linear Eigenvalue Problems

Jessica Bosch and Chen Greif

ABSTRACT. We review numerical methods for computing eigenvalues of matrices. We start by considering the computation of the dominant eigenpair of a general dense matrix using the power method, and then generalize to orthogonal iterations and the QR iteration with shifts. We also consider divide-and-conquer algorithms for tridiagonal matrices. The second part of this survey involves the computation of eigenvalues of large and sparse matrices. The Lanczos and Arnoldi methods are developed and described within the context of Krylov subspace eigensolvers. We also briefly present the idea of the Jacobi–Davidson method.

## 1. Introduction

Eigenvalue problems form one of the central problems in Numerical Linear Algebra. They arise in many areas of sciences and engineering. In this survey, we study linear eigenvalue problems. The standard algebraic eigenvalue problem has the form

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}.$$

We consider real eigenvalue problems, i.e.,  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . In some places we will use the notion of complex matrices, as they are crucial in mathematical as well computational aspects of eigenvalue solvers. There is a wide range of publications dealing with numerical methods for solving eigenvalue problems, e.g., [39, 66, 22, 47, 62, 55, 32, 54]. Typically, eigensolvers are classified into methods for symmetric (or Hermitian) and nonsymmetric (or non-Hermitian) matrices, or methods for small, dense matrices and large, sparse matrices.

This survey reviews popular methods for computing eigenvalues of a given matrix. It follows a minicourse presented by the second author at the 2015 Summer School on “Geometric and Computational Spectral Theory” at the Université de Montréal, and can be viewed as a set of comprehensive lecture notes.

When it comes to numerical computation of eigenvalues, it is reasonable to classify eigensolvers by the size and the nonzero pattern of the matrix. As opposed to the solution of linear systems, where it is possible to obtain a solution within a finite number of steps, most eigenvalue computations (except trivial cases such as a diagonal or a triangular matrix) require an iterative process. For matrices that

are not particularly large and do not have a specific nonzero structure, eigensolvers are often based on matrix decompositions. One may be interested in a small number of the eigenvalues and/or eigenvectors, or all of them, and there are methods that are available for accomplishing the stated goal. On the other hand, when the matrix is large and sparse, it is rare to seek the entire spectrum; in most cases we are interested in just a few eigenvalues and eigenvectors, and typical methods are based on matrix-vector products rather than matrix decompositions. Interestingly, despite the fact that all processes of eigenvalue computations are iterative, methods that are based on matrix decompositions are often referred to as *direct*, whereas methods that are based on matrix-vector products are termed *iterative*. This slight abuse of terminology is nonetheless widely understood and typically does not cause any confusion.

It is a bit ambitious to talk in general terms about a recipe for solution of eigenvalue problems, but it is legitimate to identify a few main components. A typical eigensolver starts with applying similarity transformations and transforming the matrix into one that has an appealing nonzero structure: for example tridiagonal if the original matrix was symmetric. Once this is accomplished, an iterative process is pursued, whereby repeated orthogonal similarity transformations are applied to get us closer and closer to a diagonal or triangular form. For large and sparse matrices, an additional component, generally speaking, in state of the art methods, is the transformation of the problem to a small and dense one on a projected subspace.

This survey devotes a significant amount of space to elaborating on the above principles. It is organized as follows. In Section 2, we briefly review basic concepts of Numerical Linear Algebra that are related to eigenvalue problems. We start with presenting methods for computing a few or all eigenvalues for small to moderate-sized matrices in Section 3. This is followed by a review of eigenvalue solvers for large and sparse matrices in Section 4. Conclusions complete the paper.

## 2. Background in Numerical Linear Algebra

**2.1. Theoretical basics.** We begin our survey with a review of basic concepts in Numerical Linear Algebra. We introduce some notation used throughout the survey.

Let  $\mathbf{A} \in \mathbb{C}^{m \times n}$ . The *kernel* or *nullspace* of  $\mathbf{A}$  is given as

$$\ker(\mathbf{A}) = \{\mathbf{x} \in \mathbb{C}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\}.$$

Another important subspace, which is often related to the kernel, is the *range* of  $\mathbf{A}$ , which is given as

$$\text{ran}(\mathbf{A}) = \{\mathbf{A}\mathbf{x} : \mathbf{x} \in \mathbb{C}^n\}.$$

The *rank* of  $\mathbf{A}$  is the maximal number of linearly independent columns (or rows), i.e.,

$$\text{rank}(\mathbf{A}) = \dim(\text{ran}(\mathbf{A})).$$

It holds  $n = \text{rank}(\mathbf{A}) + \dim(\ker(\mathbf{A}))$ .  $\mathbf{A}$  is called *rank-deficient* if  $\text{rank}(\mathbf{A}) < \min\{m, n\}$ .

In what follows, we consider real matrices  $\mathbf{A} \in \mathbb{R}^{n \times n}$  if not stated otherwise.

DEFINITION 2.1 (Invertibility).  $\mathbf{A}$  is called *invertible* or *nonsingular* if there exists a matrix  $\mathbf{B} \in \mathbb{R}^{n \times n}$  such that

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I}.$$

Here,  $\mathbf{I} \in \mathbb{R}^{n \times n}$  is the *identity matrix*. The *inverse* of  $\mathbf{A}$  is uniquely determined, and we denote it by  $\mathbf{A}^{-1}$ .

Related to the inverse and a matrix norm  $\|\cdot\|$  is the *condition number*, which is defined for a general square matrix  $\mathbf{A}$  as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|.$$

In general, if  $\kappa(\mathbf{A})$  is large<sup>1</sup>, then  $\mathbf{A}$  is said to be an *ill-conditioned* matrix. Useful matrix norms include the well-known *p-norms* or the *Frobenius norm*  $\|\cdot\|_F$ , which is given as

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{i,j}|^2},$$

where  $a_{i,j}$  is the  $(i, j)$  entry of  $\mathbf{A}$ .

Let us come to the heart of this paper. The *algebraic eigenvalue problem* has the following form:

DEFINITION 2.2 (Algebraic Eigenvalue Problem).  $\lambda \in \mathbb{C}$  is called an *eigenvalue* of  $\mathbf{A}$  if there exists a vector  $\mathbf{0} \neq \mathbf{x} \in \mathbb{C}^n$  such that

$$\mathbf{Ax} = \lambda \mathbf{x}.$$

The vector  $\mathbf{x}$  is called a (right) *eigenvector* of  $\mathbf{A}$  associated with  $\lambda$ . We call the pair  $(\lambda, \mathbf{x})$  an *eigenpair* of  $\mathbf{A}$ . The set of all eigenvalues of  $\mathbf{A}$  is called the *spectrum* of  $\mathbf{A}$  and is denoted by  $\lambda(\mathbf{A})$ .

Note from the above definition that real matrices can have complex eigenpairs. Geometrically, the action of a matrix  $\mathbf{A}$  expands or shrinks any vector lying in the direction of an eigenvector of  $\mathbf{A}$  by a scalar factor. This scalar factor is given by the corresponding eigenvalue of  $\mathbf{A}$ .

REMARK 2.3. Similarly, a left eigenvector of  $\mathbf{A}$  associated with the eigenvalue  $\lambda$  is defined as a vector  $\mathbf{0} \neq \mathbf{y} \in \mathbb{C}^n$  that satisfies  $\mathbf{y}^* \mathbf{A} = \lambda \mathbf{y}^*$ . Here,  $\mathbf{y}^* = \bar{\mathbf{y}}^T$  is the *conjugate transpose* of  $\mathbf{y}$ .

Throughout the survey, we use the term *eigenvector* for a right eigenvector.

Another way to define eigenvalues is the following:

DEFINITION 2.4 (Characteristic Polynomial). Let  $\det(\cdot)$  denote the determinant of a matrix. Then the polynomial

$$p_{\mathbf{A}}(x) = \det(\mathbf{A} - x\mathbf{I})$$

is called the *characteristic polynomial* of  $\mathbf{A}$ . It is a polynomial of degree  $n$ . The roots of  $p_{\mathbf{A}}(x)$  are the eigenvalues of  $\mathbf{A}$ .

<sup>1</sup>Of course, this depends on the definition of “large”; see, e.g., [22, Chap. 3.5].

The eigenvalues of  $\mathbf{A}$  can be used to determine the invertibility of  $\mathbf{A}$ . If  $\lambda(\mathbf{A}) = \{\lambda_1, \dots, \lambda_n\}$ , then the determinant of  $\mathbf{A}$  is equal to

$$\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i.$$

$\mathbf{A}$  is nonsingular if and only if  $\det(\mathbf{A}) \neq 0$ .

A useful concept for eigenvalues solvers is the *Rayleigh quotient*:

DEFINITION 2.5 (Rayleigh Quotient). Let  $\mathbf{0} \neq \mathbf{z} \in \mathbb{C}^n$ . The *Rayleigh quotient* of  $\mathbf{A}$  and  $\mathbf{z}$  is defined by

$$R_{\mathbf{A}}(\mathbf{z}) = \frac{\mathbf{z}^* \mathbf{A} \mathbf{z}}{\mathbf{z}^* \mathbf{z}}.$$

Note that if  $\mathbf{z}$  is an eigenvector of  $\mathbf{A}$ , then the Rayleigh quotient is the corresponding eigenvalue.

Another way to express the eigenvalue problem in Definition 2.2 is that  $(\lambda, \mathbf{x})$  is an eigenpair of  $\mathbf{A}$  if and only if  $\mathbf{0} \neq \mathbf{x} \in \ker(\mathbf{A} - \lambda \mathbf{I})$ . Based on this kernel, we can define the *eigenspace* of  $\mathbf{A}$ :

DEFINITION 2.6 (Eigenspace).

$$\mathcal{E}_{\lambda}(\mathbf{A}) = \ker(\mathbf{A} - \lambda \mathbf{I})$$

is the *eigenspace* of  $\mathbf{A}$  corresponding to  $\lambda$ .

The following concept of *invariant subspaces* bears similarities to the eigenvalue problem.

DEFINITION 2.7 (Invariant Subspace). A subspace  $\mathcal{S} \subset \mathbb{C}^n$  is said to be *invariant* under a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  ( $\mathbf{A}$ -invariant) if  $\mathbf{A}\mathcal{S} \subset \mathcal{S}$ .

DEFINITION 2.8. Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$ ,  $\mathcal{S} \subset \mathbb{C}^n$ , and  $\mathbf{S} \in \mathbb{C}^{n \times k}$  with  $k = \text{rank}(\mathbf{S}) = \dim(\mathcal{S}) \leq n$  and  $\mathcal{S} = \text{ran}(\mathbf{S})$ . Then,  $\mathcal{S}$  is  $\mathbf{A}$ -invariant if and only if there exists a matrix  $\mathbf{B} \in \mathbb{C}^{k \times k}$  such that

$$\mathbf{A}\mathbf{S} = \mathbf{S}\mathbf{B}.$$

REMARK 2.9. Using Definition 2.8, it is easy to show the following relations:

- If  $(\lambda, \mathbf{x})$  is an eigenpair of  $\mathbf{B}$ , then  $(\lambda, \mathbf{S}\mathbf{x})$  is an eigenpair of  $\mathbf{A}$ . Hence,  $\lambda(\mathbf{B}) \subset \lambda(\mathbf{A})$ .
- If  $k = n$ , then  $\mathbf{S}$  is invertible and hence

$$\mathbf{A} = \mathbf{S}\mathbf{B}\mathbf{S}^{-1}.$$

This means  $\mathbf{A}$  and  $\mathbf{B}$  are *similar* and  $\lambda(\mathbf{B}) = \lambda(\mathbf{A})$ . This concept is introduced next.

Most of the presented algorithms will transform a matrix  $\mathbf{A}$  into simpler forms, such as diagonal or triangular matrices, in order to simplify the original eigenvalue problem. Transformations that preserve the eigenvalues of matrices are called *similarity transformations*.

DEFINITION 2.10 (Similarity Transformation). Two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  are said to be *similar* if there exists a nonsingular matrix  $\mathbf{C} \in \mathbb{C}^{n \times n}$  such that

$$\mathbf{A} = \mathbf{C}\mathbf{B}\mathbf{C}^{-1}.$$

The mapping  $\mathbf{B} \rightarrow \mathbf{A}$  is called a *similarity transformation*. The similarity of  $\mathbf{A}$  and  $\mathbf{B}$  implies that they have the same eigenvalues. If  $(\lambda, \mathbf{x})$  is an eigenpair of  $\mathbf{B}$ , then  $(\lambda, \mathbf{C}\mathbf{x})$  is an eigenpair of  $\mathbf{A}$ .

The simplest form to which a matrix can be transformed is a diagonal matrix. But as we will see, this is not always possible.

DEFINITION 2.11 (Diagonalizability). If  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is similar to a diagonal matrix, then  $\mathbf{A}$  is said to be *diagonalizable*.

A similarity transformation in which  $\mathbf{C}$  is orthogonal (or unitary), i.e.,  $\mathbf{C}^T\mathbf{C} = \mathbf{I}$  (or  $\mathbf{C}^*\mathbf{C} = \mathbf{I}$ ), is called *orthogonal (or unitary) similarity transformation*. Unitary/orthogonal similarity transformations play a key role in numerical computations since  $\|\mathbf{C}\|_2 = 1$ . Considering the calculation of similarity transformations, it can be shown (cf. [22, Chap. 7.1.5]) that the roundoff error  $\mathbf{E}$  satisfies

$$\|\mathbf{E}\| \approx \epsilon_{\text{machine}} \kappa_2(\mathbf{C}) \|\mathbf{A}\|_2.$$

Here,  $\epsilon_{\text{machine}}$  is the *machine precision*<sup>2</sup> and  $\kappa_2(\mathbf{C})$  the condition number of  $\mathbf{C}$  with respect to the 2-norm. In particular,  $\kappa_2(\mathbf{C})$  is the error gain. Therefore, if the similarity transformation is unitary, we get

$$\|\mathbf{E}\| \approx \epsilon_{\text{machine}} \|\mathbf{A}\|_2$$

and hence no amplification of error.

THEOREM 2.12 (Unitary Diagonalizability; see [22, Cor. 7.1.4]).  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is *unitarily diagonalizable if and only if it is normal* ( $\mathbf{A}^*\mathbf{A} = \mathbf{A}\mathbf{A}^*$ ).

Now, let us show the connection between a similarity transformation of a matrix  $\mathbf{A}$  and its eigenpairs: It follows from Definition 2.4 and the Fundamental Theorem of Algebra that  $\mathbf{A}$  has  $n$  (not necessarily distinct) eigenvalues. If we denote the  $n$  eigenpairs by  $(\lambda_1, \mathbf{x}_1), \dots, (\lambda_n, \mathbf{x}_n)$ , i.e.  $\mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{x}_i$  for  $i = 1, \dots, n$ , we can write

$$(2.1) \quad \mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{\Lambda},$$

where  $\mathbf{\Lambda} = \text{diag}(\lambda_i)_{i=1, \dots, n} \in \mathbb{C}^{n \times n}$  is a diagonal matrix containing the eigenvalues, and  $\mathbf{X} = [\mathbf{x}_1 | \dots | \mathbf{x}_n] \in \mathbb{C}^{n \times n}$  is a matrix whose columns are formed by the eigenvectors. This looks almost as a similarity transformation. In fact, the “only” additional ingredient we need is the invertibility of  $\mathbf{X}$ . Under the assumption that  $\mathbf{X}$  is nonsingular, we obtain  $\mathbf{X}^{-1}\mathbf{A}\mathbf{X} = \mathbf{\Lambda}$ , and hence,  $\mathbf{A}$  and  $\mathbf{\Lambda}$  are similar. But when can we expect of  $\mathbf{X}$  to be nonsingular? To discuss this, we introduce some terminology:

DEFINITION 2.13 (Multiplicity). Let  $\lambda$  be an eigenvalue of  $\mathbf{A}$ .

- $\lambda$  has *algebraic multiplicity*  $m^a$ , if it is a root of multiplicity  $m^a$  of the characteristic polynomial  $p_{\mathbf{A}}$ .
- If  $m^a = 1$ , then  $\lambda$  is called *simple*. Otherwise,  $\lambda$  is said to be *multiple*.

<sup>2</sup>The machine precision is  $\epsilon_{\text{machine}} = 2^{-53} \approx 1.11 \cdot 10^{-16}$  in the double precision IEEE floating point format and  $\epsilon_{\text{machine}} = 2^{-24} \approx 5.96 \cdot 10^{-6}$  in the single precision IEEE floating point format. For more details, we refer to, e.g., [37, 66, 27].

- The *geometric multiplicity*  $m^g$  of  $\lambda$  is defined as the dimension of the associated eigenspace, i.e.,  $m^g = \dim(\mathcal{E}_\lambda(\mathbf{A}))$ . It is the maximum number of independent eigenvectors associated with  $\lambda$ .
- It holds  $m^g \leq m^a$ .
- If  $m^g < m^a$ , then  $\lambda$  and  $\mathbf{A}$  are called *defective* or *non-diagonalizable*.

Note that if all eigenvalues of  $\mathbf{A}$  are simple, then they are distinct. Now, we can state a result about the nonsingularity of the eigenvector matrix  $\mathbf{X}$  in (2.1):

**THEOREM 2.14** (Diagonal Form; see [22, Cor. 7.1.8]). *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with eigenvalues  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$ .  $\mathbf{A}$  is nondefective if and only if there exists a nonsingular matrix  $\mathbf{X} \in \mathbb{C}^{n \times n}$  such that*

$$\mathbf{X}^{-1} \mathbf{A} \mathbf{X} = \text{diag}(\lambda_i)_{i=1, \dots, n}.$$

The similarity transformation given in Theorem 2.14 transforms  $\mathbf{A}$  into a diagonal matrix whose entries reveal the eigenvalues of  $\mathbf{A}$ .

We have seen that a similarity transformation to a diagonal matrix is not always possible. Before we come to the next similarity transformation, we introduce the concept of *deflation* – the process of breaking down an eigenvalue problem into smaller eigenvalue problems.

**THEOREM 2.15** (See [22, Lemma 7.1.3]). *Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$ ,  $\mathbf{S} \in \mathbb{C}^{n \times k}$  with  $\text{rank}(\mathbf{S}) = k < n$  and  $\mathbf{B} \in \mathbb{C}^{k \times k}$  such that*

$$\mathbf{A} \mathbf{S} = \mathbf{S} \mathbf{B},$$

*i.e.,  $\text{ran}(\mathbf{S})$  is an  $\mathbf{A}$ -invariant subspace. Then, there exists a unitary  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  such that*

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{T} = \begin{bmatrix} \mathbf{T}_{11} & \mathbf{T}_{12} \\ \mathbf{0} & \mathbf{T}_{22} \end{bmatrix}$$

*and*

$$\begin{aligned} \lambda(\mathbf{T}) &= \lambda(\mathbf{T}_{11}) \cup \lambda(\mathbf{T}_{22}), \\ \lambda(\mathbf{T}_{11}) &= \lambda(\mathbf{A}) \cap \lambda(\mathbf{B}) \end{aligned}$$

*with  $\mathbf{T}_{11} \in \mathbb{C}^{k \times k}$ .*

From Theorem 2.15, we obtain a similarity transformation that transforms a matrix  $\mathbf{A}$  into an upper triangular matrix whose diagonal entries reveal the eigenvalues of  $\mathbf{A}$ . Such a decomposition always exists.

**THEOREM 2.16** (Schur Decomposition; see [22, Theor. 7.1.3]). *Given  $\mathbf{A} \in \mathbb{C}^{n \times n}$  with eigenvalues  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$ . Then, there exists a unitary matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  such that*

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{T} = \mathbf{D} + \mathbf{N},$$

*where  $\mathbf{D} = \text{diag}(\lambda_i)_{i=1, \dots, n}$ , and  $\mathbf{N} \in \mathbb{C}^{n \times n}$  is strictly upper triangular. Moreover,  $\mathbf{Q}$  can be chosen such that the eigenvalues  $\lambda_i$  appear in any order in  $\mathbf{D}$ .*

The transformation in Theorem 2.16 deals with a complex matrix  $\mathbf{Q}$  even when  $\mathbf{A}$  is real. A slight variation of the Schur decomposition shows that complex arithmetic can be avoided in this case. This is based on the fact that complex eigenvalues

always occur in complex conjugate pairs, i.e., if  $(\lambda, \mathbf{x})$  is an eigenpair of  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , then  $(\bar{\lambda}, \bar{\mathbf{x}})$  is an eigenpair of  $\mathbf{A}$ .

**THEOREM 2.17** (Real Schur Decomposition; see [22, Theor. 7.4.1]). *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with eigenvalues  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$ . Then, there exists an orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that*

$$\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{T} = \begin{bmatrix} \mathbf{T}_{1,1} & \cdots & \mathbf{T}_{1,m} \\ & \ddots & \vdots \\ & & \mathbf{T}_{m,m} \end{bmatrix},$$

where  $\mathbf{T} \in \mathbb{R}^{n \times n}$  is quasi-upper triangular. The diagonal blocks  $\mathbf{T}_{i,i}$  are either  $1 \times 1$  or  $2 \times 2$  matrices. A  $1 \times 1$  block corresponds to a real eigenvalue  $\lambda_j \in \mathbb{R}$ . A  $2 \times 2$  block corresponds to a pair of complex conjugate eigenvalues. For a complex conjugate eigenvalue pair  $\lambda_k = \mu + \nu$ ,  $\lambda_l = \mu - \nu$ ,  $\mathbf{T}_{i,i}$  has the form

$$\mathbf{T}_{i,i} = \begin{bmatrix} \mu & \nu \\ -\nu & \mu \end{bmatrix}.$$

Moreover,  $\mathbf{Q}$  can be chosen such that the diagonal blocks  $\mathbf{T}_{i,i}$  appear in any order in  $\mathbf{T}$ .

The next similarity transformation we present transforms a matrix  $\mathbf{A}$  into *upper Hessenberg* form. Such a decomposition always exists and will play an important role in eigenvalue solvers for nonsymmetric matrices.

**THEOREM 2.18** (Hessenberg Decomposition). *Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Then, there exists a unitary matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  such that*

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{H} = \begin{bmatrix} h_{1,1} & h_{1,2} & h_{1,3} & \cdots & h_{1,n} \\ h_{2,1} & h_{2,2} & h_{2,3} & \cdots & h_{2,n} \\ 0 & h_{3,2} & h_{3,3} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & h_{n-1,n} \\ 0 & \cdots & 0 & h_{n,n-1} & h_{n,n} \end{bmatrix}.$$

$\mathbf{H}$  is called an *upper Hessenberg matrix*. Further,  $\mathbf{H}$  is said to be *unreduced* if  $h_{j+1,j} \neq 0$  for all  $j = 1, \dots, n-1$ .

From a theoretical point of view, one of the most important similarity transformations is the *Jordan decomposition*, or *Jordan Canonical Form*.

**THEOREM 2.19** (Jordan Decomposition; see [22, Theor. 7.1.9]). *Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$  with exactly  $p$  distinct eigenvalues  $\lambda_1, \dots, \lambda_p \in \mathbb{C}$  for  $p \leq n$ . Then, there exists a nonsingular matrix  $\mathbf{X} \in \mathbb{C}^{n \times n}$  such that*

$$\mathbf{X}^{-1} \mathbf{A} \mathbf{X} = \begin{bmatrix} \mathbf{J}_1(\lambda_1) & & \\ & \ddots & \\ & & \mathbf{J}_p(\lambda_p) \end{bmatrix}.$$

Each block  $\mathbf{J}_i(\lambda_i)$  has the block diagonal structure

$$\mathbf{J}_i(\lambda_i) = \begin{bmatrix} \mathbf{J}_{i,1}(\lambda_i) & & \\ & \ddots & \\ & & \mathbf{J}_{i,m_i^g}(\lambda_i) \end{bmatrix} \in \mathbb{C}^{m_i^a \times m_i^a}$$

with

$$\mathbf{J}_{i,k}(\lambda_i) = \begin{bmatrix} \lambda_i & 1 & & & \\ & \ddots & \ddots & & \\ & & \lambda_i & 1 & \\ & & & \lambda_i & \\ & & & & \lambda_i \end{bmatrix} \in \mathbb{C}^{m_{i,k} \times m_{i,k}},$$

where  $m_i^a$  and  $m_i^g$  are the algebraic and geometric multiplicity of the eigenvalue  $\lambda_i$ . Each of the subblocks  $\mathbf{J}_{i,k}(\lambda_i)$  is referred to as a Jordan block.

Unfortunately, from a computational point of view, the computation of the Jordan Canonical Form is numerically unstable.

An important and practical factorization is the *QR decomposition*:

DEFINITION 2.20 (QR Decomposition; see [22, Theor. 5.2.1]). Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Then, there exists an orthogonal  $\mathbf{Q} \in \mathbb{R}^{m \times m}$  and an upper triangular  $\mathbf{R} \in \mathbb{R}^{m \times n}$  such that

$$\mathbf{A} = \mathbf{QR}.$$

This concludes the theoretical part of the background study. Next, we are getting started with computational aspects.

**2.2. First computational aspects.** This section quickly reviews aspects of perturbation theory and illustrates possible difficulties in computing eigenvalues accurately. This is followed by a brief overview of different classes of methods for solving eigenvalue problems. Details about all mentioned methods are given in the upcoming sections.

First of all, it should be clear that in general we must iterate to find eigenvalues of a matrix: According to Definition 2.4, the eigenvalues of a matrix  $\mathbf{A}$  are the roots of the characteristic polynomial  $p_{\mathbf{A}}(x)$ . In 1824, Abel proved that for polynomials of degree  $n \geq 5$ , there is no formula for its roots in terms of its coefficients that uses only the operations of addition, subtraction, multiplication, division, and taking  $k$ th roots. Hence, even if we could work in exact arithmetic, no computer would produce the exact roots of an arbitrarily polynomial in a finite number of steps. (This is different than direct methods for solving systems of linear equations such as Gaussian elimination.) Hence, computing the eigenvalues of any  $n \times n$  matrix  $\mathbf{A}$  requires an iterative process if  $n \geq 5$ .

As already indicated in the previous section, many methods are based on repeatedly performing *similarity transformations* to bring  $\mathbf{A}$  into a simpler equivalent form. This typically means generating as many zero entries in the matrix as possible. The goal is eventually to perform a Schur decomposition. If the matrix is normal, then the Schur decomposition simplifies to a diagonal matrix (not only an upper triangular matrix), and this has implications in terms of stability of numerical computations. As part of the process, we often aim to reduce the matrix into tridiagonal form (symmetric case) or upper Hessenberg form (nonsymmetric case). *Deflation*, *projection*, and other tools can be incorporated and are extremely valuable.

Now, let us focus on the reduction of a matrix  $\mathbf{A}$  to upper Hessenberg form. One way to accomplish this is the use of *Householder reflectors* (also called *Householder*



*transformations*). They can be used to zero out selected components of a vector. Hence, by performing a sequence of Householder reflections on the columns of  $\mathbf{A}$ , we can transform  $\mathbf{A}$  into a simpler form. Householder reflectors are matrices of the form

$$\mathbf{P} = \mathbf{I} - \frac{2}{\mathbf{v}^* \mathbf{v}} \mathbf{v} \mathbf{v}^*,$$

where  $\mathbf{v} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$ . Householder matrices are Hermitian ( $\mathbf{P} = \mathbf{P}^*$ ), unitary, and numerically stable. Geometrically,  $\mathbf{P}$  applied to a vector  $\mathbf{x}$  reflects it about the hyperplane  $\text{span}\{\mathbf{v}\}^\perp$ .

Assume we want to bring  $\mathbf{A}$  into upper Hessenberg form. Then, the first step is to introduce zeros into all except the first two entries of the first column of  $\mathbf{A}$ . Let us denote by  $\mathbf{x} = [a_{2,1}, \dots, a_{n,1}]^T$  the part of the first column of  $\mathbf{A}$  under consideration. We are looking for a vector  $\mathbf{v} \in \mathbb{C}^{n-1} \setminus \{\mathbf{0}\}$  such that  $\mathbf{P}\mathbf{x}$  results in a multiple of the first unit vector  $\mathbf{e}_1$ . This can be achieved with the ansatz  $\mathbf{v} = \mathbf{x} \pm \|\mathbf{x}\|_2 \mathbf{e}_1$ , since this yields

$$\mathbf{P}\mathbf{x} = \mp \|\mathbf{x}\|_2 \mathbf{e}_1.$$

Let us illustrate the action of  $\mathbf{P}$  to the first column of  $\mathbf{A}$ :

$$\begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ a_{3,1} & a_{3,2} & \cdots & a_{3,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,n} \end{bmatrix} \rightarrow \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ \mp \|\mathbf{x}\|_2 & a_{2,2} & \cdots & a_{2,n} \\ 0 & a_{3,2} & \cdots & a_{3,n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n,2} & \cdots & a_{n,n} \end{bmatrix}.$$

Note that the first step is not complete yet: Remember that, for a similarity transformation, we need to apply the Householder matrix twice, i.e.,  $\mathbf{P}^* \mathbf{A} \mathbf{P}$ . Note that the right multiplication with  $\mathbf{P}$  does not destroy the zeros:

$$\begin{bmatrix} * & * & \cdots & * \\ * & * & \cdots & * \\ * & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ * & * & \cdots & * \end{bmatrix} \xrightarrow{\mathbf{P}^*} \begin{bmatrix} * & * & \cdots & * \\ * & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & * & \cdots & * \end{bmatrix} \xrightarrow{\mathbf{P}} \begin{bmatrix} * & * & \cdots & * \\ * & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & * & \cdots & * \end{bmatrix}.$$

$\mathbf{A} \qquad \qquad \mathbf{P}^* \mathbf{A} \qquad \qquad \mathbf{P}^* \mathbf{A} \mathbf{P}$

Let us denote the Householder matrix in the first step by  $\mathbf{P}_1$ . The above procedure is repeated with the Householder matrix  $\mathbf{P}_2$  to the second column of  $\mathbf{P}_1^* \mathbf{A} \mathbf{P}_1$ , then to the third column of  $\mathbf{P}_2^* \mathbf{P}_1^* \mathbf{A} \mathbf{P}_1 \mathbf{P}_2$  with the Householder matrix  $\mathbf{P}_3$ , and so on, until we end up with a matrix in upper Hessenberg form as given in Definition 2.18. Let us denote the Householder matrix in step  $i$  by  $\mathbf{P}_i$ . After  $n - 2$  steps, we obtain the upper Hessenberg form:

$$\underbrace{\mathbf{P}_{n-2}^* \cdots \mathbf{P}_1^*}_{\mathbf{P}^*} \mathbf{A} \underbrace{\mathbf{P}_1 \cdots \mathbf{P}_{n-2}}_{\mathbf{P}} = \mathbf{H} = \begin{bmatrix} * & * & * & \cdots & * \\ * & * & * & \cdots & * \\ 0 & * & * & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & * \\ 0 & \cdots & 0 & * & * \end{bmatrix}.$$

REMARK 2.21. Here are a few additional comments about the process:

- In practice, one chooses  $\mathbf{v} = \mathbf{x} + \text{sign}(x_1)\|\mathbf{x}\|_2\mathbf{e}_1$ , where  $x_1$  is the first entry of the vector  $\mathbf{x}$  under consideration.
- Note that in each step  $i$ , the corresponding vector  $\mathbf{x}_i$ , and hence  $\mathbf{v}_i$  and  $\mathbf{P}_i$ , shrink by one in size.
- The reduction of an  $n \times n$  matrix to upper Hessenberg form via Householder reflections requires  $\mathcal{O}(n^3)$  operations.

One may ask why do we first bring  $\mathbf{A}$  to upper Hessenberg form and not immediately to triangular form using Householder reflections? In that case, the right multiplication with  $\mathbf{P}$  would destroy the zeros previously introduced:

$$\begin{array}{ccc} \begin{bmatrix} * & * & \cdots & * \\ * & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ * & * & \cdots & * \end{bmatrix} & \xrightarrow{\mathbf{P}^*} & \begin{bmatrix} * & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & * & \cdots & * \end{bmatrix} & \xrightarrow{\mathbf{P}} & \begin{bmatrix} * & * & \cdots & * \\ * & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ * & * & \cdots & * \end{bmatrix} \\ \mathbf{A} & & \mathbf{P}^*\mathbf{A} & & \mathbf{P}^*\mathbf{A}\mathbf{P} \end{array}.$$

This should not come as a surprise: we already knew from Abel (1824) that it is impossible to obtain a Schur form of  $\mathbf{A}$  in a finite number of steps; see the beginning of this Section.

REMARK 2.22. If  $\mathbf{A}$  is symmetric, the reduction to upper Hessenberg form turns into a tridiagonal matrix. That is because the right multiplication with  $\mathbf{P}_i$  also introduces zeros above the diagonal:

$$\begin{array}{ccc} \begin{bmatrix} * & * & * & \cdots & * \\ * & * & * & \cdots & * \\ * & * & * & \cdots & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ * & * & * & \cdots & * \end{bmatrix} & \xrightarrow{\mathbf{P}_1^*} & \begin{bmatrix} * & * & * & \cdots & * \\ * & * & * & \cdots & * \\ 0 & * & * & \cdots & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & * & * & \cdots & * \end{bmatrix} & \xrightarrow{\mathbf{P}_1} & \begin{bmatrix} * & * & 0 & \cdots & 0 \\ * & * & * & \cdots & * \\ 0 & * & * & \cdots & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & * & * & \cdots & * \end{bmatrix} \\ \mathbf{A} & & \mathbf{P}_1^*\mathbf{A} & & \mathbf{P}_1^*\mathbf{A}\mathbf{P}_1 \end{array}.$$

Later on, we will discuss fast algorithms for eigenvalue problems with symmetric tridiagonal matrices.

REMARK 2.23. The Hessenberg reduction via Householder reflections is backward stable, i.e., there exists a small perturbation  $\delta\mathbf{A}$  of  $\mathbf{A}$  such that

$$\hat{\mathbf{H}} = \hat{\mathbf{P}}^*(\mathbf{A} + \delta\mathbf{A})\hat{\mathbf{P}}, \quad \|\delta\mathbf{A}\|_F \leq cn^2\epsilon_{\text{machine}}\|\mathbf{A}\|_F.$$

Here,  $\hat{\mathbf{H}}$  is the computed upper Hessenberg matrix,  $\hat{\mathbf{P}} = \hat{\mathbf{P}}_1 \cdots \hat{\mathbf{P}}_{n-2}$  is a product of exactly unitary Householder matrices based on computed vectors  $\hat{\mathbf{v}}_i$ , and  $c > 0$  a constant. For more details, we refer to [66, p. 351] and [27, Sec. 19.3].

During the next sections, we will see how the upper Hessenberg form (or the tridiagonal form in case of symmetric matrices) is used within eigenvalue solvers.

Before we talk about algorithms, we need to understand when it is difficult to compute eigenvalues accurately. The following example shows that eigenvalues of a matrix are continuous (but not necessarily differentiable) functions of it.

EXAMPLE 2.24. Consider the perturbed Jordan block

$$\mathbf{A}(\varepsilon) = \begin{bmatrix} 0 & 1 & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ \varepsilon & & & 1 & \\ & & & & 0 \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

The characteristic polynomial is given as  $p_{\mathbf{A}(\varepsilon)}(x) = (-1)^n(x^n - \varepsilon)$ . Hence, the eigenvalues are  $\lambda_j(\varepsilon) = \varepsilon^{\frac{1}{n}} \exp(\frac{2\ell j\pi}{n})$  for  $j = 1, \dots, n$ . None of the eigenvalues is differentiable at  $\varepsilon = 0$ . Their rate of change at the origin is infinite. Consider for instance the case  $n = 20$  and  $\varepsilon = 10^{-16}$  (machine precision), then  $\lambda_1(\varepsilon) = 0.1507 + 0.0490i$  whereas  $\lambda(0) = 0$ .

Let us quickly address the issue of estimating the quality of computed eigenvalues. The question here is: How do eigenvalues and eigenvectors vary when the original matrix undergoes small perturbations? We start with considering the sensitivity of simple eigenvalues.

THEOREM 2.25 (See, e.g. [22, Chap. 7.2.2]). *Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$  with a simple eigenvalue  $\lambda$ , a right (unit norm) eigenvector  $\mathbf{x}$ , and a left (unit norm) eigenvector  $\mathbf{y}$ . Let  $\mathbf{A} + \delta\mathbf{A}$  be a perturbation of  $\mathbf{A}$  and  $\lambda + \delta\lambda$  the corresponding perturbed eigenvalue. Then*

$$\delta\lambda = \frac{\mathbf{y}^* \delta\mathbf{A} \mathbf{x}}{\mathbf{y}^* \mathbf{x}} + \mathcal{O}(\|\delta\mathbf{A}\|_2^2).$$

The condition number of  $\lambda$  is defined as  $s(\lambda) = \frac{1}{|\mathbf{y}^* \mathbf{x}|}$ . It can be shown that

$$s(\lambda) = \frac{1}{\cos(\theta(\mathbf{x}, \mathbf{y}))},$$

where  $\theta(\mathbf{x}, \mathbf{y})$  is the angle between  $\mathbf{x}$  and  $\mathbf{y}$ .

In general,  $\mathcal{O}(\varepsilon)$  perturbations in  $\mathbf{A}$  can induce  $\frac{\varepsilon}{s(\lambda)}$  changes in an eigenvalue. Thus, if  $s(\lambda)$  is small, then  $\lambda$  is ill-conditioned, and  $\mathbf{A}$  is “close to” a matrix with multiple eigenvalues. If  $\mathbf{A}$  is normal, then every simple eigenvalue satisfies  $s(\lambda) = 1$ , which means that these eigenvalues are well-conditioned. In the case of a multiple eigenvalue  $\lambda$ ,  $s(\lambda)$  is not unique anymore. For a defective eigenvalue  $\lambda$ , it holds in general that  $\mathcal{O}(\varepsilon)$  perturbations in  $\mathbf{A}$  can result in  $\mathcal{O}(\varepsilon^{\frac{1}{p}})$  changes in  $\lambda$ , where  $p$  denotes the size of the largest Jordan block associated with  $\lambda$ . This is the effect we have observed in Example 2.24:  $\mathbf{A}(0)$  has the effective eigenvalue zero with algebraic multiplicity  $n$  and geometric multiplicity one. Hence,  $\mathcal{O}(10^{-16})$  perturbations in  $\mathbf{A}$  can result in  $\mathcal{O}(10^{-\frac{16}{20}}) = \mathcal{O}(0.1585)$  changes in the eigenvalue. In other words, small perturbations in the input data caused a large perturbation in the output. This can lead to numerical instabilities of eigenvalue solvers.

In the following, we start with algorithms for computing a few up to all eigenvalues for small to moderate-sized matrices. Then, we continue with large and sparse matrices.

### 3. Small to moderate-sized matrices

In general, as previously stated, we may separate methods into ones that are based on matrix decompositions vs. ones that are based on matrix-vector products. The power method, which we start with in the sequel, is an important building block for both classes of methods. It is based on matrix-vector products, but it is invaluable for eigensolvers based on decompositions. We choose to include it in this section, noting that it is relevant also for eigensolvers for large and sparse matrices.

**3.1. Power method.** The power method is one of the oldest techniques for solving eigenvalue problems. It is used for computing a *dominant eigenpair*, i.e., the eigenvalue of maximum modulus of a matrix  $\mathbf{A}$  and a corresponding eigenvector. The algorithm consists of generating a sequence of matrix-vector multiplications  $\{\mathbf{A}^k \mathbf{v}_0\}_{k=0,1,\dots}$ , where  $\mathbf{v}_0$  is some nonzero initial vector.

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with  $\mathbf{A}\mathbf{x}_j = \lambda_j \mathbf{x}_j$  for  $j = 1, \dots, n$ . Assume that the eigenvectors  $\mathbf{x}_j, j = 1, \dots, n$ , are linearly independent, i.e.,  $\mathbf{A}$  is nondefective. Given  $\mathbf{0} \neq \mathbf{v}_0 \in \mathbb{C}^n$ , we can expand it using the eigenvectors of  $\mathbf{A}$  to

$$\mathbf{v}_0 = \sum_{j=1}^n \beta_j \mathbf{x}_j,$$

where  $\beta_j \in \mathbb{C}$  for  $j = 1, \dots, n$ . Applying  $\mathbf{A}$  to  $\mathbf{v}_0$  yields

$$\mathbf{A}\mathbf{v}_0 = \sum_{j=1}^n \beta_j \mathbf{A}\mathbf{x}_j = \sum_{j=1}^n \beta_j \lambda_j \mathbf{x}_j.$$

Hence, the eigenvectors corresponding to eigenvalues of larger modulus are favored. The above procedure can be repeated. In fact, for any  $k \in \mathbb{N}$ , we have

$$\mathbf{A}^k \mathbf{v}_0 = \sum_{j=1}^n \beta_j \mathbf{A}^k \mathbf{x}_j = \sum_{j=1}^n \beta_j \lambda_j^k \mathbf{x}_j.$$

In order for the following algorithm to converge, we need the following assumptions:  $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$ . Since we are interested in the eigenvalue of maximum modulus, we need some distance to the remaining eigenvalues.  $\lambda_1$  is called the *dominant eigenvalue*. We further need  $\mathbf{v}_0$  to have a component in the direction of the eigenvector corresponding to  $\lambda_1$ , i.e.,  $\beta_1 \neq 0$ . Note that this assumption is less concerning in practice since rounding errors during the iteration typically introduce components in the direction of  $\mathbf{x}_1$ . However, we need it for the following theoretical study. Due to  $\beta_1 \neq 0$ , we can write

$$\mathbf{A}^k \mathbf{v}_0 = \beta_1 \lambda_1^k \mathbf{x}_1 + \sum_{j=2}^n \beta_j \lambda_j^k \mathbf{x}_j = \beta_1 \lambda_1^k \left( \mathbf{x}_1 + \sum_{j=2}^n \frac{\beta_j}{\beta_1} \left( \frac{\lambda_j}{\lambda_1} \right)^k \mathbf{x}_j \right).$$

Since  $\lambda_1$  is a dominant eigenvalue, we get  $\left( \frac{\lambda_j}{\lambda_1} \right)^k \xrightarrow{k \rightarrow \infty} 0$  for all  $j = 2, \dots, n$ . Hence, it can be shown (cf. [47, Theor. 4.1]) that  $\mathbf{A}^k \mathbf{v}_0$ , as well as the scaled version  $\mathbf{v}_k = \frac{\mathbf{A}^k \mathbf{v}_0}{\|\mathbf{A}^k \mathbf{v}_0\|_2}$  which is used in practice to avoid overflow/underflow, converges linearly to a multiple of  $\mathbf{x}_1$  with a convergence rate proportional to  $\frac{|\lambda_2|}{|\lambda_1|}$ . A value

of  $\frac{|\lambda_2|}{|\lambda_1|} \approx 1$  indicates a slow convergence behavior. Algorithm 3.1 shows the power method. The approximated eigenvalue in step  $k$

$$\lambda_1^{(k)} = \mathbf{v}_k^T \mathbf{A} \mathbf{v}_k$$

is computed using the *Rayleigh quotient*; see Definition 2.5. This is based on the following: Given a vector  $\hat{\mathbf{x}}_1$  that approximates the eigenvector  $\mathbf{x}_1$ . Then,  $\hat{\lambda}_1 = \hat{\mathbf{x}}_1^T \mathbf{A} \hat{\mathbf{x}}_1$  is the best eigenvalue approximation in the least-squares sense, i.e.,

$$(3.1) \quad \hat{\lambda}_1 = \arg \min_{\mu} \|\mathbf{A} \hat{\mathbf{x}}_1 - \mu \hat{\mathbf{x}}_1\|_2^2.$$

We can solve this minimization problem by solving the *normal equation*

$$\begin{aligned} \hat{\mathbf{x}}_1^T \hat{\mathbf{x}}_1 \mu &= \hat{\mathbf{x}}_1^T \mathbf{A} \hat{\mathbf{x}}_1 \\ \Leftrightarrow \mu &= \frac{\hat{\mathbf{x}}_1^T \mathbf{A} \hat{\mathbf{x}}_1}{\hat{\mathbf{x}}_1^T \hat{\mathbf{x}}_1}; \end{aligned}$$

see, e.g., [46, Chap. 5.3.3]. Since we normalize the computed eigenvectors in the power method, i.e.,  $\|\hat{\mathbf{x}}_1\|_2 = 1$ , we get the desired result. The cost for  $k$  iterations

---

**Algorithm 3.1:** Power method

---

```

1 Choose  $\mathbf{v}_0 = \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$ 
2 for  $k = 1, 2, \dots$ , until termination do
3      $\tilde{\mathbf{v}} = \mathbf{A} \mathbf{v}_{k-1}$ 
4      $\mathbf{v}_k = \frac{\tilde{\mathbf{v}}}{\|\tilde{\mathbf{v}}\|_2}$ 
5      $\lambda_1^{(k)} = \mathbf{v}_k^T \mathbf{A} \mathbf{v}_k$ 
6 end
```

---

is  $\mathcal{O}(2kn^2)$  floating point operations (flops).

The power method can be applied to large, sparse, or implicit matrices. It is simple and basic but can be slow. We assumed for the convergence that  $\mathbf{A}$  is nondefective. For the case of a defective  $\mathbf{A}$ , the power method can still be applied but converges even more slowly; see, e.g., [28]. Moreover, we want to emphasize again that the power method only works if the matrix under consideration has *one* dominant eigenvalue. This excludes the case of, e.g., a dominant complex eigenvalue<sup>3</sup> or of dominant eigenvalues of opposite signs. The power method is rather used as a building block for other, more robust and general algorithms. We refer to [66, Chap. 10] for a detailed discussion of the power method. Next, we discuss a method that overcomes the mentioned difficulties.

**3.2. Inverse power method.** We have seen that the power method is in general slow. Moreover, it is good only for one well-separated dominant eigenvalue. How can we accelerate it, and what about the more general case of looking for a non-dominant eigenpair? The inverse power method uses *shift and invert* techniques to overcome these limitations of the power method. It aims to compute the eigenvalue of  $\mathbf{A}$  that is closest to a certain scalar (shift) and a corresponding eigenvector. It also enhances the convergence behavior. The price for these improvements is the

---

<sup>3</sup>As noted in Section 2.1, eigenvalues of real matrices always occur in complex conjugate pairs.

solution of a linear system in each iteration.

The idea is the following: Assume  $\mathbf{A} \in \mathbb{R}^{n \times n}$  has eigenpairs  $(\lambda_j, \mathbf{x}_j)_{j=1, \dots, n}$  with  $|\lambda_1| \geq \dots \geq |\lambda_n|$ . Let  $\alpha \in \mathbb{R}$  with  $\alpha \neq \lambda_j$  for  $j = 1, \dots, n$ . This will be the *shift* in the inverse power method. In practice, we choose  $\alpha \approx \lambda_i$  for some  $i$  depending on which (real) eigenvalue  $\lambda_i$  we want to find. Hence, in order for the method to work, we need to know approximately the value of the eigenvalue we are interested in. Then,  $\mathbf{A} - \alpha \mathbf{I}$  has eigenpairs  $(\lambda_j - \alpha, \mathbf{x}_j)_{j=1, \dots, n}$ , and  $(\mathbf{A} - \alpha \mathbf{I})^{-1}$  has eigenpairs  $(\mu_j, \mathbf{x}_j)_{j=1, \dots, n}$  with  $\mu_j = (\lambda_j - \alpha)^{-1}$ . Let  $\lambda_i$  and  $\lambda_j$  be the two eigenvalues that are closest to  $\alpha$  with  $|\lambda_i - \alpha| < |\lambda_j - \alpha|$ . Then, the two largest eigenvalues  $\mu_1$  and  $\mu_2$  of  $(\mathbf{A} - \alpha \mathbf{I})^{-1}$  are

$$\mu_1 = \frac{1}{\lambda_i - \alpha}, \quad \mu_2 = \frac{1}{\lambda_j - \alpha}.$$

Hence, the power method applied to  $(\mathbf{A} - \alpha \mathbf{I})^{-1}$  converges to  $\mu_1$  and an eigenvector of  $\mu_1$  with convergence rate

$$\frac{|\mu_2|}{|\mu_1|} = \frac{\frac{1}{|\lambda_j - \alpha|}}{\frac{1}{|\lambda_i - \alpha|}} = \frac{|\lambda_i - \alpha|}{|\lambda_j - \alpha|}.$$

We know from the previous section that we need a small value of  $\frac{|\mu_2|}{|\mu_1|}$  in order to converge fast. Hence, we desire  $|\lambda_i - \alpha| \ll |\lambda_j - \alpha|$ , which requires a “good” choice of the shift  $\alpha$ . If we are interested for instance in the dominant eigenvalue, estimations based on norms of  $\mathbf{A}$  can be used; see, e.g., [22, Chap. 2.3.2].

---

**Algorithm 3.2:** Inverse power method

---

- 1 Choose  $\mathbf{v}_0 = \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$
  - 2 **for**  $k = 1, 2, \dots$ , *until termination* **do**
  - 3     Solve  $(\mathbf{A} - \alpha \mathbf{I})\tilde{\mathbf{v}} = \mathbf{v}_{k-1}$
  - 4      $\mathbf{v}_k = \frac{\tilde{\mathbf{v}}}{\|\tilde{\mathbf{v}}\|_2}$
  - 5      $\lambda^{(k)} = \mathbf{v}_k^T \mathbf{A} \mathbf{v}_k$
  - 6 **end**
- 

As already mentioned at the beginning of this section, the price for overcoming difficulties of the power method by using a shift and invert approach is the solution of a linear system in every iteration. If  $\alpha$  is fixed, then we have to solve linear systems with one matrix and many right-hand sides: If a direct method can be applied, then we form an LU decomposition of  $\mathbf{A} - \alpha \mathbf{I}$  once. The cost for solving the two triangular systems arising from the LU decomposition is  $\mathcal{O}(n^2)$ . For huge problems, iterative methods have to be employed to solve the linear systems. This pays off only if the inverse iteration converges very fast.

In summary, we have seen that we can apply the inverse power method to find different eigenvalues using different shifts. During the whole iteration, the inverse power method uses a fixed shift  $\alpha$ . The next method involves a dynamic shift  $\alpha_k$ .

**3.3. Rayleigh quotient iteration.** The idea of the Rayleigh quotient iteration is to learn the shift as the iteration proceeds using the calculated eigenvalue

$\lambda^{(k-1)}$  from the previous step  $k-1$ . Now, each iteration is potentially more expensive since the linear systems involve different matrices in each step. However, the new algorithm may converge in many fewer iterations. In the Hermitian case, we potentially obtain a cubic convergence rate; see, e.g., [39].

Note that the matrix  $(\mathbf{A} - \lambda^{(k-1)}\mathbf{I})$  may be singular. This is the case when the shift hits an eigenvalue of  $\mathbf{A}$ . The cost for solving the linear system with  $(\mathbf{A} - \lambda^{(k-1)}\mathbf{I})$  is  $\mathcal{O}(n^3)$  if  $\mathbf{A}$  is full. For an upper Hessenberg matrix, it reduces to  $\mathcal{O}(n^2)$ , and for a tridiagonal matrix even to  $\mathcal{O}(n)$ .

Next, we discuss a technique that uses information of a computed dominant eigenpair for the approximation of a second-dominant eigenpair.

**3.4. Deflation.** Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  have eigenvalues  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$ , right eigenvectors  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , and left eigenvectors  $\mathbf{y}_1, \dots, \mathbf{y}_n$ . Note that  $(\lambda_1, \mathbf{x}_1)$  is a dominant eigenpair. Suppose we have approximated the dominant eigenvector  $\mathbf{x}_1$  of  $\mathbf{A}$  by  $\hat{\mathbf{x}}_1$  with  $\|\hat{\mathbf{x}}_1\|_2 = 1$ . Now, we are interested in approximating the next eigenvalue  $\lambda_2$ .

*Deflation* is based on a simple rank-one modification of  $\mathbf{A}$ , as follows: Compute  $\mathbf{A}_1 = \mathbf{A} - \alpha \hat{\mathbf{x}}_1 \mathbf{w}^T$ , where  $\alpha \in \mathbb{R}$  is an appropriate shift and  $\mathbf{w} \in \mathbb{R}^n$  an arbitrary vector such that  $\mathbf{w}^T \hat{\mathbf{x}}_1 = 1$ .

**THEOREM 3.1** (Wielandt; see [65]). *In the case  $\hat{\mathbf{x}}_1 = \mathbf{x}_1$ , the eigenvalues of  $\mathbf{A}_1$  are  $\lambda_1 - \alpha, \lambda_2, \dots, \lambda_n$ . Moreover, the right eigenvector  $\mathbf{x}_1$  and the left eigenvectors  $\mathbf{y}_2, \dots, \mathbf{y}_n$  are preserved.*

**PROOF.**

$$(\mathbf{A} - \alpha \mathbf{x}_1 \mathbf{w}^T) \mathbf{x}_1 = \mathbf{A} \mathbf{x}_1 - \alpha \mathbf{x}_1 \mathbf{w}^T \mathbf{x}_1 = \lambda_1 \mathbf{x}_1 - \alpha \mathbf{x}_1$$

since  $\mathbf{w}^T \mathbf{x}_1 = 1$ . For  $i = 2, \dots, n$ , we have

$$\mathbf{y}_i^* (\mathbf{A} - \alpha \mathbf{x}_1 \mathbf{w}^T) = \mathbf{y}_i^* \mathbf{A} - \alpha \mathbf{y}_i^* \mathbf{x}_1 \mathbf{w}^T = \mathbf{y}_i^* \mathbf{A} = \lambda_i \mathbf{y}_i^*$$

since  $\mathbf{y}_i^* \mathbf{x}_1 = 0$  for  $i = 2, \dots, n$ . □

Hence, a modification of  $\mathbf{A}$  to  $\mathbf{A}_1 = \mathbf{A} - \alpha \hat{\mathbf{x}}_1 \mathbf{w}^T$  displaces the dominant eigenvalue of  $\mathbf{A}$ . The rank-one modification should be chosen such that  $\lambda_2$  becomes the dominant eigenvalue of  $\mathbf{A}_1$ . We can then proceed for instance with the power method applied to  $\mathbf{A}_1$  in order to obtain an approximation of  $\lambda_2$ . This technique is called *Wielandt deflation*.

There are many ways to choose  $\mathbf{w}$ . A simple choice (due to Hotelling [29]) is to choose  $\mathbf{w} = \mathbf{y}_1$  the first left eigenvector (or an approximation of it) or  $\mathbf{w} = \mathbf{x}_1$  or rather its approximation  $\mathbf{w} = \hat{\mathbf{x}}_1$ . It can be shown (cf. [47, Chap. 4.2.2]) that if  $\mathbf{x}_1$  is nearly orthogonal to  $\mathbf{x}_2$  or if  $\frac{\lambda_1 - \lambda_2}{\alpha} \ll 1$ , the choice  $\mathbf{w} = \mathbf{x}_1$  is nearly optimal in terms of eigenvalue conditioning.

Note that we never need to form the matrix  $\mathbf{A}_1$  explicitly. This is important since  $\mathbf{A}_1$  is a dense matrix. For calculating the matrix-vector product  $\mathbf{y} = \mathbf{A}_1 \mathbf{x}$ , we just need to perform  $\mathbf{y} \leftarrow \mathbf{A} \mathbf{x}$ ,  $\beta = \alpha \mathbf{w}^T \mathbf{x}$ , and  $\mathbf{y} \leftarrow \mathbf{y} - \beta \hat{\mathbf{x}}_1$ . We can apply this procedure recursively without difficulty. However, keep in mind that, for a long

deflation process, errors accumulate.

So far, the discussed algorithms compute only one eigenpair at once, i.e., a one-dimensional invariant subspace. Next, we consider another generalization of the power and inverse power method that can be used to compute higher-dimensional invariant subspaces.

**3.5. Orthogonal iteration.** Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  have eigenpairs  $(\lambda_j, \mathbf{x}_j)_{j=1, \dots, n}$  with  $|\lambda_1| \geq \dots \geq |\lambda_n|$ . From the real Schur decomposition in Theorem 2.17, we know there exists an orthogonal  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that

$$\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{T},$$

where the diagonal blocks of  $\mathbf{T}$  correspond to the eigenvalues of  $\mathbf{A}$  in real form. Assume that the eigenvalues  $\lambda_i$ , represented in  $\mathbf{T}$  in real form, are ordered from  $\lambda_1$  to  $\lambda_n$ . Let  $1 \leq r < n$ . Then, we can do the following partitioning:

$$\mathbf{Q} = [\mathbf{Q}^{(r)}, \mathbf{Q}^{(n-r)}], \quad \mathbf{T} = \begin{bmatrix} \mathbf{T}^{(r,r)} & \mathbf{T}^{(r,n-r)} \\ \mathbf{0} & \mathbf{T}^{(n-r,n-r)} \end{bmatrix},$$

where  $\mathbf{Q}^{(r)} \in \mathbb{R}^{r \times r}$  and  $\mathbf{T}^{(r,r)} \in \mathbb{R}^{r \times r}$ . Note that  $r$  should be chosen such that the  $(r+1, r)$  entry in  $\mathbf{T}$  is zero, i.e., we do not split a complex conjugate eigenpair to  $\mathbf{T}^{(r,r)}$  and  $\mathbf{T}^{(n-r,n-r)}$ . Then,

$$\mathbf{A} \mathbf{Q}^{(r)} = \mathbf{Q}^{(r)} \mathbf{T}^{(r,r)},$$

i.e.,  $\text{ran}(\mathbf{Q}^{(r)})$  is an  $\mathbf{A}$ -invariant subspace corresponding to the  $r$  largest (in modulus) eigenvalues. Due to this property, this subspace is also called *dominant*.

Now, we are interested in computing such a dominant  $r$ -dimensional invariant subspace. Hence, instead of dealing with a matrix-vector product as in the algorithms before, we go over to a matrix-matrix product, i.e., we apply  $\mathbf{A}$  to a few vectors simultaneously. This can be achieved by the orthogonal iteration presented in Algorithm 3.3.

---

**Algorithm 3.3:** Orthogonal iteration

---

- 1 Choose  $\mathbf{Q}_0 \in \mathbb{R}^{n \times r}$  with orthonormal columns
  - 2 **for**  $k = 1, 2, \dots$ , *until termination* **do**
  - 3      $\mathbf{Z}_k = \mathbf{A} \mathbf{Q}_{k-1}$
  - 4      $\mathbf{Q}_k \mathbf{R}_k = \mathbf{Z}_k$      (QR factorization)
  - 5 **end**
- 

The QR factorization in Line 4 refers to the QR decomposition in Definition 2.20. It can be computed by, e.g., the modified Gram–Schmidt algorithm in  $\mathcal{O}(2nr^2)$  flops [22, Chap. 5.2.8], Householder reflections in  $\mathcal{O}(2r^2(n - \frac{r}{3}))$  flops [22, Chap. 5.2.2], or Givens transformations in  $\mathcal{O}(3r^2(n - \frac{r}{3}))$  flops [22, Chap. 5.2.5]. Note that the complexity can be reduced if the corresponding matrix is of upper Hessenberg form. We will discuss this further below. Note that Line 3 and 4 yield

$$\mathbf{A} \mathbf{Q}_{k-1} = \mathbf{Q}_k \mathbf{R}_k,$$



where  $\mathbf{R}_k$  is upper triangular. Now, if  $|\lambda_r| > |\lambda_{r+1}|$  and  $\mathbf{Q}_0$  has components in the desired eigendirections, then we have

$$\text{ran}(\mathbf{Q}_k) \xrightarrow{k \rightarrow \infty} \text{ran}(\mathbf{Q}^{(r)})$$

with a convergence rate proportional to  $\frac{|\lambda_{r+1}|}{|\lambda_r|}$ . For more details, we refer to [22, Chap. 7.3.2].

REMARK 3.2. By replacing the QR factorization in Line 4 of Algorithm 3.3 with  $\mathbf{Q}_k = \mathbf{Z}_k$ , we obtain the *subspace iteration*, also called *simultaneous iteration*. Under the same conditions as before, it holds

$$\text{ran}(\mathbf{Z}_k) \xrightarrow{k \rightarrow \infty} \text{ran}(\mathbf{Q}^{(r)}).$$

However, the columns of  $\mathbf{Z}_k$  form an increasingly ill-conditioned basis for  $\mathbf{A}^k \text{ran}(\mathbf{Q}_0)$  since each column of  $\mathbf{Z}_k$  converges to a multiple of the dominant eigenvector. The orthogonal iteration overcomes this difficulty by orthonormalizing the columns of  $\mathbf{Z}_k$  at each step.

From the orthogonal iteration, we can derive the QR iteration, a method for finding all eigenvalues of a matrix  $\mathbf{A}$ .

**3.6. QR iteration.** We obtain the QR iteration from the orthogonal iteration if we set  $r = n$ , i.e., we want to compute all eigenvalues, and  $\mathbf{Q}_0 = \mathbf{I}$ .

---

**Algorithm 3.4:** Prelude to QR iteration

---

```

1 Choose  $\mathbf{Q}_0 = \mathbf{I}$  (orthogonal)
2 for  $k = 1, 2, \dots$ , until termination do
3    $\mathbf{Z}_k = \mathbf{A}\mathbf{Q}_{k-1}$ 
4    $\mathbf{Q}_k \mathbf{R}_k = \mathbf{Z}_k$  (QR factorization)
5    $\mathbf{A}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k$ 
6 end
```

---

We can rewrite Algorithm 3.4 by using the following equivalence:

$$\begin{aligned}
 (3.2) \quad \mathbf{A}_{k-1} &\stackrel{\text{Line 5}}{=} \mathbf{Q}_{k-1}^T \mathbf{A} \mathbf{Q}_{k-1} \stackrel{\text{Line 3}}{=} \mathbf{Q}_{k-1}^T \mathbf{Z}_k \stackrel{\text{Line 4}}{=} \mathbf{Q}_{k-1}^T \mathbf{Q}_k \mathbf{R}_k =: \tilde{\mathbf{Q}}_k \mathbf{R}_k, \\
 \mathbf{A}_k &\stackrel{\text{Line 5}}{=} \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_{k-1} \mathbf{Q}_{k-1}^T \mathbf{Q}_k \stackrel{\text{Line 3}}{=} \mathbf{Q}_k^T \mathbf{Z}_k \mathbf{Q}_{k-1}^T \mathbf{Q}_k \\
 &\stackrel{\text{Line 4}}{=} \mathbf{R}_k \mathbf{Q}_{k-1}^T \mathbf{Q}_k \stackrel{(3.2)}{=} \mathbf{R}_k \tilde{\mathbf{Q}}_k.
 \end{aligned}$$

Note that the product of two orthogonal matrices is orthogonal. Hence,  $\tilde{\mathbf{Q}}_k$  is orthogonal. Therefore,  $\mathbf{A}_k$  is determined by a QR decomposition of  $\mathbf{A}_{k-1}$ . This form of the QR iteration is presented in Algorithm 3.5. Note that  $\mathbf{Q}_0$  does not have to be the identity matrix.

From Line 3 and 4 of Algorithm 3.5 we get

$$\mathbf{A}_k = (\mathbf{Q}_0 \cdots \mathbf{Q}_k)^T \mathbf{A} (\mathbf{Q}_0 \cdots \mathbf{Q}_k) =: \hat{\mathbf{Q}}_k^T \mathbf{A} \hat{\mathbf{Q}}_k,$$

where  $\hat{\mathbf{Q}}_k$  is orthogonal. Hence,  $\text{ran}(\hat{\mathbf{Q}}_k)$  is an  $\mathbf{A}$ -invariant subspace and  $\lambda(\mathbf{A}) = \lambda(\mathbf{A}_k)$ . If  $|\lambda_1| > \dots > |\lambda_n|$  and  $\mathbf{Q}_0$  has components in the desired eigendirections, then we have

$$\text{ran}(\hat{\mathbf{Q}}_k(:, 1:l)) \xrightarrow{k \rightarrow \infty} \text{ran}(\mathbf{Q}(:, 1:l)) \quad \forall 1 \leq l \leq n$$

---

**Algorithm 3.5:** QR iteration

---

```

1  $\mathbf{A}_0 = \mathbf{Q}_0^T \mathbf{A} \mathbf{Q}_0$  (real Schur form,  $\mathbf{Q}_0 \in \mathbb{R}^{n \times n}$  orthogonal)
2 for  $k = 1, 2, \dots$ , until termination do
3    $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_{k-1}$  (QR factorization)
4    $\mathbf{A}_k = \mathbf{R}_k \mathbf{Q}_k$ 
5 end

```

---

with a convergence rate proportional to  $\frac{|\lambda_{l+1}|}{|\lambda_l|}$ . Hence,  $\mathbf{A}_k \xrightarrow{k \rightarrow \infty} \mathbf{T}$ , where  $\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{T}$  is a real Schur decomposition of  $\mathbf{A}$ . For more details, we refer to [22, Chap. 7.3.3].

Overall, the QR iteration computes the Schur form of a matrix. As in the previous section, we considered the real Schur form here. But note that if we allow complex arithmetic, we get the same results with a (complex) Schur form. For further readings on the QR iteration we refer to, e.g., [62, 55, 32].

As already mentioned in the previous section, in this form the cost of each step of the QR iteration is  $\mathcal{O}(n^3)$ . But we can reduce the complexity if we start with  $\mathbf{A}_0$  in upper Hessenberg form. Moreover, we can speed up the convergence using shifts.

**3.7. QR iteration with shifts.** If we choose  $\mathbf{Q}_0$  such that  $\mathbf{A}_0$  is in upper Hessenberg form, the cost of each step of the QR iteration reduces to  $\mathcal{O}(n^2)$ . If  $\mathbf{A}$  is symmetric, then the cost per step is  $\mathcal{O}(n)$ . It can be shown that each  $\mathbf{A}_k$  is upper Hessenberg. This is the first modification. Second, shifts  $\zeta_k \in \mathbb{R}$  are introduced in order to accelerate the deflation process (see Theorem 2.15). Deflation occurs every time  $\mathbf{A}_k$  is reduced, i.e., at least one of its subdiagonal entries is zero. In such a case, we continue with two smaller subproblems. The matrices  $\mathbf{A}_{k-1}$  and  $\mathbf{A}_k$  in

---

**Algorithm 3.6:** QR iteration with shifts

---

```

1  $\mathbf{A}_0 = \mathbf{Q}_0^T \mathbf{A} \mathbf{Q}_0$  upper Hessenberg form (Tridiagonal if  $\mathbf{A}$  is symmetric)
2 for  $k = 1, 2, \dots$ , until termination do
3    $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_{k-1} - \zeta_k \mathbf{I}$  (QR factorization)
4    $\mathbf{A}_k = \mathbf{R}_k \mathbf{Q}_k + \zeta_k \mathbf{I}$ 
5 end

```

---

Algorithm 3.6 are orthogonally similar since

$$\mathbf{A}_k = \mathbf{R}_k \mathbf{Q}_k + \zeta_k \mathbf{I} = \mathbf{Q}_k^T (\mathbf{Q}_k \mathbf{R}_k + \zeta_k \mathbf{I}) \mathbf{Q}_k = \mathbf{Q}_k^T \mathbf{A}_{k-1} \mathbf{Q}_k,$$

and  $\mathbf{Q}_k$  from the QR decomposition is orthogonal.

Why does the shift strategy work? If  $\zeta_k$  is an eigenvalue of the unreduced Hessenberg matrix  $\mathbf{A}_{k-1}$ , then  $\mathbf{A}_{k-1} - \zeta_k \mathbf{I}$  is singular. This implies  $\mathbf{R}_k$  is singular, where the  $(n, n)$  entry of  $\mathbf{R}_k$  is zero. Then, the last row of the upper Hessenberg matrix  $\mathbf{A}_k = \mathbf{R}_k \mathbf{Q}_k + \zeta_k \mathbf{I}$  consists of zeros except for the  $(n, n)$  entry which is  $\zeta_k$ . So we have converged to the form

$$\mathbf{A}_k = \begin{bmatrix} \mathbf{A}' & \mathbf{a} \\ \mathbf{0}^T & \zeta_k \end{bmatrix},$$

and can now work on a smaller matrix (deflate) and continue the QR iteration. We can accept the  $(n, n)$  entry as  $\zeta_k$  as it is presumably a good approximation to the eigenvalue. In summary, we obtain deflation after one step in exact arithmetic if we shift by an exact eigenvalue.

If  $\zeta = \zeta_k$  for all  $k = 1, 2, \dots$ , and we order the eigenvalues  $\lambda_i$  of  $\mathbf{A}$  such that

$$|\lambda_1 - \zeta| \geq \dots \geq |\lambda_n - \zeta|,$$

then the  $p$ th subdiagonal entry in  $\mathbf{A}_k$  converges to zero with rate  $\frac{|\lambda_{p+1} - \zeta|}{|\lambda_p - \zeta|}$ . Of course, we need  $|\lambda_{p+1} - \zeta| < |\lambda_p - \zeta|$  in order to get any convergence result.

In practice, deflation occurs whenever a subdiagonal entry  $a_{p+1,p}^{(k)}$  of  $\mathbf{A}_k$  is small enough, e.g., if

$$|a_{p+1,p}^{(k)}| \leq c\epsilon_{\text{machine}}(|a_{p,p}^{(k)}| + |a_{p+1,p+1}^{(k)}|)$$

for a small constant  $c > 0$ .

Let us quickly summarize some shift strategies: The single-shift strategy uses  $\zeta_k = a_{n,n}^{(k-1)}$ . It can be shown (cf. [22, Chap. 7.5.3]) that the convergence  $a_{n,n-1}^{(k)} \xrightarrow{k \rightarrow \infty} 0$  is even quadratic. When we deal with complex eigenvalues, then  $\zeta_k = a_{n,n}^{(k-1)}$  tends to be a poor approximation. Then, the double-shift strategy is preferred which performs two single-shift steps in succession, i.e., Lines 3–4 in Algorithm 3.6 are repeated a second time with a second shift. Using implicit QR factorizations, one double-shift step can be implemented with  $\mathcal{O}(n^2)$  flops ( $\mathcal{O}(n)$  flops in the symmetric case); see e.g., [22, Chap. 7.5.5]. This technique was first described by Francis [18, 19] and refers to a *Francis QR step*.

The overall QR algorithm requires  $\mathcal{O}(n^3)$  flops. For more details about the QR iteration, we refer to, e.g., [42, 33, 35, 60, 63, 64]. Further readings concerning shift strategies include [17, 15, 61].

We know that the Hessenberg reduction of a symmetric matrix leads to a tridiagonal matrix. In the following, we review methods for this special case.

**3.8. Algorithms for symmetric (tridiagonal) matrices.** Before we consider eigenvalue problems for the special case of symmetric tridiagonal matrices, we review one of the oldest methods for symmetric matrices  $\mathbf{A}$  — *Jacobi's method*. For general symmetric eigenvalue problems, we refer the reader to [11, Chap. 5] — it contains important theoretic concepts, e.g., gaps of eigenvalues and the related perturbation theory, and gives a nice overview of direct eigenvalue solvers.

3.8.1. *Jacobi's method.* Jacobi's method is one of the oldest algorithms [30] for eigenvalue problems with a cost of  $\mathcal{O}(cn^3)$  flops with a large constant  $c$ . However, it is still of current interest due to its parallelizability and accuracy [12].

The method is based on a sequence of orthogonal similarity transformations

$$(3.3) \quad \dots \mathbf{Q}_3^T \mathbf{Q}_2^T \mathbf{Q}_1^T \mathbf{A} \mathbf{Q}_1 \mathbf{Q}_2 \mathbf{Q}_3 \dots,$$



i.e.,  $\mathbf{A}_{k+1}$  moves closer to diagonal form with each Jacobi step. In order to maximize the reduction in (3.4),  $p$  and  $q$  should be chosen such that  $|a_{p,q}^{(k)}|$  is maximal. With this choice, we get after  $k$  Jacobi steps (cf. [11, Theor. 5.11])

$$\text{off}(\mathbf{A}_k)^2 \leq \left(1 - \frac{2}{n(n-1)}\right)^k \text{off}(\mathbf{A}_0)^2,$$

i.e., convergence at a linear rate. This scheme is the original version from Jacobi in 1846 and is referred to as *classical Jacobi algorithm*. It even can be shown that the asymptotic convergence rate is quadratic (cf. [11, Theor. 5.12]); see [49, 58]. While the cost for an update is  $\mathcal{O}(n)$  flops, the search for the optimal  $(p, q)$  costs  $\mathcal{O}(n^2)$  flops. For a simpler method, we refer the reader to the *cyclic Jacobi method*; see e.g. [66, p. 270] or [22, Chap. 8.5.3]. In general, the cost of the cyclic Jacobi method is considerably higher than the cost of the symmetric QR iteration. However, it is easily parallelizable.

Again we want to emphasize that we do not need to tridiagonalize in Jacobi's method. In the following, we discuss two methods that need to start by reducing a symmetric matrix  $\mathbf{A}$  to tridiagonal form: *bisection* and *divide-and-conquer*. Another method is *MR<sup>3</sup>* or *MRRR* (Algorithm of Multiple Relatively Robust Representations) [13] — a sophisticated variant of the inverse iteration, which is efficient when eigenvalues are close to each other.

3.8.2. *Bisection.* For the rest of Section 3, we consider eigenvalue problems for symmetric tridiagonal matrices of the form

$$(3.5) \quad \mathbf{A} = \begin{bmatrix} a_1 & b_1 & & & & \\ & b_1 & a_2 & & & \\ & & b_2 & a_3 & \ddots & \\ & & & \ddots & \ddots & b_{n-1} \\ & & & & b_{n-1} & a_n \end{bmatrix}.$$

We begin with bisection, a method that can be used to find a subset of eigenvalues, e.g., the largest/smallest eigenvalue or eigenvalues within an interval. Let  $\mathbf{A}^{(k)} = \mathbf{A}(1 : k, 1 : k)$  be the leading  $k \times k$  principal submatrix of  $\mathbf{A}$  with characteristic polynomial

$$p^{(k)}(x) = \det(\mathbf{A}^{(k)} - x\mathbf{I})$$

for  $k = 1, \dots, n$ . If  $b_i \neq 0$  for  $i = 1, \dots, n-1$ , then

$$\det(\mathbf{A}^{(k)}) = a_k \det(\mathbf{A}^{(k-1)}) - b_{k-1}^2 \det(\mathbf{A}^{(k-2)}),$$

which yields

$$p^{(k)}(x) = (a_k - x)p^{(k-1)}(x) - b_{k-1}^2 p^{(k-2)}(x)$$

with  $p^{(-1)}(x) = 0$  and  $p^{(0)}(x) = 1$ . Hence,  $p^{(n)}(x)$  can be evaluated in  $\mathcal{O}(n)$  flops. Given  $y < z \in \mathbb{R}$  with  $p^{(n)}(y)p^{(n)}(z) < 0$  (Hence, there exists a  $w \in (y, z)$  with  $p^{(n)}(w) = 0$ ), we can use the *method of bisection* (see, e.g., [22, Chap. 8.4.1]) to find an approximate root of  $p^{(n)}(x)$  and hence an approximate eigenvalue of  $\mathbf{A}$ . The method of bisection converges linearly in the sense that the error is approximately halved at each step.

Assume that the eigenvalues  $\lambda_j(\mathbf{A}^{(k)})$  of  $\mathbf{A}^{(k)}$  are ordered as

$$\lambda_1(\mathbf{A}^{(k)}) \geq \dots \geq \lambda_k(\mathbf{A}^{(k)}).$$

For computing, e.g.,  $\lambda_k(\mathbf{A})$  for a given  $k$  or the largest eigenvalue that is smaller than a given  $\mu \in \mathbb{R}$ , then we need the following theorem:

**THEOREM 3.3** (Sturm Sequence Property; see [22, Theor. 8.4.1]). *If  $\mathbf{A}$  is unreduced, i.e.,  $b_i \neq 0$  for  $i = 1, \dots, n-1$ , then the eigenvalues of  $\mathbf{A}^{(k-1)}$  strictly separate the eigenvalues of  $\mathbf{A}^{(k)}$ :*

$$\begin{aligned} \lambda_k(\mathbf{A}^{(k)}) &< \lambda_{k-1}(\mathbf{A}^{(k-1)}) < \lambda_{k-1}(\mathbf{A}^{(k)}) \\ &< \lambda_{k-2}(\mathbf{A}^{(k-1)}) < \lambda_{k-2}(\mathbf{A}^{(k)}) < \dots \\ &< \lambda_2(\mathbf{A}^{(k)}) < \lambda_1(\mathbf{A}^{(k-1)}) < \lambda_1(\mathbf{A}^{(k)}). \end{aligned}$$

Moreover, if  $a(\mu)$  denotes the number of sign changes in the sequence

$$\{p^{(0)}(\mu), p^{(1)}(\mu), \dots, p^{(n)}(\mu)\},$$

where  $p^{(k)}(\mu)$  has the opposite sign from  $p^{(k-1)}(\mu)$  if  $p^{(k)}(\mu) = 0$ , then  $a(\mu)$  equals the number of  $\mathbf{A}$ 's eigenvalues that are less than  $\mu$ .

In order to find an initial interval for the method of bisection, we make use of the following simplified version of the *Gershgorin theorem*:

**THEOREM 3.4** (Gershgorin). *If  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric, then*

$$\lambda(\mathbf{A}) \subseteq \cup_{i=1}^n [a_{i,i} - r_i, a_{i,i} + r_i],$$

where  $r_i = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{i,j}|$ .

The more general version of the Gershgorin theorem can be found, e.g., in [22, Theor. 8.1.3]. Suppose we want to compute  $\lambda_k(\mathbf{A})$ , where  $\mathbf{A}$  is symmetric tridiagonal as in (3.5). Then, from Theorem 3.4, we get  $\lambda_k(\mathbf{A}) \in [y, z]$  with

$$y = \min_{1 \leq i \leq n} a_i - |b_i| - |b_{i-1}|, \quad z = \max_{1 \leq i \leq n} a_i + |b_i| + |b_{i-1}|$$

and  $b_0 = b_n = 0$ . Hence, with this choice of  $y$  and  $z$ , we can reformulate the method of bisection to converge to  $\lambda_k(\mathbf{A})$ ; see, e.g., [22, Chap. 8.4.2]. Another version of this scheme can be used to compute subsets of eigenvalues of  $\mathbf{A}$ ; see [3]. For a variant that computes specific eigenvalues, we refer to [39, p. 46].

The cost of bisection is  $\mathcal{O}(nk)$  flops, where  $k$  is the number of desired eigenvalues. Hence, it can be much faster than the QR iteration if  $k \ll n$ . Once the desired eigenvalues are found, we can use the *inverse power method* (Section 3.2) to find the corresponding eigenvectors. The inverse power method costs in the best case (well-separated eigenvalues)  $\mathcal{O}(nk)$  flops. In the worst case (many clustered eigenvalues), the cost is  $\mathcal{O}(nk^2)$  flops and the accuracy of the computed eigenvectors is not guaranteed. Next, we review a method that is better suited for finding all (or most) eigenvalues and eigenvectors, especially when the eigenvalues may be clustered.

3.8.3. *Divide-and-conquer.* The idea of divide-and-conquer is to recursively divide the eigenvalue problem into smaller subproblems until we reach matrices of dimension one, for which the eigenvalue problem is trivial. The method was first introduced in 1981 [9] while its parallel version was developed in 1987 [14].

The starting point is to write the symmetric tridiagonal matrix  $\mathbf{A}$  in (3.5) as a sum of a block diagonal matrix of two tridiagonal matrices  $\mathbf{T}_1$  and  $\mathbf{T}_2$ , plus a rank-1 correction:

$$\mathbf{A} = \begin{bmatrix} \mathbf{T}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_2 \end{bmatrix} + b_m \mathbf{v} \mathbf{v}^T,$$

where  $\mathbf{v} \in \mathbb{R}^n$  is a column vector whose  $m$ th and  $(m+1)$ st entry is equal to one ( $1 \leq m \leq n-1$ ) and all remaining entries are zero. Suppose we have the real Schur decompositions of  $\mathbf{T}_1$  and  $\mathbf{T}_2$ , i.e.,  $\mathbf{T}_i = \mathbf{Q}_i \mathbf{D}_i \mathbf{Q}_i^T$  for  $i = 1, 2$  with  $\mathbf{Q}_1 \in \mathbb{R}^{m \times m}$  and  $\mathbf{Q}_2 \in \mathbb{R}^{(n-m) \times (n-m)}$  orthogonal.

$$\mathbf{A} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2 \end{bmatrix} \left( \begin{bmatrix} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_2 \end{bmatrix} + b_m \mathbf{u} \mathbf{u}^T \right) \begin{bmatrix} \mathbf{Q}_1^T & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2^T \end{bmatrix},$$

where

$$\mathbf{u} = \begin{bmatrix} \mathbf{Q}_1^T & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2^T \end{bmatrix} \mathbf{v} = \begin{bmatrix} \text{last column of } \mathbf{Q}_1^T \\ \text{first column of } \mathbf{Q}_2^T \end{bmatrix}.$$

Hence,  $\lambda(\mathbf{A}) = \lambda(\mathbf{D} + b_m \mathbf{u} \mathbf{u}^T)$  where  $\mathbf{D} = \begin{bmatrix} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_2 \end{bmatrix} = \text{diag}(d_i)_{i=1, \dots, n}$  is a diagonal matrix. Hence, the problem now reduces to finding the eigenvalues of a diagonal plus a rank-1 matrix. This can be further simplified to finding the eigenvalues of the identity matrix plus a rank-1 matrix, using simply the characteristic polynomial: In particular, under the assumption that  $\mathbf{D} - \lambda \mathbf{I}$  is nonsingular, and using

$$\det(\mathbf{D} + b_m \mathbf{u} \mathbf{u}^T - \lambda \mathbf{I}) = \det(\mathbf{D} - \lambda \mathbf{I}) \det(\mathbf{I} + b_m (\mathbf{D} - \lambda \mathbf{I})^{-1} \mathbf{u} \mathbf{u}^T),$$

we obtain

$$\lambda \in \lambda(\mathbf{A}) \Leftrightarrow \det(\mathbf{I} + b_m (\mathbf{D} - \lambda \mathbf{I})^{-1} \mathbf{u} \mathbf{u}^T) = 0.$$

The matrix  $\mathbf{I} + b_m (\mathbf{D} - \lambda \mathbf{I})^{-1} \mathbf{u} \mathbf{u}^T$  is of special structure, and its determinant can be computed using the following lemma:

LEMMA 3.5 (See [11, Lemma 5.1]). *Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . Then*

$$\det(\mathbf{I} + \mathbf{x} \mathbf{y}^T) = 1 + \mathbf{y}^T \mathbf{x}.$$

In our case, we get

$$\det(\mathbf{I} + b_m (\mathbf{D} - \lambda \mathbf{I})^{-1} \mathbf{u} \mathbf{u}^T) = 1 + b_m \sum_{i=1}^n \frac{u_i^2}{d_i - \lambda} \equiv f(\lambda),$$

and the eigenvalues of  $\mathbf{A}$  are the roots of the *secular equation*  $f(\lambda) = 0$ . This can be solved, e.g., by Newton's method, which converges in practice in a bounded number of steps per eigenvalue. Note that solving the secular equation needs caution due to possible deflations ( $d_i = d_{i+1}$  or  $u_i = 0$ ) or small values of  $u_i$ . For more details on the function  $f$  and on solving the secular equation, we refer to [11, Chap. 5.3.3]. The cost for computing all eigenvalues is  $\mathcal{O}(n^2 \log(n))$  flops by using the above strategy of recursively dividing the eigenvalue problem into smaller subproblems.

Hence, the pure eigenvalue computation (without eigenvectors) is more expensive than in the QR iteration. However, the eigenvectors can be computed more cheaply:

LEMMA 3.6 (See [11, Lemma 5.2]). *If  $\lambda \in \lambda(\mathbf{D} + b_m \mathbf{u}\mathbf{u}^T)$ , then  $(\mathbf{D} - \lambda\mathbf{I})^{-1}\mathbf{u}$  is a corresponding eigenvector.*

The cost for computing all eigenvectors is  $\mathcal{O}(n^2)$  flops. However, the eigenvector computation from Lemma 3.6 is not numerically stable. If  $\lambda$  is too close to a diagonal entry  $d_i$ , we obtain large roundoff errors since we divide by  $d_i - \lambda$ . If two eigenvalues  $\lambda_i$  and  $\lambda_j$  are very close, the orthogonality of the computed eigenvectors can get lost. For a numerical stable computation, we refer to [24], which requires in practice  $\mathcal{O}(cn^3)$  flops where  $c \ll 1$ .

Here, we finish the discussion of eigenvalue problems for small to moderate-sized matrices. We now move to discuss problems where the matrix is large and sparse.

#### 4. Large and sparse matrices

In this section,  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is considered to be large and sparse. Sparse matrices are matrices with very few nonzero entries. Sparse often means that there are  $\mathcal{O}(1)$  nonzero entries per row. We note that matrices that are not necessarily sparse but give rise to very fast matrix-vector products (for example, via the Fast Fourier Transform) often also allow for applying the methods discussed in this section.

The meaning of “large matrices” is relative. Let us say we consider matrices of size millions. In order to take advantage of the large number of zero entries, special storage schemes are required; see, e.g., [47, Chap. 2]. We will assume that it is not easy to form a matrix decomposition such as the QR factorization. In particular, similarity transformations would destroy the sparsity. Hence, we will mainly rely on matrix-vector products, which are often computable in  $\mathcal{O}(n)$  flops instead of  $\mathcal{O}(n^2)$ .

In this chapter, we review methods for computing a few eigenpairs of  $\mathbf{A}$ . In fact, in practice, one often needs the  $k$  smallest/largest eigenvalues or the  $k$  eigenvalues closest to  $\mu \in \mathbb{C}$  for a small  $k$  and their corresponding eigenvectors. In the following, we introduce *orthogonal projection methods*, from which we can derive the state-of-the-art *Krylov methods*, which make use of cheap matrix-vector products. Note that projection methods even play a role for the methods discussed in Section 3.

**4.1. Orthogonal projection methods.** Suppose we want to find an approximation  $(\hat{\lambda}, \hat{\mathbf{x}})$  of an eigenpair  $(\lambda, \mathbf{x})$  of  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . The idea of *projection techniques* is to extract  $\hat{\mathbf{x}}$  from some subspace  $\mathcal{K}$ . This is called the *subspace of approximants* or the *right subspace*. The uniqueness of  $\hat{\mathbf{x}}$  is typically realized via the imposition of orthogonality conditions. We denote by

$$\mathbf{r} = \mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}$$

the residual vector. It is a measure for the quality of the approximate eigenpair  $(\hat{\lambda}, \hat{\mathbf{x}})$ . The orthogonality conditions consist of constraining the residual  $\mathbf{r}$  to be orthogonal to some subspace  $\mathcal{L}$ , i.e.,

$$\mathbf{r} \perp \mathcal{L}.$$



$\mathcal{L}$  is called the *left subspace*. This framework is commonly known as the *Petrov-Galerkin conditions* in diverse areas of mathematics, e.g., the finite element method. The case  $\mathcal{L} = \mathcal{K}$  leads to the *Galerkin conditions* and gives an *orthogonal projection*, which we discuss next. The case where  $\mathcal{L}$  is different from  $\mathcal{K}$  is called *oblique projection*, and we quickly have a look into this framework at the end of this section.

Let us assume that  $\mathbf{A}$  is symmetric. Let  $\mathcal{K}$  be a  $k$ -dimensional subspace of  $\mathbb{R}^n$ . An *orthogonal projection* technique onto  $\mathcal{K}$  seeks an approximate eigenpair  $(\hat{\lambda}, \hat{\mathbf{x}})$  such that  $\hat{\mathbf{x}} \in \mathcal{K}$  and

$$\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}} \perp \mathcal{K},$$

or equivalently

$$(4.1) \quad (\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \mathbf{v}) = 0 \quad \forall \mathbf{v} \in \mathcal{K}.$$

Let  $\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$  be an orthonormal basis of  $\mathcal{K}$  and  $\mathbf{Q}_k = [\mathbf{q}_1 | \dots | \mathbf{q}_k] \in \mathbb{R}^{n \times k}$ . Then, (4.1) becomes

$$(\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \mathbf{q}_i) = 0 \quad \forall i = 1, \dots, k.$$

If we express  $\hat{\mathbf{x}}$  in terms of the basis of  $\mathcal{K}$ , i.e.,  $\hat{\mathbf{x}} = \mathbf{Q}_k \mathbf{y}$ , we get

$$(\mathbf{A}\mathbf{Q}_k \mathbf{y} - \hat{\lambda}\mathbf{Q}_k \mathbf{y}, \mathbf{q}_i) = 0 \quad \forall i = 1, \dots, k,$$

and due to  $\mathbf{Q}_k^T \mathbf{Q}_k = \mathbf{I}$ , we obtain

$$\mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k \mathbf{y} = \hat{\lambda} \mathbf{y}.$$

This is the basis for *Krylov subspace methods*, which we discuss in the next section. The matrix  $\mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k \in \mathbb{R}^{k \times k}$  will often be smaller than  $\mathbf{A}$  and is either upper Hessenberg (nonsymmetric case) or tridiagonal (symmetric case). The *Rayleigh-Ritz procedure* presented in Algorithm 4.1 computes such a Galerkin approximation. The  $\theta_i$  are called *Ritz values* and  $\hat{\mathbf{x}}_i$  are the *Ritz vectors*. We will see that the Ritz

---

**Algorithm 4.1:** Rayleigh–Ritz procedure

---

- 1 Compute an orthonormal basis  $\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$  of the subspace  $\mathcal{K}$ . Set  $\mathbf{Q}_k = [\mathbf{q}_1 | \dots | \mathbf{q}_k]$ .
  - 2 Compute  $\mathbf{T}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k$ .
  - 3 Compute  $j$  eigenvalues of  $\mathbf{T}_k$ , say  $\theta_1, \dots, \theta_j$ .
  - 4 Compute the corresponding eigenvectors  $\mathbf{v}_j$  of  $\mathbf{T}_k$ . Then, the corresponding approximate eigenvectors of  $\mathbf{A}$  are  $\hat{\mathbf{x}}_j = \mathbf{Q}_k \mathbf{v}_j$ .
- 

values and Ritz vectors are the best approximate eigenpairs in the least-squares sense. But first, let us put the presented framework into a similarity transformation of  $\mathbf{A}$ : Suppose  $\mathbf{Q} = [\mathbf{Q}_k, \mathbf{Q}_u] \in \mathbb{R}^{n \times n}$  is an orthogonal matrix with  $\mathbf{Q}_k \in \mathbb{R}^{n \times k}$  being the matrix above that spans the subspace  $\mathcal{K}$ . We introduce

$$\mathbf{T} := \mathbf{Q}^T \mathbf{A} \mathbf{Q} = \begin{bmatrix} \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k & \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_u \\ \mathbf{Q}_u^T \mathbf{A} \mathbf{Q}_k & \mathbf{Q}_u^T \mathbf{A} \mathbf{Q}_u \end{bmatrix} =: \begin{bmatrix} \mathbf{T}_k & \mathbf{T}_{uk} \\ \mathbf{T}_{ku} & \mathbf{T}_u \end{bmatrix}.$$

Let  $\mathbf{T}_k = \mathbf{V} \Theta \mathbf{V}^T$  be the eigendecomposition of  $\mathbf{T}_k$ . Note that for  $k = 1$ ,  $\mathbf{T}_1$  is just the Rayleigh quotient (see Definition 2.5).

Now, we can answer the question on the "best" approximation to an eigenvector in  $\mathcal{K}$ . Similar to the observation in Section 3.1 that the Rayleigh quotient is the best eigenvalue approximation in the least-squares sense, we have the following useful result.

**THEOREM 4.1** (See [11, Theor. 7.1]). *The minimum of  $\|\mathbf{A}\mathbf{Q}_k - \mathbf{Q}_k\mathbf{R}\|_2$  over all  $k \times k$  symmetric matrices  $\mathbf{R}$  is attained by  $\mathbf{R} = \mathbf{T}_k$ , in which case  $\|\mathbf{A}\mathbf{Q}_k - \mathbf{Q}_k\mathbf{R}\|_2 = \|\mathbf{T}_{ku}\|_2$ . Let  $\mathbf{T}_k = \mathbf{V}\mathbf{\Theta}\mathbf{V}^T$  be the eigendecomposition of  $\mathbf{T}_k$ . The minimum of  $\|\mathbf{A}\mathbf{P}_k - \mathbf{P}_k\mathbf{D}\|_2$  over all  $n \times k$  orthogonal matrices  $\mathbf{P}_k$  ( $\mathbf{P}_k^T\mathbf{P}_k = \mathbf{I}$ ) where  $\text{span}(\mathbf{P}_k) = \text{span}(\mathbf{Q}_k)$  and over diagonal matrices  $\mathbf{D}$  is also  $\|\mathbf{T}_{ku}\|_2$  and is attained by  $\mathbf{P}_k = \mathbf{Q}_k\mathbf{V}$  and  $\mathbf{D} = \mathbf{\Lambda}$ .*

In practice, the columns of  $\mathbf{Q}_k$  will be computed by, e.g., the *Lanczos algorithm* or *Arnoldi algorithm*, which we discuss in Section 4.3 and 4.4.

Now, let us have a quick look at the *oblique projection* technique in which  $\mathcal{L}$  is different from  $\mathcal{K}$ . Let  $\mathcal{K}$  and  $\mathcal{L}$  be  $k$ -dimensional subspaces of  $\mathbb{R}^n$ . An *oblique projection* technique onto  $\mathcal{K}$  seeks an approximate eigenpair  $(\hat{\lambda}, \hat{\mathbf{x}})$  such that  $\hat{\mathbf{x}} \in \mathcal{K}$  and

$$\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}} \perp \mathcal{L},$$

or equivalently

$$(4.2) \quad (\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \mathbf{v}) = 0 \quad \forall \mathbf{v} \in \mathcal{L}.$$

Let  $\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$  be an orthonormal basis of  $\mathcal{K}$ ,  $\{\mathbf{p}_1, \dots, \mathbf{p}_k\}$  an orthonormal basis of  $\mathcal{L}$ ,  $\mathbf{Q}_k = [\mathbf{q}_1 | \dots | \mathbf{q}_k] \in \mathbb{R}^{n \times k}$ , and  $\mathbf{P}_k = [\mathbf{p}_1 | \dots | \mathbf{p}_k] \in \mathbb{R}^{n \times k}$ . Further, we assume biorthogonality, i.e.,  $\mathbf{P}_k^T\mathbf{Q}_k = \mathbf{I}$ . Then, (4.2) becomes

$$(\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \mathbf{p}_i) = 0 \quad \forall i = 1, \dots, k.$$

If we express  $\hat{\mathbf{x}}$  in terms of the basis of  $\mathcal{K}$ , i.e.,  $\hat{\mathbf{x}} = \mathbf{Q}_k\mathbf{y}$ , we get

$$(\mathbf{A}\mathbf{Q}_k\mathbf{y} - \hat{\lambda}\mathbf{Q}_k\mathbf{y}, \mathbf{p}_i) = 0 \quad \forall i = 1, \dots, k,$$

and due to  $\mathbf{P}_k^T\mathbf{Q}_k = \mathbf{I}$ , we obtain

$$\mathbf{P}_k^T\mathbf{A}\mathbf{Q}_k\mathbf{y} = \hat{\lambda}\mathbf{y}.$$

Oblique projection techniques form the basis for the *non-Hermitian Lanczos process* [25, 26, 41], which belongs to the class of *Krylov subspace solvers*. Krylov subspace solvers form the topic of the next section. For a further discussion on the oblique projection technique, we refer to, e.g., [47, Chap. 4.3.3].

**4.2. Krylov subspace methods.** Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Krylov subspace methods are used to solve linear systems or eigenvalue problems of sparse matrices. They only require that  $\mathbf{A}$  be accessible via a "black-box" subroutine which describes the application of  $\mathbf{A}$  to a vector. A  $k$ -dimensional *Krylov subspace* associated with a matrix  $\mathbf{A}$  and a vector  $\mathbf{v}$  is the subspace given by

$$\mathcal{K}_k(\mathbf{A}; \mathbf{v}) = \text{span}\{\mathbf{v}, \mathbf{A}\mathbf{v}, \mathbf{A}^2\mathbf{v}, \dots, \mathbf{A}^{k-1}\mathbf{v}\}.$$

The corresponding *Krylov matrix* is denoted by

$$\mathbf{K}_k(\mathbf{A}; \mathbf{v}) = [\mathbf{v} | \mathbf{A}\mathbf{v} | \mathbf{A}^2\mathbf{v} | \dots | \mathbf{A}^{k-1}\mathbf{v}].$$

Using a Krylov subspace as *right subspace*  $\mathcal{K}$  in projection methods has proven to be efficient. Various Krylov subspace methods arose from different choices of the *left subspaces*  $\mathcal{L}$ . The Krylov subspace  $\mathcal{K}_k(\mathbf{A}; \mathbf{v})$  arises naturally if we refer to it as the subspace generated by  $k - 1$  steps of the power iteration (see Section 3.1) with initial guess  $\mathbf{v}$ . Similarly, for the inverse power iteration (see Section 3.2), we obtain the subspace

$$\mathcal{K}_k((\mathbf{A} - \alpha\mathbf{I})^{-1}; \mathbf{v}).$$

Both iterations produce a sequence of vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  that span a Krylov subspace and take  $\mathbf{v}_k$  as the approximate eigenvector. Now, rather than taking  $\mathbf{v}_k$ , it is natural to use the whole sequence  $\mathbf{v}_1, \dots, \mathbf{v}_k$  in searching for the eigenvector. In fact, we saw in the previous section (Theorem 4.1 for the symmetric case) that we can even use  $\mathcal{K}_k$  to compute the  $k$  best approximate eigenvalues and eigenvectors. There are three basic algorithms for generating a basis for the Krylov subspace: the *Lanczos process* for symmetric matrices, which we discuss next, the *Arnoldi process* for nonsymmetric matrices (Section 4.4), and the *nonsymmetric Lanczos process*. The latter computes matrices  $\mathbf{Q}$  and  $\mathbf{P}$  with  $\mathbf{P}^T\mathbf{Q} = \mathbf{I}$  such that  $\mathbf{P}^T\mathbf{A}\mathbf{Q}$  is tridiagonal; see, e.g., [25, 26, 41]. Moreover, there exist *block versions* of the Arnoldi and Lanczos process; see, e.g., [8, 48], which may exploit the block structure of a matrix in some situations. They are basically an acceleration technique of the *subspace iteration*, similar to the way the subspace iteration generalizes the power methods.

**4.3. The Lanczos process.** Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be symmetric. The *Lanczos process* computes an orthogonal basis for the Krylov subspace  $\mathcal{K}_k(\mathbf{A}; \mathbf{v})$  for some initial vector  $\mathbf{v}$ , and approximates the eigenvalues of  $\mathbf{A}$  by the Ritz values.

Recall that the Hessenberg reduction of a symmetric matrix  $\mathbf{A}$  reduces to a tridiagonal matrix, i.e., there exists an orthogonal  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that

$$(4.3) \quad \mathbf{T} = \mathbf{Q}^T \mathbf{A} \mathbf{Q} = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \alpha_3 & \ddots & \\ & & \ddots & \ddots & \beta_{n-1} \\ & & & \beta_{n-1} & \alpha_n \end{bmatrix}.$$

The connection between the tridiagonalization of  $\mathbf{A}$  and the QR factorization of  $\mathbf{K}_k(\mathbf{A}; \mathbf{q}_1)$ , where  $\mathbf{q}_1 = \mathbf{Q}\mathbf{e}_1$  is given as follows:

**THEOREM 4.2** (See [22, Theor. 8.3.1]). *Let (4.3) be the tridiagonal decomposition of a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with  $\mathbf{q}_1 = \mathbf{Q}\mathbf{e}_1$ . Then:*

- (1)  $\mathbf{Q}^T \mathbf{K}_n(\mathbf{A}; \mathbf{q}_1) = \mathbf{R}$  is upper triangular.
- (2) If  $\mathbf{R}$  is nonsingular, then  $\mathbf{T}$  is unreduced.
- (3) If  $k = \arg \min_{j=1, \dots, n} \{r_{j,j} = 0\}$ , then  $k - 1 = \arg \min_{j=1, \dots, n-1} \{\beta_j = 0\}$ .

It follows from (1) in Theorem 4.2 that  $\mathbf{Q}\mathbf{R}$  is the QR factorization of  $\mathbf{K}_n(\mathbf{A}; \mathbf{q}_1)$ . In order to preserve the sparsity, we need an alternative to similarity transformations in order to compute the tridiagonalization. Let us write  $\mathbf{Q} = [\mathbf{q}_1 | \dots | \mathbf{q}_n]$ . Considering the  $k$ th column of  $\mathbf{A}\mathbf{Q} = \mathbf{Q}\mathbf{T}$ , we obtain the following three-term recurrence:

$$(4.4) \quad \mathbf{A}\mathbf{q}_k = \beta_{k-1}\mathbf{q}_{k-1} + \alpha_k\mathbf{q}_k + \beta_k\mathbf{q}_{k+1}.$$

Since the columns of  $\mathbf{Q}$  are orthonormal, multiplying (4.4) from the left by  $\mathbf{q}_k$  yields

$$\alpha_k = \mathbf{q}_k^T \mathbf{A} \mathbf{q}_k.$$

This leads to the method in Algorithm 4.2 developed by Lanczos in 1950 [34].

---

**Algorithm 4.2:** Lanczos process

---

```

1 Given  $\mathbf{q}_0 = \mathbf{0}$ ,  $\mathbf{q}_1 = \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$ ,  $\beta_0 = 0$ 
2 for  $k = 1, 2, \dots$  do
3    $\mathbf{z}_k = \mathbf{A} \mathbf{q}_k$ 
4    $\alpha_k = \mathbf{q}_k^T \mathbf{z}_k$ 
5    $\mathbf{z}_k = \mathbf{z}_k - \beta_{k-1} \mathbf{q}_{k-1} - \alpha_k \mathbf{q}_k$ 
6    $\beta_k = \|\mathbf{z}_k\|_2$ 
7   if  $\beta_k = 0$  then
8     quit
9   end
10   $\mathbf{q}_{k+1} = \frac{\mathbf{z}_k}{\beta_k}$ 
11 end
```

---

The vectors  $\mathbf{q}_k$  computed by the Lanczos algorithm are called *Lanczos vectors*. The Lanczos process stops before the complete tridiagonalization if  $\mathbf{q}_1$  is contained in an exact  $\mathbf{A}$ -invariant subspace:

**THEOREM 4.3** (See [22, Theor. 10.1.1]). *The Lanczos Algorithm 4.2 runs until  $k = m$ , where*

$$m = \text{rank}(\mathbf{K}_n(\mathbf{A}, \mathbf{q}_1)).$$

Moreover, for  $k = 1, \dots, m$ , we have

$$(4.5) \quad \mathbf{A} \mathbf{Q}_k = \mathbf{Q}_k \mathbf{T}_k + \beta_k \mathbf{q}_{k+1} \mathbf{e}_k^T,$$

where  $\mathbf{T}_k = \mathbf{T}(1:k, 1:k)$ ,  $\mathbf{Q}_k = [\mathbf{q}_1 | \dots | \mathbf{q}_k]$  has orthonormal columns with

$$\text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_k\} = \mathcal{K}_k(\mathbf{A}, \mathbf{q}_1).$$

In particular,  $\beta_m = 0$ , and hence

$$\mathbf{A} \mathbf{Q}_m = \mathbf{Q}_m \mathbf{T}_m.$$

The eigenvalues of the tridiagonal  $\mathbf{T}_m$  can then be computed via, e.g., the QR iteration. A corresponding eigenvector can be obtained by using the inverse power iteration with the approximated eigenvalue as shift.

We can show that the quality of the approximation after  $k$  Lanczos steps depends on  $\beta_k$  and on parts of the eigenvectors of  $\mathbf{T}_k$  (cf. [22, Chap. 10.1.4]): Therefore, let  $(\theta, \mathbf{y})$  be an eigenpair of  $\mathbf{T}_k$ . Applying (4.5) to  $\mathbf{y}$  yields

$$\begin{aligned} \mathbf{A} \mathbf{Q}_k \mathbf{y} &= \mathbf{Q}_k \mathbf{T}_k \mathbf{y} + \beta_k \mathbf{q}_{k+1} \mathbf{e}_k^T \mathbf{y} \\ &= \theta \mathbf{Q}_k \mathbf{y} + \beta_k \mathbf{q}_{k+1} \mathbf{e}_k^T \mathbf{y} \end{aligned}$$

and hence the following error estimation (cf. Theorem 4.1)

$$\|\mathbf{A} \mathbf{Q}_k \mathbf{y} - \theta \mathbf{Q}_k \mathbf{y}\|_2 = |\beta_k| |\mathbf{e}_k^T \mathbf{y}|.$$

Hence, we want to accomplish  $\beta_k = 0$  fast. Regarding the convergence theory, we refer to [31, 38, 43] and [22, Chap. 10.1.5]. In summary, the Ritz values converge fast to the extreme eigenvalues. Using shift and invert strategies (as in the inverse power method in Section 3.2), we can obtain convergence to interior eigenvalues. In practice, rounding errors have a significant effect on the behavior of the Lanczos iteration. If the computed  $\beta_k$  are close to zero, then the Lanczos vectors lose their orthogonality. Reorthogonalization strategies provide a remedy; see, e.g., [38, 21, 40, 50, 6, 67]. Nevertheless, we know from the last section that the Ritz values and vectors are good approximations.

So far, we have assumed that  $\mathbf{A}$  is symmetric. Next, we consider the nonsymmetric case.

**4.4. The Arnoldi process.** For a nonsymmetric  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , we know that there exists a Hessenberg decomposition, i.e., there exists an orthogonal  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that

$$(4.6) \quad \mathbf{H} = \mathbf{Q}^T \mathbf{A} \mathbf{Q} = \begin{bmatrix} h_{1,1} & h_{1,2} & h_{1,3} & \cdots & h_{1,n} \\ h_{2,1} & h_{2,2} & h_{2,3} & \cdots & h_{2,n} \\ 0 & h_{3,2} & h_{3,3} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & h_{n-1,n} \\ 0 & \cdots & 0 & h_{n,n-1} & h_{n,n} \end{bmatrix}.$$

The connection between the Hessenberg reduction of  $\mathbf{A}$  and the QR factorization of  $\mathbf{K}_k(\mathbf{A}; \mathbf{q}_1)$ , where  $\mathbf{q}_1 = \mathbf{Q} \mathbf{e}_1$  is given as follows (cf. the symmetric case in Theorem 4.2)

**THEOREM 4.4** (See [22, Theor. 7.4.3]). *Suppose  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is orthogonal and let  $\mathbf{q}_1 = \mathbf{Q} \mathbf{e}_1$  and  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Then,  $\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{H}$  is an unreduced upper Hessenberg matrix if and only if  $\mathbf{Q}^T \mathbf{K}_n(\mathbf{A}; \mathbf{q}_1) = \mathbf{R}$  is nonsingular and upper triangular.*

It follows from Theorem 4.4 that  $\mathbf{Q} \mathbf{R}$  is the QR factorization of  $\mathbf{K}_n(\mathbf{A}; \mathbf{q}_1)$ . As before, in order to preserve the sparsity, we need an alternative to similarity transformations in order to compute the Hessenberg reduction. Let us write  $\mathbf{Q} = [\mathbf{q}_1 | \dots | \mathbf{q}_n]$ . Considering the  $k$ th column of  $\mathbf{A} \mathbf{Q} = \mathbf{Q} \mathbf{H}$ , we obtain the following recurrence:

$$(4.7) \quad \mathbf{A} \mathbf{q}_k = \sum_{i=1}^{k+1} h_{i,k} \mathbf{q}_i.$$

Since the columns of  $\mathbf{Q}$  are orthonormal, multiplying (4.7) from the left by  $\mathbf{q}_i$  yields

$$h_{i,k} = \mathbf{q}_i^T \mathbf{A} \mathbf{q}_k$$

for  $i = 1, \dots, k$ . This leads to the method in Algorithm 4.3 developed by Arnoldi in 1951 [1]. It can be viewed as an extension of the Lanczos process to nonsymmetric matrices. Note that, in contrast to the symmetric case, we have no three-term recurrence anymore. Hence, we have to store all computed vectors  $\mathbf{q}_k$ . These vectors are called *Arnoldi vectors*. The Arnoldi process can be seen as a modified Gram-Schmidt orthogonalization process (cf. [22, Chap. 5.2.8]) since in each step  $k$  we orthogonalize  $\mathbf{A} \mathbf{q}_k$  against all previous  $\mathbf{q}_i$  — this requires  $\mathcal{O}(kn)$  flops. Hence,

**Algorithm 4.3:** Arnoldi process

---

```

1 Given  $\mathbf{q}_1 = \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$ 
2 for  $k = 1, 2, \dots$  do
3    $\mathbf{z}_k = \mathbf{A}\mathbf{q}_k$ 
4   for  $i = 1, \dots, k$  do
5      $h_{i,k} = \mathbf{q}_i^T \mathbf{z}_k$ 
6      $\mathbf{z}_k = \mathbf{z}_k - h_{i,k} \mathbf{q}_i$ 
7   end
8    $h_{k+1,k} = \|\mathbf{z}_k\|_2$ 
9   if  $h_{k+1,k} = 0$  then
10    quit
11  end
12   $\mathbf{q}_{k+1} = \frac{\mathbf{z}_k}{h_{k+1,k}}$ 
13 end

```

---

the computational cost grows rapidly with the number of steps. After  $k$  steps of the Arnoldi Algorithm 4.3, we have

$$(4.8) \quad \mathbf{A}\mathbf{Q}_k = \mathbf{Q}_k \mathbf{H}_k + h_{k+1,k} \mathbf{q}_{k+1} \mathbf{e}_k^T,$$

where  $\mathbf{H}_k = \mathbf{H}(1:k, 1:k)$  and

$$(4.9) \quad \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_k\} = \mathcal{K}_k(\mathbf{A}, \mathbf{q}_1).$$

We can show that the quality of the approximation depends on the magnitude of  $h_{k+1,k}$  and on parts of the eigenvectors of  $\mathbf{H}_k$  (cf. [22, Chap. 10.5.1]): Therefore, let  $(\theta, \mathbf{y})$  be an eigenpair of  $\mathbf{H}_k$ . Applying (4.8) to  $\mathbf{y}$  yields

$$\begin{aligned} \mathbf{A}\mathbf{Q}_k \mathbf{y} &= \mathbf{Q}_k \mathbf{H}_k \mathbf{y} + h_{k+1,k} \mathbf{q}_{k+1} \mathbf{e}_k^T \mathbf{y} \\ &= \theta \mathbf{Q}_k \mathbf{y} + h_{k+1,k} \mathbf{q}_{k+1} \mathbf{e}_k^T \mathbf{y} \end{aligned}$$

and hence the following error estimation (cf. Theorem 4.1)

$$\|\mathbf{A}\mathbf{Q}_k \mathbf{y} - \theta \mathbf{Q}_k \mathbf{y}\|_2 = |h_{k+1,k}| |\mathbf{e}_k^T \mathbf{y}|.$$

Hence, we want  $h_{k+1,k} = 0$  fast. As the Lanczos process, the Arnoldi process has a (lucky) breakdown at step  $k = m$  if  $h_{m+1,m} = 0$  since  $\mathcal{K}_m(\mathbf{A}, \mathbf{q}_1)$  is an  $\mathbf{A}$ -invariant subspace in this case. In the following, we discuss accelerating techniques for the Arnoldi and Lanczos process.

**4.5. Restarted Arnoldi and Lanczos.** Note that with each Arnoldi step we have to store one additional Arnoldi vector. A remedy is restarting the Arnoldi process with carefully chosen restarts after a certain maximum of steps is reached. Acceleration techniques (mainly of a polynomial nature) generate an initial guess with small components in the unwanted parts of the spectrum. The strategies we present are called *polynomial acceleration* or *filtering techniques*. They exploit the powers of a matrix similar as the power method in the sense that they generate iterations of the form

$$\mathbf{z}_r = p_r(\mathbf{A}) \mathbf{z}_0,$$

where  $p_r$  is a polynomial of degree  $r$ . In the case of the power method, we have  $p_r(t) = t^r$ . Filtering methods have been successfully combined with the subspace iteration. When combined with the Arnoldi process, they are often called *implicitly restarted methods*, which we discuss next. Selecting a good polynomial often relies on some knowledge of the eigenvalues or related quantities (e.g., Ritz values).

Suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is diagonalizable and has eigenpairs  $\{(\lambda_i, \mathbf{x}_i)\}_{i=1, \dots, n}$  with  $\lambda_1 \geq \dots \geq \lambda_n$ . Let

$$\mathbf{q}_1 = \sum_{i=1}^n \alpha_i \mathbf{x}_i$$

be an initial guess for the Arnoldi process. After running  $r$  steps of the Arnoldi process, we do a restart. We may seek a new initial vector from the span of the Arnoldi vectors  $\mathbf{q}_1, \dots, \mathbf{q}_r$ , which has, due to (4.9), the form

$$\begin{aligned} \mathbf{q}_+ &= \sum_{j=1}^r \beta_j \mathbf{A}^{j-1} \mathbf{q}_1 = \sum_{j=1}^r \beta_j \sum_{i=1}^n \alpha_i \mathbf{A}^{j-1} \mathbf{x}_i = \sum_{j=1}^r \beta_j \sum_{i=1}^n \alpha_i \lambda_i^{j-1} \mathbf{x}_i \\ &= \sum_{i=1}^n \alpha_i p_{r-1}(\lambda_i) \mathbf{x}_i. \end{aligned}$$

Suppose we are interested in the eigenvalue  $\lambda_j$ . If  $|\alpha_j p_{r-1}(\lambda_j)| \gg |\alpha_l p_{r-1}(\lambda_l)|$  for all  $l \neq j$ , then  $\mathbf{q}_+$  has large components in the eigendirection  $\mathbf{x}_j$ . Note that the  $\alpha_i$  are unknown. Hence, with an appropriate constructed polynomial, we can amplify the components in the desired parts of the spectrum. For instance, we are seeking for a polynomial that satisfies  $p_{r-1}(\lambda_j) = 1$  and  $|p_{r-1}(\lambda_l)| \ll 1$  for all  $l \neq j$ . However, the eigenvalues  $\lambda_i$  are unknown as well. Hence we need some approximation. Let  $\Omega$  be a domain (e.g., an ellipse) that contains  $\lambda(\mathbf{A}) \setminus \{\lambda_j\}$ , and suppose we have an estimate of  $\lambda_j$ . Then, we can aim to solve

$$\min_{\substack{p_{r-1} \in P_{r-1}, \\ p_{r-1}(\lambda_j) = 1}} \max_{t \in \Omega} |p_{r-1}(t)|.$$

Suitable polynomials include the shifted and scaled Chebyshev polynomials, and in the symmetric case, we can exploit the three-term recurrence for fast computation; see, e.g., [45].

An alternative to Chebyshev polynomials is the following: Given  $\{\theta_i\}_{i=1, \dots, r-1}$ , then one natural idea is to set

$$(4.10) \quad p_{r-1}(t) = (t - \theta_1)(t - \theta_2) \cdots (t - \theta_{r-1});$$

see [22, Chap. 10.5.2]. If  $\lambda_i \approx \theta_l$  for some  $l$ , then  $\mathbf{q}_+$  has small components in the eigendirection  $\mathbf{x}_i$ . Hence, the  $\theta_i$  are all unwanted values. For  $\theta_i$  we can use the Ritz values, which presumably approximate the eigenvalues of  $\mathbf{A}$ . For further heuristics, we refer to [44].

The above strategies are *explicit restarting* techniques, which use only one vector for the restart. The following *implicit restarting* strategy uses  $k$  vectors from the previous Arnoldi process for the new restarted Arnoldi process and throws away the remaining  $r - k =: p$  vectors. The procedure was developed in 1992 [53]. It implicitly determines a polynomial of the form (4.10) using the QR iteration with

shifts. Suppose we have performed  $r$  steps of the Arnoldi iteration with starting vector  $\mathbf{q}_1$ . Due to (4.8), we have

$$(4.11) \quad \mathbf{A}\mathbf{Q}_r = \mathbf{Q}_r\mathbf{H}_r + h_{r+1,r}\mathbf{q}_{r+1}\mathbf{e}_r^T,$$

where  $\mathbf{H}_r \in \mathbb{R}^{r \times r}$  is upper Hessenberg,  $\mathbf{Q}_r \in \mathbb{R}^{n \times r}$  has orthonormal columns, and  $\mathbf{Q}_r\mathbf{e}_1 = \mathbf{q}_1$ . Next, we apply  $p$  steps of the QR iteration with shifts  $\theta_1, \dots, \theta_p$  (Algorithm 3.6), i.e., in step  $i$  we compute

$$(4.12) \quad \mathbf{V}_i\mathbf{R}_i = \mathbf{H}^{(i-1)} - \theta_i\mathbf{I},$$

$$(4.13) \quad \mathbf{H}^{(i)} = \mathbf{R}_i\mathbf{V}_i + \theta_i\mathbf{I},$$

where  $\mathbf{H}^{(0)} = \mathbf{H}_r$ . After  $p$  steps, we have

$$\mathbf{H}^{(p)} = \mathbf{R}_p\mathbf{V}_p + \theta_p\mathbf{I} = \mathbf{V}_p^T(\mathbf{V}_p\mathbf{R}_p + \theta_p\mathbf{I})\mathbf{V}_p = \mathbf{V}_p^T\mathbf{H}^{(p-1)}\mathbf{V}_p = \dots = \mathbf{V}^T\mathbf{H}^{(0)}\mathbf{V}$$

with  $\mathbf{V} = \mathbf{V}_1 \cdots \mathbf{V}_p$ . We use the notation

$$(4.14) \quad \mathbf{H}_+ := \mathbf{H}^{(p)} = \mathbf{V}^T\mathbf{H}^{(0)}\mathbf{V} = \mathbf{V}^T\mathbf{H}_r\mathbf{V}.$$

The relationship to a polynomial of the form (4.10) is the following:

**THEOREM 4.5** (See [22, Theor. 10.5.1]). *If  $\mathbf{V} = \mathbf{V}_1 \cdots \mathbf{V}_p$  and  $\mathbf{R} = \mathbf{R}_p \cdots \mathbf{R}_1$  are defined by (4.12)–(4.13), then*

$$\mathbf{V}\mathbf{R} = (\mathbf{H}_r - \theta_1\mathbf{I}) \cdots (\mathbf{H}_r - \theta_p\mathbf{I}).$$

Using (4.14), we get in (4.11)

$$(4.15) \quad \mathbf{A}\mathbf{Q}_r = \mathbf{Q}_r\mathbf{V}\mathbf{H}_+\mathbf{V}^T + h_{r+1,r}\mathbf{q}_{r+1}\mathbf{e}_r^T.$$

Multiplying (4.15) from the right by  $\mathbf{V}$  yields

$$(4.16) \quad \mathbf{A}\mathbf{Q}_+ = \mathbf{Q}_+\mathbf{H}_+ + h_{r+1,r}\mathbf{q}_{r+1}\mathbf{e}_r^T\mathbf{V}$$

with  $\mathbf{Q}_+ = \mathbf{Q}_r\mathbf{V}$ . It can be shown that  $\mathbf{V}_1, \dots, \mathbf{V}_p$  from the shifted QR iteration are upper Hessenberg. Hence,  $\mathbf{V}(r, 1 : r - p - 1) = \mathbf{0}^T$  and therefore  $\mathbf{e}_r^T\mathbf{V} = [0 \cdots 0 \alpha * \cdots *]$  is a row vector of length  $r$  whose first  $r - p - 1$  entries are zero. Now, using the notation  $\mathbf{Q}_+ = [\hat{\mathbf{Q}}_{r-p}, \hat{\mathbf{Q}}_p]$  with  $\hat{\mathbf{Q}}_{r-p} \in \mathbb{R}^{n \times (r-p)}$  we can write (4.16) as

$$(4.17) \quad \mathbf{A}[\hat{\mathbf{Q}}_{r-p}, \hat{\mathbf{Q}}_p] = [\hat{\mathbf{Q}}_{r-p}, \hat{\mathbf{Q}}_p] \begin{bmatrix} \hat{\mathbf{H}}_{r-p} & * \\ \beta\mathbf{e}_1\mathbf{e}_{r-p}^T & * \end{bmatrix} + h_{r+1,r}\mathbf{q}_{r+1} \underbrace{[0 \cdots 0]}_{r-p-1} \alpha * \cdots *.$$

Now, we throw away the last  $p$  columns in (4.17) and obtain an  $(r-p)$ -step Arnoldi decomposition

$$\begin{aligned} \mathbf{A}\hat{\mathbf{Q}}_{r-p} &= \hat{\mathbf{Q}}_{r-p}\hat{\mathbf{H}}_{r-p} + \beta\hat{\mathbf{Q}}_p\mathbf{e}_1\mathbf{e}_{r-p}^T + h_{r+1,r}\mathbf{q}_{r+1} \underbrace{[0 \cdots 0]}_{r-p-1} \alpha \\ &= \hat{\mathbf{Q}}_{r-p}\hat{\mathbf{H}}_{r-p} + \left( \beta\hat{\mathbf{Q}}_p\mathbf{e}_1 + \alpha h_{r+1,r}\mathbf{q}_{r+1} \right) \mathbf{e}_{r-p}^T \\ &=: \hat{\mathbf{Q}}_{r-p}\hat{\mathbf{H}}_{r-p} + \hat{\mathbf{v}}_{r+1}\mathbf{e}_{r-p}^T. \end{aligned}$$

This is the Arnoldi recursion we would have obtained by restarting the Arnoldi process with the starting vector  $\mathbf{q}_+ = \mathbf{Q}_+\mathbf{e}_1$ . Hence, we do not need to restart the Arnoldi process from step one but rather from step  $r - p + 1$ . For further details, we refer to [22, Chap. 10.5.3] and the references therein.



REMARK 4.6. It can be shown (cf. [22, Chap. 10.5.3]) that

$$\mathbf{q}_+ = c(\mathbf{A} - \theta_1 \mathbf{I}) \cdots (\mathbf{A} - \theta_p \mathbf{I}) \mathbf{Q}_r \mathbf{e}_1$$

for some scalar  $c$  and is hence of the form (4.10).

For further reading on the Arnoldi process we refer to, e.g., [54, 47, 62, 55].

Next we present another acceleration technique which is very popular for solving linear systems.

**4.6. Preconditioning.** In the following, we quickly review the preconditioning concept for solving large and sparse systems of linear equations of the general form

$$(4.18) \quad \mathbf{A}\mathbf{z} = \mathbf{b}.$$

Here,  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is the given coefficient matrix,  $\mathbf{z} \in \mathbb{R}^n$  is the unknown solution vector, and  $\mathbf{b} \in \mathbb{R}^n$  is the given right-hand side vector. In order for the Equation (4.18) to have a unique solution, we assume that  $\mathbf{A}$  is nonsingular. Systems of the form (4.18) arise after the discretization of a continuous problem like partial differential equations such as the time-harmonic Maxwell equations. Other applications arise in incompressible magnetohydrodynamics as well as constrained optimization. As already mentioned in Section 4.2, Krylov subspace solvers are also used for solving linear systems. In fact, they are state-of-the-art iterative solvers. However, they are usually only efficient in combination with an accelerator, which is called a *preconditioner*. The aim of a preconditioner is to enhance the convergence of the iterative solver. In our case, we want to accelerate the speed of convergence of Krylov subspace solvers. The basic idea is to construct a nonsingular matrix  $\mathbf{P} \in \mathbb{R}^{n \times n}$  and solve

$$(4.19) \quad \mathbf{P}^{-1} \mathbf{A}\mathbf{z} = \mathbf{P}^{-1} \mathbf{b}$$

instead of  $\mathbf{A}\mathbf{z} = \mathbf{b}$ . In order for  $\mathbf{P}$  to be efficient, it should approximate  $\mathbf{A}$ , and at the same time, the action of  $\mathbf{P}^{-1}$  should require little work. The construction process of  $\mathbf{P}$  should incorporate the goal of eigenvalue clustering. That means,  $\mathbf{P}^{-1} \mathbf{A}$  is aimed to have a few number of eigenvalues or eigenvalue clusters. This is bases on the following: In a nutshell, for linear systems the residual of a Krylov subspace solver  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{z}_k$  satisfies  $\mathbf{r}_k = p_k(\mathbf{A})\mathbf{r}_0$ , and one approach would be to minimize the norm of the residual, which amounts to requiring that  $\|p_k(\lambda_i)\mathbf{v}_i\|_2$  be as small as possible for all  $i = 1, \dots, n$ . Here,  $\{(\lambda_i, \mathbf{v}_i)\}_{i=1, \dots, n}$  are the eigenpairs of  $\mathbf{A}$ . Therefore, replacing  $\mathbf{A}$  by  $\mathbf{P}^{-1} \mathbf{A}$  such that  $\mathbf{P}^{-1} \mathbf{A}$  has more clustered eigenvalues is one way to go. This typically results in outstanding performances of Krylov subspace solvers. For an overview of iterative solvers and preconditioning techniques, we refer to [20, 23, 46, 16, 5, 2, 4, 59].

Preconditioning plays an important role in eigenvalue problems as well. Taken in the same spirit as seeking an operator that improves the spectrum, we can think of the inverse power iteration (see Section 3.2) as a preconditioning approach: The operator  $(\mathbf{A} - \theta \mathbf{I})^{-1}$  has a much better spectrum than  $\mathbf{A}$  for a suitable chosen shift  $\theta$ . So, we can run Arnoldi on  $(\mathbf{A} - \theta \mathbf{I})^{-1}$  rather than  $\mathbf{A}$  since the eigenvectors of  $\mathbf{A}$  and  $(\mathbf{A} - \theta \mathbf{I})^{-1}$  are identical. Another idea is to incorporate polynomial preconditioning, i.e., replace  $\mathbf{A}$  by  $p_k(\mathbf{A})$ . As a guideline, we want to transform

the  $k$  wanted eigenvalues of  $\mathbf{A}$  to  $k$  eigenvalues of  $p_k(\mathbf{A})$  that are much larger than the other eigenvalues, so as to accelerate convergence. Preconditioning also plays a role in solving *generalized eigenvalue problems*

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{B}\mathbf{x}.$$

They can be solved, e.g., by the *Jacobi–Davidson method*, whose idea we briefly discuss in Section 4.8 for solving the standard algebraic eigenvalue problem  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ . The discussion of generalized eigenproblems is out of the scope of this survey. We refer the reader to [56, 57, 62] for a background to these problems.

**4.7. Davidson method.** Davidson’s method is basically a preconditioned version of the Lanczos process, but the amount of work increases similarly to Arnoldi, due to increased orthogonalization requirements. Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathcal{K}_k = \mathcal{K}_k(\mathbf{A}; \mathbf{v})$  be a Krylov subspace with respect to some vector  $\mathbf{v}$ . Let  $\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$  be an orthonormal basis of  $\mathcal{K}_k$ . In the orthogonal projection technique, we are seeking for an  $\hat{\mathbf{x}} \in \mathcal{K}_k$  such that

$$\left(\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \mathbf{q}_i\right) = 0 \quad \forall i = 1, \dots, k,$$

Suppose, we have a Ritz pair  $(\theta_i, \mathbf{u}_i)$ . Then, the residual is given by

$$\mathbf{r}_i = \mathbf{A}\mathbf{u}_i - \theta_i\mathbf{u}_i = (\mathbf{A} - \theta_i\mathbf{I})\mathbf{u}_i.$$

Now, we can improve the eigenpair approximation by precondition the residual, i.e., by solving

$$(\mathbf{P} - \theta_i\mathbf{I})\mathbf{t} = \mathbf{r}_i,$$

and define  $\mathbf{t}$  as a new search direction, enriching the subspace. That is,  $\mathbf{t}$  is orthogonalized against all basis vectors  $\mathbf{q}_1, \dots, \mathbf{q}_k$ , and the resulting vector  $\mathbf{q}_{k+1}$  enriches  $\mathcal{K}_k$  to  $\mathcal{K}_{k+1}$ .

Davidson [10] originally proposed to precondition with the diagonal matrix of  $\mathbf{A}$ , i.e.,  $\mathbf{P} = \text{diag}(\mathbf{A})$ , since he dealt with a diagonal dominant matrix  $\mathbf{A}$ . Additionally, diagonal preconditioning offers a computationally cheap iteration. For the use of other preconditioners, we refer to [7]. Further references on Davidson’s method include [54, 55, 36].

Next, we consider an extension of Davidson’s method, which has the potential of working better for matrices that are not diagonally dominant.

**4.8. Jacobi–Davidson method.** The idea is to extend the strategy of preconditioning the residual. If  $(\hat{\lambda}, \hat{\mathbf{x}})$  with  $\|\hat{\mathbf{x}}\|_2 = 1$  is an approximate eigenpair of  $\mathbf{A}$ , then the residual is  $\mathbf{r} = \mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}$ . Now, we look for  $(\hat{\lambda} + \delta\hat{\lambda}, \hat{\mathbf{x}} + \delta\hat{\mathbf{x}})$  to improve the eigenpair. We write

$$\mathbf{A}(\hat{\mathbf{x}} + \delta\hat{\mathbf{x}}) = (\hat{\lambda} + \delta\hat{\lambda})(\hat{\mathbf{x}} + \delta\hat{\mathbf{x}}),$$

which is equivalent to

$$(\mathbf{A} - \hat{\lambda}\mathbf{I})\delta\hat{\mathbf{x}} - \delta\hat{\lambda}\hat{\mathbf{x}} = -\mathbf{r} + \delta\hat{\lambda}\delta\hat{\mathbf{x}}.$$

By neglecting the second-order term, we obtain

$$(\mathbf{A} - \hat{\lambda}\mathbf{I})\delta\hat{\mathbf{x}} - \delta\hat{\lambda}\hat{\mathbf{x}} = -\mathbf{r}.$$

This is an underdetermined system and a constraint must be added, e.g.,  $\|\hat{\mathbf{x}} + \delta\hat{\mathbf{x}}\|_2 = 1$ . With  $\|\hat{\mathbf{x}}\|_2 = 1$  and neglecting the second-order term, this condition becomes

$$\hat{\mathbf{x}}^T \delta\hat{\mathbf{x}} = 0.$$

If  $\hat{\lambda} = \hat{\mathbf{x}}^T \mathbf{A} \hat{\mathbf{x}}$ , then we obtain  $\delta\hat{\mathbf{x}}$  by solving the projected system

$$\begin{aligned} (\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) (\mathbf{A} - \hat{\lambda}\mathbf{I}) (\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) \delta\hat{\mathbf{x}} &= - (\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) (\mathbf{r} - \delta\hat{\lambda}\hat{\mathbf{x}}) \\ &= - (\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) \mathbf{r} \\ &= - (\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) (\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}) \\ &= - (\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) \mathbf{A}\hat{\mathbf{x}} \\ &= - (\mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}) = -\mathbf{r} \end{aligned}$$

subject to the constraint  $\hat{\mathbf{x}}^T \delta\hat{\mathbf{x}} = 0$ . As in the previous section, we replace  $\mathbf{A}$  by a preconditioner  $\mathbf{P}$ , such that we have to solve an approximate projected system

$$(\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) (\mathbf{P} - \hat{\lambda}\mathbf{I}) (\mathbf{I} - \hat{\mathbf{x}}\hat{\mathbf{x}}^T) \delta\hat{\mathbf{x}} = -\mathbf{r}$$

subject to the constraint  $\hat{\mathbf{x}}^T \delta\hat{\mathbf{x}} = 0$ .

The connection of the described method to Jacobi is given in Remark 4.7.

REMARK 4.7. Given an approximate eigenpair  $(\hat{\lambda}, \hat{\mathbf{x}})$  of  $\mathbf{A}$ , Jacobi [30] proposed to solve an eigenvalue problem  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  by finding a correction  $\mathbf{t}$  such that

$$\mathbf{A}(\hat{\mathbf{x}} + \mathbf{t}) = \lambda(\hat{\mathbf{x}} + \mathbf{t}), \quad \hat{\mathbf{x}} \perp \mathbf{t}.$$

This is called the *Jacobi Orthogonal Component Correction (JOCC)*.

The Jacobi–Davidson framework can also be connected with Newton’s method; see, e.g., [47, Chap. 8.4].

The debate over the advantaged and disadvantages of Jacobi–Davidson versus other approaches such as the Arnoldi process (with shift and invert) is delicate. Sleijpen and van der Vorst [51] relate it to whether the new direction has a strong component in previous directions. It is a fairly technical argument, and not much theory is available. For more details about the Jacobi–Davidson method, we refer to [51, 52, 54].

## 5. Conclusions

The numerical solution of eigenvalue problems is an extremely active area of research. Eigenvalues are very important in many areas of applications, and challenges keep arising. The survey covers only some basic principles, which have established themselves as the fundamental building blocks of eigenvalue solvers. We have left out some important recent advances, which are extremely important but also rather technical. Generalized eigenvalue problems are also very important, but there is not enough room to cover them in this survey.

One of the main messages of this survey is the distinction between important *mathematical* observations about eigenvalues, and practical *computational* considerations. Objects such as the Jordan Canonical Form or determinants are classical mathematical tools, but they cannot be easily utilized in practical computations. On the

other hand, sparsity of the matrix and the availability of matrix decompositions are often overlooked when a pure mathematical discussion of the problem ensues, but they are absolutely essential in the design of numerical methods.

Altogether, this topic is satisfyingly rich and challenging. Efficiently and accurately computing eigenvalues and eigenvectors of matrices continues to be one of the most important problems in mathematical sciences.

### References

1. W. E. Arnoldi, *The principle of minimized iteration in the solution of the matrix eigenvalue problem*, Quart. Appl. Math. **9** (1951), 17–29.
2. O. Axelsson, *A survey of preconditioned iterative methods for linear systems of algebraic equations*, BIT **25** (1985), no. 1, 165–187.
3. W. Barth, R. S. Martin, and J. H. Wilkinson, *Handbook Series Linear Algebra: Calculation of the eigenvalues of a symmetric tridiagonal matrix by the method of bisection*, Numer. Math. **9** (1967), no. 5, 386–393.
4. M. Benzi, *Preconditioning techniques for large linear systems: A survey*, J. Comput. Phys. **182** (2002), no. 2, 418–477.
5. M. Benzi, G. H. Golub, and J. Liesen, *Numerical solution of saddle point problems*, Acta Numer. **14** (2005), 1–137.
6. D. Calvetti, L. Reichel, and D. C. Sorensen, *An implicitly restarted Lanczos method for large symmetric eigenvalue problems*, Electron. Trans. Numer. Anal. **2** (1994), no. 1, 1–21.
7. M. Crouzeix, B. Philippe, and M. Sadkane, *The Davidson method*, SIAM J. Sci. Comput. **15** (1994), no. 1, 62–76.
8. J. Cullum and W. E. Donath, *A block Lanczos algorithm for computing the  $q$  algebraically largest eigenvalues and a corresponding eigenspace of large, sparse, real symmetric matrices*, 1974 IEEE Conference on Decision and Control, 1974, pp. 505–509.
9. J. J. M. Cuppen, *A divide and conquer method for the symmetric tridiagonal eigenproblem*, Numer. Math. **36** (1981), no. 2, 177–195.
10. E. R. Davidson, *The iterative calculation of a few of the lowest eigenvalues and corresponding eigenvectors of large real-symmetric matrices*, J. Comput. Phys. **17** (1975), no. 1, 87–94.
11. J. Demmel, *Applied numerical linear algebra*, SIAM, 1997.
12. J. Demmel and K. Veselić, *Jacobi’s method is more accurate than QR*, SIAM J. Matrix Anal. Appl. **13** (1992), no. 4, 1204–1245.
13. I. S. Dhillon and B. N. Parlett, *Multiple representations to compute orthogonal eigenvectors of symmetric tridiagonal matrices*, Linear Algebra Appl. **387** (2004), 1–28.
14. J. J. Dongarra and D. C. Sorensen, *A fully parallel algorithm for the symmetric eigenvalue problem*, SIAM J. Sci. Stat. Comp. **8** (1987), no. 2, s139–s154.
15. A. A. Dubrulle and G. H. Golub, *A multishift QR iteration without computation of the shifts*, Numer. Algorithms **7** (1994), no. 2, 173–181.
16. H. C. Elman, D. J. Silvester, and A. J. Wathen, *Finite elements and fast iterative solvers: With applications in incompressible fluid dynamics*, Numer. Math. Sci. Comput., Oxford Univ. Press, Oxford, 2005.
17. J. Erxiong, *A note on the double-shift QL algorithm*, Linear Algebra Appl. **171** (1992), 121–132.
18. J. G. F. Francis, *The QR transformation: A unitary analogue to the LR transformation—Part 1*, Comput. J. **4** (1961), no. 3, 265–271.
19. ———, *The QR transformation—Part 2*, The Comput. J. **4** (1962), no. 4, 332–345.
20. R. W. Freund, G. H. Golub, and N. M. Nachtigal, *Iterative solution of linear systems*, Acta Numer. **1** (1992), 57–100.
21. G. H. Golub, R. R. Underwood, and J. H. Wilkinson, *The Lanczos algorithm for the symmetric  $Ax = \lambda Bx$  problem*, Tech. report, Dep. Comput. Sci., Stanford Univ., Stanford, CA, 1972.
22. G. H. Golub and C. F. van Loan, *Matrix computations*, 4th ed., Johns Hopkins Stud. Math. Sci., Johns Hopkins Univ. Press, Baltimore, MD, 2013.
23. A. Greenbaum, *Iterative methods for solving linear systems*, Frontiers Appl. Math., vol. 17, SIAM, Philadelphia, PA, 1997.

24. M. Gu and S. C. Eisenstat, *A divide-and-conquer algorithm for the symmetric tridiagonal eigenproblem*, SIAM J. Matrix Anal. Appl. **16** (1995), no. 1, 172–191.
25. M. H. Gutknecht, *A completed theory of the unsymmetric Lanczos process and related algorithms. I*, SIAM J. Matrix Anal. Appl. **13** (1992), no. 2, 594–639.
26. ———, *A completed theory of the unsymmetric Lanczos process and related algorithms. II*, SIAM J. Matrix Anal. Appl. **15** (1994), no. 1, 15–58.
27. N. Higham, *Accuracy and stability of numerical algorithms*, second ed., SIAM, 2002.
28. L. Hogben, *Elementary linear algebra*, West, St. Paul, MN, 1987.
29. H. Hotelling, *Analysis of a complex of statistical variables into principal components*, J. Educ. Psychol. **24** (1933), 417–441.
30. C. G. J. Jacobi, *Über ein leichtes Verfahren die in der Theorie der Säcularstörungen vorkommenden Gleichungen numerisch aufzulösen*, J. Reine Angew. Math. **30** (1846), 51–94.
31. S. Kaniel, *Estimates for some computational techniques in linear algebra*, Math. Comp. **20** (1966), 369–378.
32. D. Kressner, *Numerical methods for general and structured eigenvalue problems*, Lect. Notes Comput. Sci. Eng., vol. 46, Springer, Berlin, 2005.
33. V. N. Kublanovskaja, *On some algorithms for the solution of the complete eigenvalue problem*, Ž. Vyčisl. Mat. i Mat. Fiz. **1** (1961), 555–570.
34. C. Lanczos, *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*, J. Res. Nat. Bur. Stand. **45** (1950), 255–282.
35. R. S. Martin, G. Peters, and J. H. Wilkinson, *Handbook Series Linear Algebra: The QR algorithm for real Hessenberg matrices*, Numer. Math. **14** (1970), no. 3, 219–231.
36. R. B. Morgan and D. S. Scott, *Generalizations of Davidson’s method for computing eigenvalues of sparse symmetric matrices*, SIAM J. Sci. Stat. Comp. **7** (1986), no. 3, 817–825.
37. M. Overton, *Numerical computing with IEEE floating point arithmetic*, SIAM, 2001.
38. C. C. Paige, *The computation of eigenvalues and eigenvectors of very large sparse matrices*, PhD thesis, Univ. London, 1971.
39. B. Parlett, *The symmetric eigenvalue problem*, SIAM, 1998.
40. B. N. Parlett and D. S. Scott, *The Lanczos algorithm with selective orthogonalization*, Math. Comp. **33** (1979), no. 145, 217–238.
41. B. N. Parlett, D. R. Taylor, and Z. A. Liu, *A look-ahead Lanczos algorithm for unsymmetric matrices*, Math. Comp. **44** (1985), no. 169, 105–124.
42. H. Rutishauser, *Solution of eigenvalue problems with the LR-transformation*, Nat. Bur. Standards Appl. Math. Ser. (1958), no. 49, 47–81.
43. Y. Saad, *On the rates of convergence of the Lanczos and the block-Lanczos methods*, SIAM J. Numer. Anal. **17** (1980), no. 5, 687–706.
44. ———, *Variations on Arnoldi’s method for computing eigenelements of large unsymmetric matrices*, Linear Algebra Appl. **34** (1980), 269–295.
45. ———, *Chebyshev acceleration techniques for solving nonsymmetric eigenvalue problems*, Math. Comp. **42** (1984), no. 166, 567–588.
46. ———, *Iterative methods for sparse linear systems*, 2nd ed., SIAM, Philadelphia, PA, 2003.
47. ———, *Numerical methods for large eigenvalue problems*, 2nd ed., SIAM, Philadelphia, PA, 2011.
48. M. Sadkane, *A block Arnoldi-Chebyshev method for computing the leading eigenpairs of large sparse unsymmetric matrices*, Numer. Math. **64** (1993), no. 1, 181–193.
49. A. Schönhage, *Zur quadratischen Konvergenz des Jacobi-Verfahrens*, Numer. Math. **6** (1964), 410–412.
50. H. D. Simon, *Analysis of the symmetric Lanczos algorithm with reorthogonalization methods*, Linear Algebra Appl. **61** (1984), 101–131.
51. 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** (1996), no. 2, 401–425.
52. ———, *A Jacobi–Davidson iteration method for linear eigenvalue problems*, SIAM Rev. **42** (2000), no. 2, 267–293.
53. D. C. Sorensen, *Implicit application of polynomial filters in a k-step Arnoldi method*, SIAM J. Matrix Anal. Appl. **13** (1992), no. 1, 357–385.
54. ———, *Numerical methods for large eigenvalue problems*, Acta Numer. **11** (2002), 519–584.
55. G. Stewart, *Matrix algorithms: Volume II: Eigensystems*, SIAM, 2001.

56. G. W. Stewart, *Introduction to matrix computations*, Academic Press [A subsidiary of Harcourt Brace Jovanovich, Publishers], New York-London, 1973, Computer Science and Applied Mathematics.
57. G. W. Stewart and J. G. Sun, *Matrix perturbation theory*, Computer Science and Scientific Computing, Academic Press, Inc., Boston, MA, 1990.
58. H. P. M. van Kempen, *On the quadratic convergence of the special cyclic Jacobi method.*, Numer. Math. **9** (1966), 19–22.
59. A. J. Wathen, *Preconditioning*, Acta Numer. **24** (2015), 329–376.
60. D. S. Watkins, *Understanding the QR algorithm*, SIAM Rev. **24** (1982), no. 4, 427–440.
61. ———, *The transmission of shifts and shift blurring in the QR algorithm*, Linear Algebra Appl. **241** (1996), 877–896.
62. ———, *The matrix eigenvalue problem: GR and Krylov subspace methods*, SIAM, 2007.
63. ———, *The QR algorithm revisited*, SIAM Rev. **50** (2008), no. 1, 133–145.
64. ———, *Francis’s algorithm*, Amer. Math. Monthly **118** (2011), no. 5, 387–403.
65. H. Wielandt, *Das Iterationsverfahren bei nicht selbstadjungierten linearen Eigenwertaufgaben*, Math. Z. **50** (1944), 93–143.
66. J. H. Wilkinson, *The algebraic eigenvalue problem*, Monogr. Numer. Anal., Clarendon Press, Oxford, 1965.
67. K. Wu and H. Simon, *Thick-restart Lanczos method for large symmetric eigenvalue problems*, SIAM J. Matrix Anal. Appl. **22** (2000), no. 2, 602–616.

DEPARTMENT OF COMPUTER SCIENCE, THE UNIVERSITY OF BRITISH COLUMBIA, VANCOUVER, BC, V6T 1Z4, CANADA

*E-mail address:* jbosch@cs.ubc.ca, greif@cs.ubc.ca