

Chapter 6

Numerical computation of eigenvalues

Introduction

Calculating the eigenvalues and eigenvectors of a matrix is a task often encountered in scientific and engineering applications. Eigenvalue problems naturally arise in quantum physics, solid mechanics, structural engineering and molecular dynamics, to name just a few applications. The aim of this chapter is to present an overview of the standard methods for calculating eigenvalues and eigenvectors numerically. We focus predominantly on the case of a Hermitian matrix $A \in \mathbf{C}^{n \times n}$, which is technically simpler and arises in many applications. The reader is invited to go through the background material in [Appendix A.5](#) before reading this chapter. The rest of this chapter is organized as follows

- In [Section 6.1](#), we make general remarks concerning the calculation of eigenvalues.
- In [Section 6.2](#), we present standard methods based on a simple vector iteration.
- In [Section 6.3](#), we present a method for calculating several eigenvectors simultaneously, based on iterating a subspace.
- In [Section 6.4](#), we present method for constructing an approximation of the eigenvectors in a given subspace of \mathbf{C}^n .

6.1 Numerical methods for eigenvalue problems: general remarks

As mentioned in [Appendix A.5](#), a complex number $\lambda \in \mathbf{C}$ is an eigenvalue of $A \in \mathbf{C}^{n \times n}$ if and only if λ is a root of the characteristic polynomial $p_A: \mathbf{C} \rightarrow \mathbf{C}$ of A , which is given by

$$p_A(\lambda) = \det(A - \lambda I).$$

One may, therefore, calculate the eigenvalues of A by calculating the roots of the polynomial p_A using, for example, one of the methods presented in [Chapter 5](#). While feasible for small matrices, this approach is not viable for large matrices, because the number of floating point operations required for calculating calculating the coefficients of the characteristic polynomial scales as the factorial of n .

In view of the prohibitive computational cost required for calculating the characteristic polynomial, other methods are required for solving large eigenvalue problems numerically. All the methods that we study in this chapter are of iterative nature. While some of them are aimed at calculating all the eigenpairs of the matrix A , other methods enable to calculate only a small number of eigenpairs at a lower computational cost, which is often desirable. Indeed, calculating all the eigenvalues of a large matrix is computationally expensive; on a personal computer, the following Julia code takes well over a second to terminate:

```
import LinearAlgebra
A = rand(2000, 2000)
LinearAlgebra.eigen(A)
```

In many applications, the matrix A is sparse, and in this case it is important to use algorithms for eigenvalue problems that do not destroy the sparsity structure. Note that the eigenvectors of a sparse matrix are generally not sparse.

To conclude this section, we introduce some notation used throughout this chapter. For a diagonalizable matrix A , we denote the eigenvalues by $\lambda_1, \dots, \lambda_n$, with $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$. The associated normalized eigenvectors are denoted by $\mathbf{v}_1, \dots, \mathbf{v}_n$. Therefore, it holds that

$$\mathbf{V}^{-1}\mathbf{A}\mathbf{V} = \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n), \quad \text{where } \mathbf{V} = \begin{pmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_n \end{pmatrix}.$$

6.2 Simple vector iterations

In this section, we present simple iterative methods aimed at calculating just one eigenvector of the matrix A , which we assume to be diagonalizable for simplicity.

6.2.1 The power iteration

The power iteration is the simplest method for calculating the eigenpair associated with the eigenvalue of A with largest modulus. Since the eigenvectors of A span \mathbf{C}^n , any vector \mathbf{x}_0 may be decomposed as

$$\mathbf{x}_0 = \alpha_1 \mathbf{v}_1 + \dots + \alpha_n \mathbf{v}_n. \quad (6.1)$$

The idea of the power iteration is to repeatedly left-multiply this vector by the matrix A , in order to amplify the coefficient of \mathbf{v}_0 relative to the other ones. Indeed, notice that

$$\mathbf{A}^k \mathbf{x}_0 = \lambda_1^k \alpha_1 \mathbf{v}_1 + \dots + \lambda_n^k \alpha_n \mathbf{v}_n.$$

If λ_1 is strictly greater in modulus than the other eigenvalues, and if $\alpha_1 \neq 0$, then for large k the vector $\mathbf{A}^k \mathbf{x}_0$ is approximately aligned, in a sense made precise below, with the eigenvector \mathbf{v}_1 . In order to avoid overflow errors at the numerical level, the iterates are normalized at each iteration. The power iteration is presented in [Algorithm 8](#).

To precisely quantify the convergence of the power method, we introduce the notion of *acute*

Algorithm 8 Power iteration

```

 $x \leftarrow x_0$ 
for  $i \in \{1, 2, \dots\}$  do
   $x \leftarrow Ax$ 
   $x \leftarrow x/\|x\|$ 
end for

```

angle between vectors of \mathbf{C}^n .

$$\begin{aligned} \angle(\mathbf{x}, \mathbf{y}) &= \arccos\left(\frac{|\mathbf{x}^* \mathbf{y}|}{\sqrt{\mathbf{x}^* \mathbf{x}} \sqrt{\mathbf{y}^* \mathbf{y}}}\right) \\ &= \arcsin\left(\frac{\|(1 - P_{\mathbf{y}})\mathbf{x}\|}{\|\mathbf{x}\|}\right), \quad P_{\mathbf{y}} := \frac{\mathbf{y}\mathbf{y}^*}{\mathbf{y}^* \mathbf{y}}. \end{aligned}$$

This definition generalizes the familiar notion of angle for vectors in \mathbf{R}^2 or \mathbf{R}^3 , and we note that the angle function satisfies $\angle(e^{i\theta_1}\mathbf{x}, e^{i\theta_2}\mathbf{y}) = \angle(\mathbf{x}, \mathbf{y})$ as well as $\angle(\mathbf{x}, \mathbf{y}) \in [0, \pi/2]$. We can then prove the following convergence result.

Proposition 6.1 (Convergence of the power iteration). *Suppose that A is diagonalizable and that $|\lambda_1| > |\lambda_2|$. Then, for every initial guess with $\alpha_1 \neq 0$, the sequence $(\mathbf{x}_k)_{k \geq 0}$ generated by the power iteration satisfies*

$$\lim_{k \rightarrow \infty} \angle(\mathbf{x}_k, \mathbf{v}_1) = 0.$$

Proof. By construction, it holds that

$$\mathbf{x}_k = \frac{\lambda_1^k \alpha_1 \mathbf{v}_1 + \dots + \lambda_n^k \alpha_n \mathbf{v}_n}{\|\lambda_1^k \alpha_1 \mathbf{v}_1 + \dots + \lambda_n^k \alpha_n \mathbf{v}_n\|} = e^{i\theta_k} \frac{\mathbf{v}_1 + \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \mathbf{v}_2 + \dots + \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \mathbf{v}_n}{\left\| \mathbf{v}_1 + \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \mathbf{v}_2 + \dots + \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \mathbf{v}_n \right\|}, \quad (6.2)$$

where

$$e^{i\theta_k} := \frac{\lambda_1^k \alpha_1}{|\lambda_1^k \alpha_1|}.$$

It follows from (6.2) that $e^{-i\theta_k} \mathbf{x}_k \rightarrow \mathbf{v}_1 / \|\mathbf{v}_1\| = \mathbf{v}_1$ in the limit as $k \rightarrow \infty$, where we employed the fact that $\|\mathbf{v}_1\| = 1$. Using the definition of the angle between two vectors in \mathbf{C}^n , and the continuity with respect to either argument of the \mathbf{C}^n Euclidean inner product and of the arccos function, we obtain that

$$\begin{aligned} \angle(\mathbf{x}_k, \mathbf{v}_1) &= \arccos\left(\frac{|\mathbf{v}_1^* \mathbf{x}_k|}{\sqrt{\mathbf{v}_1^* \mathbf{v}_1} \sqrt{\mathbf{x}_k^* \mathbf{x}_k}}\right) = \arccos(|\mathbf{v}_1^* \mathbf{x}_k|) \\ &= \arccos\left(\left| \mathbf{v}_1^* \left(e^{-i\theta_k} \mathbf{x}_k \right) \right| \right) \xrightarrow[k \rightarrow \infty]{} \arccos(1) = 0, \end{aligned}$$

which concludes the proof. \square

An inspection of the proof also reveals that the dominant term in the error, asymptotically

in the limit as $k \rightarrow \infty$, is the one with coefficient $\frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1}$. Therefore, we deduce that

$$\angle(\mathbf{x}_k, \mathbf{v}_1) = \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right).$$

The convergence is slow if $|\lambda_2/\lambda_1|$ is close to one, and fast if $|\lambda_2| \ll |\lambda_1|$. Once an approximation of the eigenvector \mathbf{v}_1 has been calculated, the corresponding eigenvalue λ_1 can be estimated from the *Rayleigh quotient*:

$$\rho_A: \mathbf{C}_*^n \rightarrow \mathbf{C}: \mathbf{x} \mapsto \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}}. \quad (6.3)$$

For any eigenvector \mathbf{v} of \mathbf{A} , the corresponding eigenvalue is equal to $\rho_A(\mathbf{v})$. In order to study the error on the eigenvalue λ_1 for the power iteration, we assume for simplicity that \mathbf{A} is Hermitian and that the eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ are orthonormal. Substituting (6.2) in the Rayleigh quotient (6.3), we obtain

$$\rho_A(\mathbf{x}_k) = \frac{\lambda_1 + \left|\frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1}\right|^2 \lambda_2 + \dots + \left|\frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1}\right|^2 \lambda_n}{1 + \left|\frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1}\right|^2 + \dots + \left|\frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1}\right|^2}.$$

Therefore, by reducing to a common denominator we deduce

$$\begin{aligned} |\rho_A(\mathbf{x}_k) - \lambda_1| &= \left| \frac{\lambda_1 + \left|\frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1}\right|^2 \lambda_2 + \dots + \left|\frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1}\right|^2 \lambda_n}{1 + \left|\frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1}\right|^2 + \dots + \left|\frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1}\right|^2} - \lambda_1 \right| \\ &\leq \left|\frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1}\right|^2 |\lambda_2 - \lambda_1| + \dots + \left|\frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1}\right|^2 |\lambda_n - \lambda_1| = \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right). \end{aligned}$$

The convergence of the eigenvalue in the particular case of a Hermitian matrix is faster than for a general matrix in $\mathbf{C}^{n \times n}$. For general matrices, it is possible to show using a similar argument that the error is of order $\mathcal{O}(|\lambda_2/\lambda_1|^k)$ in the limit as $k \rightarrow \infty$.

Essential convergence. It is useful at this point to introduce the concept of *essential convergence*. A sequence (\mathbf{x}_k) in \mathbf{C}^n is said to *converge essentially* to a vector \mathbf{x}_∞ if there exists a sequence of complex numbers $(e^{i\phi_k})$ such that the sequence $(e^{i\phi_k} \mathbf{x}_k)$ converges to \mathbf{x}_∞ . Equivalently, the sequence (\mathbf{x}_k) converges essentially to \mathbf{x}_∞ if $\angle(\mathbf{x}_k, \mathbf{x}_\infty)$ converges to 0. Proving this equivalence is the goal of [Exercise 6.11](#). Reformulated in this new terminology, [Proposition 6.1](#) states that the sequence (\mathbf{x}_k) obtained from the power iteration converges essentially to \mathbf{v}_1 .

6.2.2 Inverse iteration

The power iteration is simple but enables to calculate only the dominant eigenvalue of the matrix \mathbf{A} , i.e. the eigenvalue of largest modulus. In addition, the convergence of the method is slow when $|\lambda_2| \approx |\lambda_1|$.

The inverse iteration enables a more efficient calculation of not only the dominant eigenvalue but also the other eigenvalues of \mathbf{A} . It is based on applying the power iteration to $(\mathbf{A} - \mu \mathbf{I})^{-1}$,

where $\mu \in \mathbf{C}$ is a shift. The eigenvalues of $(\mathbf{A} - \mu\mathbf{I})^{-1}$ are given by $(\lambda_1 - \mu)^{-1}, \dots, (\lambda_n - \mu)^{-1}$, with associated eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_n$. If $0 < |\lambda_J - \mu| < |\lambda_j - \mu|$ for all $j \neq J$, then the dominant eigenvalue of the matrix $(\mathbf{A} - \mu\mathbf{I})^{-1}$ is $(\lambda_J - \mu)^{-1}$, and so the power iteration applied to this matrix yields an approximation of the eigenvector \mathbf{v}_J . In other words, the inverse iteration with shift μ enables to calculate an approximation of the eigenvector of \mathbf{A} corresponding to the eigenvalue nearest μ . The inverse iteration is presented in [Algorithm 9](#). In practice, the inverse matrix $(\mathbf{A} - \mu\mathbf{I})^{-1}$ need not be calculated, and it is often preferable to solve a linear system at each iteration.

Algorithm 9 Inverse iteration

```

 $\mathbf{x} \leftarrow \mathbf{x}_0$ 
for  $i \in \{1, 2, \dots\}$  do
  Solve  $(\mathbf{A} - \mu\mathbf{I})\mathbf{y} = \mathbf{x}$ 
   $\mathbf{x} \leftarrow \mathbf{y}/\|\mathbf{y}\|$ 
end for
 $\lambda \leftarrow \mathbf{x}^* \mathbf{A} \mathbf{x} / \mathbf{x}^* \mathbf{x}$ 
return  $\mathbf{x}, \lambda$ 

```

An application of [Proposition 6.1](#) immediately gives the following convergence result for the inverse iteration.

Proposition 6.2 (Convergence of the inverse iteration). *Assume that $\mathbf{A} \in \mathbf{C}^n$ is diagonalizable and that there exist J and K such that*

$$0 < |\lambda_J - \mu| < |\lambda_K - \mu| \leq |\lambda_j - \mu| \quad \forall j \neq J.$$

Assume also that $\alpha_J \neq 0$, where α_J is the coefficient of \mathbf{v}_J in the expansion of \mathbf{x}_0 given in (6.1). Then the iterates of the inverse iteration satisfy

$$\lim_{k \rightarrow \infty} \angle(\mathbf{x}_k, \mathbf{v}_J) = 0.$$

More precisely,

$$\angle(\mathbf{x}_k, \mathbf{v}_J) = \mathcal{O}\left(\left|\frac{\lambda_J - \mu}{\lambda_K - \mu}\right|^k\right).$$

[Proposition 6.2](#) states that \mathbf{x}_k converges essentially to \mathbf{v}_J . Notice that the closer μ is to λ_J , the faster the inverse iteration converges. Note also that with $\mu = 0$, the inverse iteration enables to calculate the eigenvalue of \mathbf{A} of smallest modulus.

6.2.3 Rayleigh quotient iteration

Since the inverse iteration is fast when μ is close to an eigenvalue λ_J , it is natural to wonder whether the method can be improved by progressively updating μ as the simulation progresses. Specifically, an approximation of the eigenvalue associated with the current vector may be employed in place of μ . This leads to the Rayleigh quotient iteration, presented in [Algorithm 10](#).

Algorithm 10 Inverse iteration

```

 $x \leftarrow x_0$ 
for  $i \in \{1, 2, \dots\}$  do
   $\mu \leftarrow x^*Ax/x^*x$ 
  Solve  $(A - \mu I)y = x$ 
   $x \leftarrow y/\|y\|$ 
end for
 $\lambda \leftarrow x^*Ax/x^*x$ 
return  $x, \lambda$ 

```

It is possible to show that, when A is Hermitian, the Rayleigh quotient iteration converges to an eigenvector for almost every initial guess x_0 . Furthermore, if convergence to an eigenvector occurs, then μ converges cubically to the corresponding eigenvalue. See [11] and the references therein for more details.

6.3 Methods based on a subspace iteration

The subspace iteration resembles the power iteration but it is more general: not just one but several vectors are updated at each iteration.

6.3.1 Simultaneous iteration

Let $X_0 = (x_1 \ \dots \ x_p)$ denote an initial set of linearly independent vectors. Before we present the simultaneous iteration, we recall a statement concerning the QR decomposition of a matrix, which is related to the Gram–Schmidt orthonormalization process. We recall that the Gram–Schmidt method enables to construct, starting from an ordered set of vectors $\{x_1, \dots, x_p\}$ in \mathbf{C}^n , a new set of vectors $\{q_1, \dots, q_p\}$ which are *orthonormal* and span the same subspace of \mathbf{C}^n as the original vectors.

Proposition 6.3 (Reduced QR decomposition). *Assume that $X \in \mathbf{C}^{n \times p}$ has linearly independent columns. Then there exist a matrix $Q \in \mathbf{C}^{n \times p}$ with orthonormal columns and an upper triangular matrix $R \in \mathbf{C}^{p \times p}$ such that the following factorization holds:*

$$X = QR. \tag{6.4}$$

This decomposition is known as a reduced QR decomposition if $p < n$, or simply QR decomposition if $p = n$, in which case X is a square matrix and Q is a unitary matrix. The decomposition is unique if we require that the diagonal elements of R are real and positive.

Proof. The statement is clear when $p = 1$. Reasoning by induction, we assume that the result is true up to $p - 1$, and prove that it then also holds true for p . We wish to show that there is a unique matrix $Q \in \mathbf{C}^{n \times p}$ with orthonormal columns and a unique upper triangular matrix $R \in \mathbf{C}^{p \times p}$ with real and positive diagonal elements such that (6.4) is satisfied. To this

end, we decompose the matrices \mathbf{Q} and \mathbf{R} as follows:

$$\mathbf{Q} = \begin{pmatrix} \mathbf{Q}_{p-1} & \mathbf{q} \end{pmatrix}, \quad \mathbf{R} = \begin{pmatrix} \mathbf{R}_{p-1} & \mathbf{r} \\ \mathbf{0}_{p-1}^T & r \end{pmatrix}. \quad (6.5)$$

Here $\mathbf{Q}_{p-1} \in \mathbf{C}^{n \times (p-1)}$ is a matrix with orthonormal columns, $\mathbf{R} \in \mathbf{C}^{(p-1) \times (p-1)}$ is an upper triangular matrix with positive real diagonal elements, $\mathbf{q} \in \mathbf{C}^n$ is a normalized vector orthogonal to all the columns of \mathbf{Q}_{p-1} , $\mathbf{r} \in \mathbf{C}^{p-1}$ is a vector and $r \in \mathbf{R}_{>0}$ is a scalar. Let us also denote by $\mathbf{X}_{p-1} \in \mathbf{C}^{n \times (p-1)}$ the matrix containing the $p-1$ first columns of \mathbf{X} , and by $\mathbf{x}_p \in \mathbf{C}^n$ the p -th column of \mathbf{X} . Substituting (6.5) into (6.4), we then obtain

$$\begin{pmatrix} \mathbf{X}_{p-1} & \mathbf{x}_p \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_{p-1} \mathbf{R}_{p-1} & \mathbf{Q}_{p-1} \mathbf{r} + \mathbf{q} r \end{pmatrix}, \quad (6.6)$$

By the induction hypothesis, there exist a unique choice of matrices \mathbf{Q}_{p-1} and \mathbf{R}_{p-1} with the required structure such that $\mathbf{X}_{p-1} = \mathbf{Q}_{p-1} \mathbf{R}_{p-1}$. Comparing the last column of both sides in (6.6), we obtain

$$\mathbf{x}_p = \mathbf{Q}_{p-1} \mathbf{r} + \mathbf{q} r. \quad (6.7)$$

Left-multiplying both sides by \mathbf{Q}_{p-1}^* and employing the orthogonality between \mathbf{q} and the columns of \mathbf{Q}_{p-1} , we deduce that necessarily $\mathbf{r} = \mathbf{Q}_{p-1}^* \mathbf{x}_p$. It then follows from (6.7) that

$$\mathbf{q} = \frac{1}{r} (\mathbf{x}_p - \mathbf{Q}_{p-1} \mathbf{Q}_{p-1}^* \mathbf{x}_p), \quad r = \|\mathbf{x}_p - \mathbf{Q}_{p-1} \mathbf{Q}_{p-1}^* \mathbf{x}_p\|.$$

It is simple to check that \mathbf{q} is indeed orthogonal to the columns of \mathbf{Q} , which concludes the proof. Note that $\mathbf{Q}_{p-1} \mathbf{Q}_{p-1}^* \mathbf{x}_p$ is the orthogonal projection of \mathbf{x}_p onto the subspace spanned by the column of \mathbf{Q}_{p-1} . \square

Note that the columns of the matrix \mathbf{Q} of the decomposition coincide with the vectors that would be obtained by applying the Gram–Schmidt method to the columns of the matrix \mathbf{X} . In fact, the Gram–Schmidt process is one of several methods by which the QR decomposition can be calculated in practice.

Algorithm 11 Simultaneous iteration

```

 $\mathbf{X} \leftarrow \mathbf{X}_0$ 
for  $k \in \{1, 2, \dots\}$  do
     $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A} \mathbf{X}_{k-1}$  (QR decomposition).
     $\mathbf{X}_k \leftarrow \mathbf{Q}_k$ .
end for

```

The simultaneous iteration method is presented in [Algorithm 11](#). Like the normalization in the power iteration [Algorithm 8](#), the QR decomposition performed at each step in [Algorithm 11](#) enables to avoid overflow errors. Notice that when $p = 1$, the simultaneous iteration reduces to the power iteration. We emphasize that the factorization step at each iteration does not influence the subspace spanned by the columns of \mathbf{X} . Therefore, this subspace after k iterations coincides with that spanned by the columns of the matrix $\mathbf{A}^k \mathbf{X}_0$. In fact, in exact arithmetic, it would be equivalent to perform the QR decomposition only once as a final step, after the **for**

loop. Indeed, denoting by $Q_k R_k$ the QR decomposition of $A X_{k-1}$, we have

$$\begin{aligned} X_k &= A X_{k-1} R_k^{-1} = A^2 X_{k-2} R_{k-1}^{-1} R_k^{-1} = \dots = A^k X_0 R_1^{-1} \dots R_k^{-1} \\ &\Leftrightarrow X_k (R_k \dots R_1) = A^k X_0. \end{aligned}$$

Since X_k has orthonormal columns and $R_k \dots R_1$ is an upper triangular matrix (see [Exercise 4.3](#)) with real positive elements on the diagonal (check this!), it follows that X_k can be obtained by QR factorization of $A^k X_0$. In order to show the convergence of the simultaneous iteration, we begin by proving the following preparatory lemma.

Lemma 6.4 (Continuity of the QR decomposition). *If $Q_k R_k \rightarrow QR$, where Q is orthogonal and R is upper triangular with positive real entries on the diagonal, then $Q_k \rightarrow Q$.*

Proof. We reason by contradiction and assume there is $\varepsilon > 0$ and a subsequence $(Q_{k_n})_{n \geq 0}$ such that $\|Q_{k_n} - Q\| \geq \varepsilon$ for all n . Since the set of unitary matrices is a compact subset of $\mathbf{C}^{n \times n}$, there exists a further subsequence $(Q_{k_{nm}})_{m \geq 0}$ that converges to a limit Q_∞ which is also a unitary matrix and at least ε away in norm from Q . But then

$$R_{k_{nm}} = Q_{k_{nm}}^{-1} (Q_{k_{nm}} R_{k_{nm}}) = Q_{k_{nm}}^* (Q_{k_{nm}} R_{k_{nm}}) \xrightarrow{m \rightarrow \infty} Q_\infty^* (QR) =: R_\infty.$$

Since R_k is upper triangular with positive diagonal elements for all k , clearly R_∞ is also upper triangular with positive diagonal elements. But then $Q_\infty R_\infty = QR$, and by uniqueness of the decomposition we deduce that $Q = Q_\infty$, which is a contradiction. \square

Before presenting the convergence theorem, we introduce the following terminology: we say that $X_k \in \mathbf{C}^{n \times p}$ converges essentially to a matrix X_∞ if each column of X_k converges essentially to the corresponding column of X_∞ . We prove the convergence in the Hermitian case for simplicity. In the general case of $A \in \mathbf{C}^{n \times n}$, it cannot be expected that X_k converges essentially to V , because the columns of X_k are orthogonal but eigenvectors may not be orthogonal. In this case, the columns of X_k converge not to the eigenvectors but to the so-called Schur vectors of A ; see [11] for more information.

Theorem 6.5 (Convergence of the simultaneous iteration ⓐ). *Assume that $A \in \mathbf{C}^{n \times n}$ is Hermitian, that $X_0 \in \mathbf{C}^{n \times p}$ has linearly independent columns, and finally that the subspace spanned by the column of X_0 satisfies*

$$\text{col}(X_0) \cap \text{Span}\{\mathbf{v}_{p+1}, \dots, \mathbf{v}_n\} = \emptyset. \quad (6.8)$$

If it holds that

$$\lambda_1 > \lambda_2 > \dots > \lambda_p > \lambda_{p+1} \geq \lambda_{p+2} \geq \dots \geq \lambda_n, \quad (6.9)$$

then X_k converges essentially to $V_1 := \begin{pmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_p \end{pmatrix}$.

Proof. Let $\mathbf{B} = \mathbf{V}^{-1}\mathbf{X}_0 \in \mathbf{C}^{n \times p}$, so that $\mathbf{X}_0 = \mathbf{V}\mathbf{B}$, and note that $\mathbf{A}^k\mathbf{X}_0 = \mathbf{V}\mathbf{D}^k\mathbf{B}$. We denote by $\mathbf{B}_1 \in \mathbf{C}^{p \times p}$ and $\mathbf{B}_2 \in \mathbf{C}^{(n-p) \times p}$ the upper $p \times p$ and lower $(n-p) \times p$ blocks of \mathbf{B} , respectively. The matrix \mathbf{B}_1 is nonsingular, otherwise the assumption (6.8) would not hold. Indeed, if there was a nonzero vector $\mathbf{z} \in \mathbf{C}^p$ such that $\mathbf{B}_1\mathbf{z} = 0$, then

$$\mathbf{X}_0\mathbf{z} = \mathbf{V} \begin{pmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{pmatrix} \mathbf{z} = \begin{pmatrix} \mathbf{V}_1 & \mathbf{V}_2 \end{pmatrix} \begin{pmatrix} \mathbf{0} \\ \mathbf{B}_2\mathbf{z} \end{pmatrix} = \mathbf{V}_2\mathbf{B}_2\mathbf{z}.$$

implying that $\mathbf{X}_0\mathbf{z} \in \text{col}(\mathbf{X}_0)$ is a linear combination of the vectors in $\mathbf{V}_2 = (\mathbf{v}_{p+1} \ \dots \ \mathbf{v}_n)$, which contradicts the assumption. We also denote by \mathbf{D}_1 and \mathbf{D}_2 the $p \times p$ upper-left and the $(n-p) \times (n-p)$ lower-right blocks of \mathbf{D} , respectively. From the expression of $\mathbf{A}^k\mathbf{X}_0$, we have

$$\begin{aligned} \mathbf{A}^k\mathbf{X}_0 &= \begin{pmatrix} \mathbf{V}_1 & \mathbf{V}_2 \end{pmatrix} \begin{pmatrix} \mathbf{D}_1^k & \\ & \mathbf{D}_2^k \end{pmatrix} \begin{pmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{pmatrix} = \mathbf{V}_1\mathbf{D}_1^k\mathbf{B}_1 + \mathbf{V}_2\mathbf{D}_2^k\mathbf{B}_2, \\ &= \left(\mathbf{V}_1 + \mathbf{V}_2\mathbf{D}_2^k\mathbf{B}_2\mathbf{B}_1^{-1}\mathbf{D}_1^{-k} \right) \mathbf{D}_1^k\mathbf{B}_1. \end{aligned} \quad (6.10)$$

The second term in the bracket on the right-hand side converges to zero in the limit as $k \rightarrow \infty$ by (6.9). Let $\widetilde{\mathbf{Q}}_k\widetilde{\mathbf{R}}_k$ denote the reduced QR decomposition of the bracketed term. By Lemma 6.4, which we proved for the standard QR decomposition but also holds for the reduced one, we deduce from $\widetilde{\mathbf{Q}}_k\widetilde{\mathbf{R}}_k \rightarrow \mathbf{V}_1$ that $\widetilde{\mathbf{Q}}_k \rightarrow \mathbf{V}_1$. Rearranging (6.10), we have

$$\mathbf{A}^k\mathbf{X}_0 = \widetilde{\mathbf{Q}}_k(\widetilde{\mathbf{R}}_k\mathbf{D}_1^k\mathbf{B}_1).$$

Since the matrix between brackets is a $p \times p$ square invertible matrix, this equation implies that $\text{col}(\mathbf{A}^k\mathbf{X}_0) = \text{col}(\widetilde{\mathbf{Q}}_k)$. Denoting by $\mathbf{Q}_k\mathbf{R}_k$ the QR decomposition of $\mathbf{A}_k\mathbf{X}_0$, we therefore have $\text{col}(\mathbf{Q}_k) = \text{col}(\widetilde{\mathbf{Q}}_k)$, and so the projectors on these subspaces are equal. We recall that, for a set of orthonormal vectors $\mathbf{r}_1, \dots, \mathbf{r}_p$ gathered in a matrix $\mathbf{R} = (\mathbf{r}_1 \ \dots \ \mathbf{r}_p)$, the projector on $\text{col}(\mathbf{R}) = \text{Span}\{\mathbf{r}_1, \dots, \mathbf{r}_p\} \subset \mathbf{C}^n$ is the square $n \times n$ matrix

$$\mathbf{R}\mathbf{R}^* = \mathbf{r}_1\mathbf{r}_1^* + \dots + \mathbf{r}_p\mathbf{r}_p^*.$$

Consequently, the equality of the projectors implies $\mathbf{Q}_k\mathbf{Q}_k^* = \widetilde{\mathbf{Q}}_k\widetilde{\mathbf{Q}}_k^*$. Now, we want to establish the essential convergence of \mathbf{Q}_k to \mathbf{V}_1 . To this end, we reason by induction, relying on the fact that the first k columns of \mathbf{X}_0 undergo a simultaneous iteration independent of the other columns. For example, the first column simply undergoes a power iteration, and so it converges essentially to \mathbf{v}_1 . Assume now that the columns 1 to $p-1$ of \mathbf{Q}_k converge essentially to $\mathbf{v}_1, \dots, \mathbf{v}_{p-1}$. Then the p -th column of \mathbf{Q}_k at iteration k , which we denote by $\mathbf{q}_p^{(k)}$, satisfies

$$\begin{aligned} \mathbf{q}_p^{(k)}\mathbf{q}_p^{(k)*} &= \mathbf{Q}_k\mathbf{Q}_k^* - \mathbf{q}_1^{(k)}\mathbf{q}_1^{(k)*} - \dots - \mathbf{q}_{p-1}^{(k)}\mathbf{q}_{p-1}^{(k)*} = \widetilde{\mathbf{Q}}_k\widetilde{\mathbf{Q}}_k^* - \mathbf{q}_1^{(k)}\mathbf{q}_1^{(k)*} - \dots - \mathbf{q}_{p-1}^{(k)}\mathbf{q}_{p-1}^{(k)*} \\ &\xrightarrow[k \rightarrow \infty]{} \mathbf{V}_1\mathbf{V}_1^* - \mathbf{v}_1\mathbf{v}_1^* - \dots - \mathbf{v}_{p-1}\mathbf{v}_{p-1}^* = \mathbf{v}_p\mathbf{v}_p^*. \end{aligned}$$

Therefore, noting that $|a| = \sqrt{a\bar{a}}$ for every $a \in \mathbf{C}$, we deduce

$$|\mathbf{v}_p^* \mathbf{q}_p^{(k)}| = \sqrt{\mathbf{v}_p^* \mathbf{q}_p^{(k)} \mathbf{q}_p^{(k)*} \mathbf{v}_p} \xrightarrow{k \rightarrow \infty} \sqrt{\mathbf{v}_p^* \mathbf{v}_p \mathbf{v}_p^* \mathbf{v}_p} = 1,$$

which shows that $\angle(\mathbf{q}_p^{(k)}, \mathbf{v}_p)$ converges to 0. Finally, observing that

$$\left\| e^{-i\theta_k} \mathbf{q}_p^{(k)} - \mathbf{v}_p \right\|^2 = 2 - 2|\mathbf{v}_p^* \mathbf{q}_p^{(k)}| \xrightarrow{k \rightarrow \infty} 0, \quad \theta_k = \frac{\mathbf{v}_p^* \mathbf{q}_p^{(k)}}{|\mathbf{v}_p^* \mathbf{q}_p^{(k)}|},$$

we conclude that $\mathbf{q}_p^{(k)}$ converges essentially to \mathbf{v}_p . □

In addition to this convergence result, it is possible to show that the error satisfies

$$\angle(\text{col}(\mathbf{X}_k), \text{col}(\mathbf{V}_1)) = \mathcal{O}\left(\left|\frac{\lambda_{p+1}}{\lambda_p}\right|^k\right).$$

Here, the angle between two subspaces \mathcal{A} and \mathcal{B} of \mathbf{C}^n is defined as

$$\angle(\mathcal{A}, \mathcal{B}) = \max_{\mathbf{a} \in \mathcal{A} \setminus \{\mathbf{0}\}} \left(\min_{\mathbf{b} \in \mathcal{B} \setminus \{\mathbf{0}\}} \angle(\mathbf{a}, \mathbf{b}) \right).$$

6.3.2 The QR algorithm

The QR algorithm, which is based on the QR decomposition, is one of the most famous algorithms for calculating *all* the eigenpairs of a matrix. We first present the algorithm and then relate it to the simultaneous iteration in [Section 6.3.1](#). The method is presented in [Algorithm 12](#).

Algorithm 12 QR algorithm

```

X0 = A
for  $i \in \{1, 2, \dots\}$  do
    QkRk = Xk-1 (QR decomposition)
    Xk = RkQk
end for
    
```

Successive iterates of the QR algorithm are related by the equation

$$\mathbf{X}_k = \mathbf{Q}_k^{-1} \mathbf{X}_{k-1} \mathbf{Q}_k = \mathbf{Q}_k^* \mathbf{X}_{k-1} \mathbf{Q}_k = \dots = (\mathbf{Q}_1 \dots \mathbf{Q}_k)^* \mathbf{X}_0 (\mathbf{Q}_1 \dots \mathbf{Q}_k) \quad (6.11)$$

Therefore, all the iterates are related by a unitary similarity transformation, and so they all have the same eigenvalues as $\mathbf{X}_0 = \mathbf{A}$. Rearranging (6.11), we have

$$(\mathbf{Q}_1 \dots \mathbf{Q}_k) \mathbf{X}_k = \mathbf{A} (\mathbf{Q}_1 \dots \mathbf{Q}_k),$$

and so, introducing $\tilde{\mathbf{Q}}_k = \mathbf{Q}_1 \dots \mathbf{Q}_k$ and noting that $\mathbf{X}_k = \mathbf{Q}_{k+1} \mathbf{R}_{k+1}$ by the algorithm, we deduce

$$\tilde{\mathbf{Q}}_{k+1} \mathbf{R}_{k+1} = \mathbf{A} \tilde{\mathbf{Q}}_k.$$

This reveals that the matrix sequence $(\tilde{Q}_k)_{k \geq 1}$ undergoes a simultaneous iteration and so, assuming that A is Hermitian with n distinct nonzero eigenvalues, we deduce that $\tilde{Q}_k \rightarrow V$ essentially in the limit as $k \rightarrow \infty$, by [Theorem 6.5](#). As a consequence, by [\(6.11\)](#), it holds that $X_k \rightarrow V^* X_0 V = D$; in other words, the matrix X_k converges to a diagonal matrix with the eigenvalues of A on the diagonal.

6.4 Projection methods

In this section, we begin by presenting a general method for constructing an approximation of the eigenvectors of A in a given subspace \mathcal{U} of \mathbf{C}^n . We then discuss a particular choice for the subspace \mathcal{U} as a Krylov subspace, which is very useful in practice.

Assume that $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ is an orthonormal basis of \mathcal{U} . Then for any vector $\mathbf{v} \in \mathbf{C}^n$, the vector of \mathcal{U} that is closest to \mathbf{v} in the Euclidean distance is given by the orthogonal projection

$$P_{\mathcal{U}}\mathbf{v} := \mathbf{U}\mathbf{U}^*\mathbf{v} = (\mathbf{u}_1\mathbf{u}_1^* + \dots + \mathbf{u}_p\mathbf{u}_p^*)\mathbf{v}.$$

In practice, the eigenvectors of A are unknown, and so it is impossible to calculate approximations using this formula. The Rayleigh–Ritz method, which we present hereafter, is an alternative and practical method for constructing approximations of the eigenvectors and eigenvalues. In general, the subspace \mathcal{U} does not contain any eigenvector of A , and so the problem

$$A\mathbf{v} = \lambda\mathbf{v}, \quad \mathbf{v} \in \mathcal{U} \tag{6.12}$$

does not admit a solution. Let us denote by \mathbf{U} the matrix with columns $\mathbf{u}_1, \dots, \mathbf{u}_p$. Since any vector $\mathbf{v} \in \mathcal{U}$ is equal to $\mathbf{U}\mathbf{z}$ for some vector $\mathbf{z} \in \mathbf{C}^p$, equation [\(6.12\)](#) is equivalent to the problem

$$A\mathbf{U}\mathbf{z} = \lambda\mathbf{U}\mathbf{z},$$

which is a system of n equations with $p < n$ unknowns. The Rayleigh–Ritz method is based on the idea that, in order to obtain a problem with as many unknowns as there are equations, we can multiply this equation by \mathbf{U}^* , which leads to the problem

$$B\mathbf{z} := (\mathbf{U}^*A\mathbf{U})\mathbf{z} = \lambda\mathbf{z}. \tag{6.13}$$

This is standard eigenvalue problem for the matrix $\mathbf{U}^*A\mathbf{U} \in \mathbf{C}^{p \times p}$, which is much easier to solve than the original problem if $p \ll n$. Equivalently, equation [\(6.13\)](#) may be formulated as follows: find $\mathbf{v} \in \mathcal{U}$ such that

$$\mathbf{u}^*(A\mathbf{v} - \lambda\mathbf{v}), \quad \forall \mathbf{u} \in \mathcal{U}. \tag{6.14}$$

The solutions to [\(6.13\)](#) and [\(6.14\)](#) are related by the equation $\mathbf{v} = \mathbf{U}\mathbf{z}$. Of course, the eigenvalues of B in problem [\(6.13\)](#), which are called the Ritz values of A relative to \mathcal{U} , are in general different from those of A . Once an eigenvector \mathbf{y} of B has been calculated, an approximate eigenvector of A , called a *Ritz vector* of A relative to \mathcal{U} , is obtained from the equation $\hat{\mathbf{v}} = \mathbf{U}\mathbf{y}$. The Rayleigh–Ritz algorithm is presented in full in [Algorithm 13](#).

Algorithm 13 Rayleigh–Ritz

 Choose $\mathcal{U} \subset \mathbf{C}^n$

 Construct a matrix \mathbf{U} whose columns are orthonormal and span \mathcal{U}

 Find the eigenvalues $\hat{\lambda}_i$ and eigenvectors $\mathbf{y}_i \in \mathbf{C}^p$ of $\mathbf{B} := \mathbf{U}^* \mathbf{A} \mathbf{U}$

 Calculate the corresponding Ritz vectors $\hat{\mathbf{v}}_i = \mathbf{U} \mathbf{y}_i \in \mathbf{C}^n$.

It is clear that if $\mathbf{v}_i \in \mathcal{U}$, then λ_i is an eigenvalue of \mathbf{B} in (6.13). In fact, we can show the following more general statement.

Proposition 6.6. *If \mathcal{U} is an invariant subspace of \mathbf{A} , meaning that $\mathbf{A}\mathcal{U} \subset \mathcal{U}$, then each Ritz vector of \mathbf{A} relative to \mathcal{U} is an eigenvector of \mathbf{A} .*

Proof. Let $\mathbf{U} \in \mathbf{C}^{n \times p}$ and $\mathbf{W} \in \mathbf{C}^{n \times (n-p)}$ be matrices whose columns form orthonormal bases of \mathcal{U} and \mathcal{U}^\perp , respectively. Here \mathcal{U}^\perp denotes the orthogonal complement of \mathcal{U} with respect to the Euclidean inner product. Then, since $\mathbf{W}^* \mathbf{A} \mathbf{U} = \mathbf{0}$ by assumption, it holds that

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} \mathbf{U}^* \mathbf{A} \mathbf{U} & \mathbf{U}^* \mathbf{A} \mathbf{W} \\ \mathbf{W}^* \mathbf{A} \mathbf{U} & \mathbf{W}^* \mathbf{A} \mathbf{W} \end{pmatrix} = \begin{pmatrix} \mathbf{U}^* \mathbf{A} \mathbf{U} & \mathbf{U}^* \mathbf{A} \mathbf{W} \\ \mathbf{0} & \mathbf{W}^* \mathbf{A} \mathbf{W} \end{pmatrix}, \quad \mathbf{Q} = \begin{pmatrix} \mathbf{U} & \mathbf{W} \end{pmatrix}.$$

If $(\mathbf{y}, \hat{\lambda})$ is an eigenvector of $\mathbf{U}^* \mathbf{A} \mathbf{U}$, then

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} (\mathbf{U}^* \mathbf{A} \mathbf{U}) \mathbf{y} \\ \mathbf{0} \end{pmatrix} = \hat{\lambda} \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} =: \hat{\lambda} \mathbf{x},$$

and so $(\mathbf{x}, \hat{\lambda})$ is an eigenpair of $\mathbf{Q}^* \mathbf{A} \mathbf{Q}$. But then $(\mathbf{Q} \mathbf{x}, \hat{\lambda}) = (\mathbf{U} \mathbf{y}, \hat{\lambda})$ is an eigenpair of \mathbf{A} , which proves the statement. \square

If \mathcal{U} is close to being an invariant subspace of \mathbf{A} , then it is expected that the Ritz vectors and Ritz values of \mathbf{A} relative to \mathcal{U} will provide good approximations of some of the eigenpairs of \mathbf{A} . Quantifying this approximation is difficult, so we only present without proof the following error bound. See [10] for more information.

Proposition 6.7. *Let \mathbf{A} be a full rank Hermitian matrix and \mathcal{U} a p -dimensional subspace of \mathbf{C}^n . Then there exists eigenvalues $\lambda_{i_1}, \dots, \lambda_{i_p}$ of \mathbf{A} which satisfy*

$$\forall j \in \{1, \dots, p\}, \quad |\lambda_{i_j} - \hat{\lambda}_j| \leq \|(\mathbf{I} - \mathbf{P}_{\mathcal{U}}) \mathbf{A} \mathbf{P}_{\mathcal{U}}\|_2.$$

In the case where \mathbf{A} is Hermitian, it is possible to show that the Ritz values are bounded from above by the eigenvalues of \mathbf{A} . The proof of this result relies on the Courant–Fisher theorem for characterizing the eigenvalues of a Hermitian matrix, which is recalled in Theorem A.5 in the appendix.

Proposition 6.8. *If $A \in \mathbf{C}^{n \times n}$ is Hermitian, then*

$$\forall i \in \{1, \dots, p\}, \quad \widehat{\lambda}_i \leq \lambda_i$$

Proof. By the Courant–Fisher theorem, it holds that

$$\widehat{\lambda}_i = \max_{S \subset \mathbf{C}^p, \dim(S)=i} \left(\min_{\mathbf{x} \in S \setminus \{0\}} \frac{\mathbf{x}^* \mathbf{B} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} \right)$$

Letting $\mathbf{y} = \mathbf{U} \mathbf{x}$ and then $\mathcal{R} = \mathbf{U} S$, we deduce that

$$\begin{aligned} \widehat{\lambda}_i &= \max_{S \subset \mathbf{C}^p, \dim(S)=i} \left(\min_{\mathbf{y} \in \mathbf{U} S \setminus \{0\}} \frac{\mathbf{y}^* \mathbf{A} \mathbf{y}}{\mathbf{y}^* \mathbf{y}} \right) \\ &= \max_{\mathcal{R} \subset \mathcal{U}, \dim(\mathcal{R})=i} \left(\min_{\mathbf{y} \in \mathcal{R} \setminus \{0\}} \frac{\mathbf{y}^* \mathbf{A} \mathbf{y}}{\mathbf{y}^* \mathbf{y}} \right) \leq \max_{\mathcal{R} \subset \mathbf{C}^n, \dim(\mathcal{R})=i} \left(\min_{\mathbf{y} \in \mathcal{R} \setminus \{0\}} \frac{\mathbf{y}^* \mathbf{A} \mathbf{y}}{\mathbf{y}^* \mathbf{y}} \right) = \lambda_i, \end{aligned}$$

where we used the Courant–Fisher for the matrix A in the last equality. \square

This projection approach is sometimes combined with a simultaneous subspace iteration: an approximation \mathbf{X}_k of the p first eigenvector is first calculated using [Algorithm 11](#), and then the matrix \mathbf{X}_k is used in place of \mathbf{U} in [Algorithm 13](#).

6.4.1 Projection method in a Krylov subspace

The power iteration constructs at iteration k an approximation of \mathbf{v}_1 in the one-dimensional subspace spanned by the vector $A^k \mathbf{x}_0$, and only the previous iteration \mathbf{x}_k is employed to construct \mathbf{x}^{k+1} . One may wonder whether, by employing all the previous iterates rather than only the previous one, a better approximation of \mathbf{v}_1 can be constructed. More precisely, instead of looking for an approximation in the subspace $\text{Span}\{A^k \mathbf{x}_0\}$, would it be useful to extend the search area to the Krylov subspace

$$\mathcal{K}_{k+1}(A, \mathbf{x}_0) := \text{Span}\{\mathbf{x}_0, A\mathbf{x}_0, \dots, A^k \mathbf{x}_0\}?$$

The answer to this question is positive, and the resulting method is often much faster than the power iteration. This is achieved by employing the Rayleigh–Ritz projection method [Algorithm 13](#) with the choice $\mathcal{U} = \mathcal{K}_{k+1}(A, \mathbf{x}_0)$. Applying this method requires to calculate an orthonormal basis of the Krylov subspace and to calculate the reduced matrix $\mathbf{U}^* \mathbf{A} \mathbf{U}$. The *Arnoldi method* enables to achieve these two goals simultaneously.

6.4.2 The Arnoldi iteration

This Arnoldi iteration is based on the Gram–Schmidt process and presented in [Algorithm 14](#). The iteration breaks down if $h_{j+1,j} = 0$, which indicates that $\mathcal{A} \mathbf{u}_j$ belongs to the Krylov subspace $\text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_j\} = \mathcal{K}_j(A, \mathbf{u}_1)$, implying that $\mathcal{K}_{j+1}(A, \mathbf{u}_1) = \mathcal{K}_j(A, \mathbf{u}_1)$. In this case, the subspace $\mathcal{K}_j(A, \mathbf{u}_1)$ is an invariant subspace of A because, by [Exercise 6.2](#), we have

$$A \mathcal{K}_j(A, \mathbf{u}_1) \subset \mathcal{K}_{j+1}(A, \mathbf{u}_1) = \mathcal{K}_j(A, \mathbf{u}_1).$$

Algorithm 14 Arnoldi iteration for constructing an orthonormal basis of $\mathcal{K}_p(\mathbf{A}, \mathbf{u}_1)$

Choose \mathbf{u}_1 with unit norm.
for $j \in \{1, \dots, p\}$ **do**
 $\mathbf{u}_{j+1} \leftarrow \mathbf{A}\mathbf{u}_j$
 for $i \in \{1, \dots, j\}$ **do**
 $h_{i,j} \leftarrow \mathbf{u}_i^* \mathbf{u}_{j+1}$
 $\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1} - h_{i,j} \mathbf{u}_i$
 end for
 $h_{j+1,j} \leftarrow \|\mathbf{u}_{j+1}\|$
 $\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1}/h_{j+1,j}$
end for

Therefore, applying the Rayleigh–Ritz with $\mathcal{U} = \text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_j\}$ yields exact eigenpairs [Proposition 6.6](#). If the iteration does not break down then, by construction, the vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ at the end of the algorithm are orthonormal. It is also simple to show by induction that they form a basis of $\mathcal{K}_p(\mathbf{A}, \mathbf{u}_1)$. The scalar coefficients $h_{i,j}$ can be collected in a matrix square $p \times p$ matrix

$$\mathbf{H} = \begin{pmatrix} h_{1,1} & h_{1,2} & h_{1,3} & \cdots & h_{1,p} \\ h_{2,1} & h_{2,2} & h_{2,3} & \cdots & h_{2,p} \\ 0 & h_{3,2} & h_{3,3} & \cdots & h_{3,p} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & h_{p,p-1} & h_{p,p} \end{pmatrix}.$$

This matrix contains only zeros under the first subdiagonal; such a matrix is called a *Hessenberg* matrix. Inspecting the algorithm, we notice that the j -th column contains the coefficients of the projection of the vector $\mathbf{A}\mathbf{u}_j$ onto the basis $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$. In other words,

$$\mathbf{U}^* \mathbf{A} \mathbf{U} = \mathbf{H}, \quad (6.15)$$

We have thus shown that the Arnoldi algorithm enables to construct both an orthonormal basis of a Krylov subspace and the associated reduced matrix. In fact, we have the following equation

$$\mathbf{A} \mathbf{U} = \mathbf{U} \mathbf{H} + h_{p+1,p} (\mathbf{v}_{p+1} \mathbf{e}_p^*), \quad \mathbf{e}_p = \begin{pmatrix} 0 \\ \vdots \\ 1 \end{pmatrix} \in \mathbf{C}^p. \quad (6.16)$$

The Arnoldi algorithm, coupled with the Rayleigh–Ritz method, has very good convergence properties in the limit as $p \rightarrow \infty$, in particular for eigenvalues with a large modulus. The following result shows that the residual $\mathbf{r} = \mathbf{A}\hat{\mathbf{v}} - \hat{\lambda}\mathbf{v}$ associated with a Ritz vector can be estimated inexpensively. Specifically, the norm of the residual is equal to the last component of the associated eigenvector of \mathbf{H} multiplied by $h_{p+1,p}$.

Proposition 6.9 (Formula for the residual ). *Let \mathbf{y}_i be an eigenvector of \mathbf{H} associated with*

the eigenvalues $\widehat{\lambda}_i$, and let $\widehat{\mathbf{v}}_i = \mathbf{U}\mathbf{y}_i$ denote the corresponding eigenvector. Then

$$\mathbf{A}\widehat{\mathbf{v}}_i - \widehat{\lambda}_i\mathbf{v}_i = h_{p+1,p}(\mathbf{y}_i)_p\mathbf{v}_{p+1}.$$

Consequently, it holds that

$$\|\mathbf{A}\widehat{\mathbf{v}}_i - \widehat{\lambda}_i\mathbf{v}_i\| = |h_{p+1,p}(\mathbf{y}_i)_p|.$$

Proof. Multiplying both sides of (6.16) by \mathbf{y}_i , we obtain

$$\mathbf{A}\mathbf{U}\mathbf{y}_i = \mathbf{U}\mathbf{H}\mathbf{y}_i + h_{p+1,p}(\mathbf{v}_{p+1}\mathbf{e}_p^*)\mathbf{y}_i.$$

Using the definition of $\widehat{\mathbf{v}}_i$ and rearranging the equation, we have

$$\mathbf{A}\widehat{\mathbf{v}}_i - \widehat{\lambda}_i\mathbf{y}_i = h_{p+1,p}(\mathbf{v}_{p+1}\mathbf{e}_p^*)\mathbf{y}_i,$$

which immediately gives the result. \square

In practice, the larger the dimension p of the subspace \mathcal{U} employed in the Rayleigh–Ritz method, the more memory is required for storing an orthonormal basis of \mathcal{U} . In addition, for large values of p , computing the reduced matrix (6.15) and its eigenpairs becomes computationally expensive; the computational cost of computing the matrix \mathbf{H} scales as $\mathcal{O}(p^2)$. To remedy these potential issues, the algorithm can be restarted periodically. For example, Algorithm 15 can be employed as an alternative to the power iteration in order to find the eigenvector associated with the eigenvalue with largest modulus.

Algorithm 15 Restarted Arnoldi iteration

Choose $\mathbf{u}_1 \in \mathbf{C}^n$ and $p \ll n$

for $i \in \{1, 2, \dots\}$ **do**

 Perform p iterations of the Arnoldi iteration and construct \mathcal{U} ;

 Calculate the Ritz vector $\widehat{\mathbf{v}}_1$ associated with the largest Ritz value relative to \mathcal{U} ;

 If this vector is sufficiently accurate, then stop. Otherwise, restart with $\mathbf{u}_1 = \widehat{\mathbf{v}}_1$.

end for

6.4.3 The Lanczos iteration

The Lanczos iteration may be viewed as a simplified version of the Arnoldi iteration in the case where the matrix \mathbf{A} is Hermitian. Let us denote by $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ the orthonormal vectors generated by the Arnoldi iteration. When \mathbf{A} is Hermitian, it holds that

$$h_{i,j} = \mathbf{u}_i^*(\mathbf{A}\mathbf{u}_j) = (\mathbf{A}\mathbf{u}_i)^*\mathbf{u}_j = \overline{h_{j,i}}.$$

Therefore, the matrix \mathbf{H} is Hermitian. This is not surprising, since we showed that $\mathbf{H} = \mathbf{U}^*\mathbf{A}\mathbf{U}$ and the matrix \mathbf{A} is Hermitian. Since \mathbf{H} is also of Hessenberg form, we deduce that \mathbf{H} is tridiagonal. An inspection of Algorithm 14 shows that the subdiagonal entries of \mathbf{H} are real. Since \mathbf{A} is Hermitian, the diagonal entries $h_{i,i} = \mathbf{u}_i^*(\mathbf{A}\mathbf{u}_i)$ are also real, and so we conclude that

all the entries of the matrix \mathbf{H} are in fact real. This matrix is of the form

$$\mathbf{H} = \begin{pmatrix} \alpha_1 & \beta_2 & & & & \\ \beta_2 & \alpha_2 & \beta_3 & & & \\ & \beta_3 & \ddots & \ddots & & \\ & & \ddots & \ddots & \beta_p & \\ & & & & \beta_p & \alpha_p \end{pmatrix}$$

Adapting the Arnoldi iteration to this setting leads to [Algorithm 16](#).

Algorithm 16 Lanczos iteration for constructing an orthonormal basis of $\mathcal{K}_p(\mathbf{A}, \mathbf{u}_1)$

Choose \mathbf{u}_1 with unit norm.

$\beta_1 \leftarrow 0, \mathbf{u}_0 \leftarrow \mathbf{0} \in \mathbf{C}^n$

for $j \in \{1, \dots, p\}$ **do**

$\mathbf{u}_{j+1} \leftarrow \mathbf{A}\mathbf{u}_j - \beta_j\mathbf{u}_{j-1}$

$\alpha_j \leftarrow \mathbf{u}_j^* \mathbf{u}_{j+1}$

$\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1} - \alpha_j \mathbf{u}_j$

$\beta_{j+1} \leftarrow \|\mathbf{u}_{j+1}\|$

$\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1}/\beta_{j+1}$

end for

6.5 Exercises

⚙️ **Exercise 6.1.** *PageRank is an algorithm for assigning a rank to the vertices of a directed graph. It is used by many search engines, notably Google, for sorting search results. In this context, the directed graph encodes the links between pages of the World Wide Web: the vertices of the directed graph are webpages, and there is an edge going from page i to page j if page i contains a hyperlink to page j .*

Let us consider a directed graph $G(V, E)$ with vertices $V = \{1, \dots, n\}$ and edges E . The graph can be represented by its adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$, whose entries are given by

$$a_{ij} = \begin{cases} 1 & \text{if there is an edge from } i \text{ to } j, \\ 0 & \text{otherwise.} \end{cases}$$

Let r_i denote the “value” assigned to vertex i . The idea of PageRank, in its simplest form, is to assign values to the vertices by solving the following system of equations;

$$\forall i \in V, \quad r_i = \sum_{j \in \mathcal{N}(i)} \frac{r_j}{o_j}. \quad (6.17)$$

where o_j is the outdegree of vertex j , i.e. the number of edges leaving from j . Here the sum is over the set of nodes $\mathcal{N}(i)$, which denotes all the “incoming” neighbors of i , i.e. those that have an edge pointing towards node i .

- Read the Wikipedia page on PageRank to familiarize yourself with the algorithm.

- Let $\mathbf{r} = (r_1 \ \dots \ r_n)^T$. Show using (6.17) that \mathbf{r} satisfies

$$\mathbf{r} = \mathbf{A}^T \begin{pmatrix} \frac{1}{o_1} & & \\ & \ddots & \\ & & \frac{1}{o_n} \end{pmatrix} \mathbf{r} =: \mathbf{A}^T \mathbf{O}^{-1} \mathbf{r}.$$

In other words, \mathbf{r} is an eigenvector with eigenvalue 1 of the matrix $\mathbf{M} = \mathbf{A}^T \mathbf{O}^{-1}$.

- Show that \mathbf{M} is a left-stochastic matrix, i.e. that each column sums to 1.
- Prove that the eigenvalues of any matrix $\mathbf{B} \in \mathbf{R}^{n \times n}$ coincide with those of \mathbf{B}^T . You may use the fact that $\det(\mathbf{B}) = \det(\mathbf{B}^T)$.
- Using the previous items, show that 1 is an eigenvalue and that $\rho(\mathbf{M}) = 1$. For the second part, find a subordinate matrix norm such that $\|\mathbf{M}\| = 1$.
- Implement PageRank in order to rank pages from a 2013 snapshot of English Wikipedia. You can use either the simplified version of the algorithm given in (6.17) or the improved version with a damping factor described on Wikipedia. In the former case, the following are both sensible stopping criteria:

$$\frac{\|\mathbf{M}\hat{\mathbf{r}} - \hat{\mathbf{r}}\|_1}{\|\hat{\mathbf{r}}\|_1} < 10^{-15} \quad \text{or} \quad \frac{\|\mathbf{M}\hat{\mathbf{r}} - \hat{\lambda}\hat{\mathbf{r}}\|_1}{\|\hat{\mathbf{r}}\|_1} < 10^{-15}, \quad \hat{\lambda} = \frac{\hat{\mathbf{r}}^T \mathbf{M} \hat{\mathbf{r}}}{\hat{\mathbf{r}}^T \hat{\mathbf{r}}},$$

where $\hat{\mathbf{v}}$ is an approximation of the eigenvector corresponding to the dominant eigenvalue. A dataset is available on the course website to complete this part. This dataset contains a subset of the data publicly available [here](#), and was generated from the full dataset by retaining only the 5% best rated articles. After decompressing the archive, you can load the dataset into Julia by using the following commands:

```
import CSV
import DataFrames

# Data (nodes and edges)
nodes = CSV.read("names.csv", DataFrames.DataFrame)
edges = CSV.read("edges.csv", DataFrames.DataFrame)

# Convert data to matrices
nodes = Matrix(nodes)
edges = Matrix(edges)
```

After you have assigned a rank to all the pages, print the 10 pages with the highest ranks. My code returns the following entries:

- | | | |
|-------------------|---------------------|-----------|
| 1. United States | 5. France | 9. Canada |
| 2. United Kingdom | 6. Germany | 10. India |
| 3. World War II | 7. English language | |
| 4. Latin | 8. China | |

- **Extra credit:** Write a function `search(keyword)` that can be employed for searching the database. Here is an example of what it could return:

```
julia> search("New York")
481-element Vector{String}:
 "New York City"
 "New York"
 "The New York Times"
 "New York Stock Exchange"
 "New York University"
 ...
```

⚙️ **Exercise 6.2.** Show the following properties of the Krylov subspace $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$.

- $\mathcal{K}_p(\mathbf{A}, \mathbf{x}) \subset \mathcal{K}_{p+1}(\mathbf{A}, \mathbf{x})$.
- $A\mathcal{K}_p(\mathbf{A}, \mathbf{x}) \subset \mathcal{K}_{p+1}(\mathbf{A}, \mathbf{x})$.
- The Krylov subspace $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$ is invariant under rescaling: for all $\alpha \in \mathbf{C}$,

$$\mathcal{K}_p(\mathbf{A}, \mathbf{x}) = \mathcal{K}_p(\alpha\mathbf{A}, \mathbf{x}) = \mathcal{K}_p(\mathbf{A}, \alpha\mathbf{x}).$$

- The Krylov subspace $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$ is invariant under shift of the matrix \mathbf{A} : for all $\alpha \in \mathbf{C}$,

$$\mathcal{K}_p(\mathbf{A}, \mathbf{x}) = \mathcal{K}_p(\mathbf{A} - \alpha\mathbf{I}, \mathbf{x}).$$

- Similarity transformation: If $\mathbf{T} \in \mathbf{C}^{n \times n}$ is nonsingular, then

$$\mathcal{K}_p(\mathbf{T}^{-1}\mathbf{A}\mathbf{T}, \mathbf{T}^{-1}\mathbf{x}) = \mathbf{T}^{-1}\mathcal{K}_p(\mathbf{A}, \mathbf{x}).$$

⚙️ **Exercise 6.3.** The minimal polynomial of a matrix $\mathbf{A} \in \mathbf{C}^{n \times n}$ is the monic polynomial p of lowest degree such that $p(\mathbf{A}) = \mathbf{0}$. Prove that, if \mathbf{A} is Hermitian with $m \leq n$ distinct eigenvalues, then the minimal polynomial is given by

$$p(t) = \prod_{i=1}^m (t - \lambda_i).$$

⚙️ **Exercise 6.4.** The minimal polynomial for a general matrix $\mathbf{A} \in \mathbf{C}^{n \times n}$ is given by

$$p(t) = \prod_{i=1}^m (t - \lambda_i)^{s_i}.$$

where s_i is the size of the largest Jordan block associated with the eigenvalue λ_i in the normal Jordan form of A . Verify that $p(A) = 0$.

⚙ **Exercise 6.5.** Let d denote the degree of the minimal polynomial of A . Show that

$$\forall p \geq d, \quad \mathcal{K}_{p+1}(A, \mathbf{x}) = \mathcal{K}_p(A, \mathbf{x}).$$

Deduce that, for $p \geq n$, the subspace $\mathcal{K}_p(A, \mathbf{x})$ is an invariant subspace of A .

⚙ **Exercise 6.6.** Let $A \in \mathbf{C}^{n \times n}$. Show that $\mathcal{K}_n(A, \mathbf{x})$ is the smallest invariant subspace of A that contains \mathbf{x} .

□ **Exercise 6.7.** Consider the matrix

$$M = \begin{pmatrix} 0 & 1 & 2 & 0 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 0 & 2 \\ 0 & 0 & 2 & 0 \end{pmatrix}$$

- Find the dominant eigenvalue of M by using the power iteration.
- Find the eigenvalue of M closest to 1 by using the inverse iteration.
- Find the other two eigenvalues of M by using a method of your choice.

⚙ **Exercise 6.8** (A posteriori error bound). Assume that $A \in \mathbf{C}^{n \times n}$ is Hermitian, and that $\widehat{\mathbf{v}}$ is a normalized approximation of an eigenvector which satisfies

$$\|\widehat{\mathbf{z}}\| := \|A\widehat{\mathbf{v}} - \widehat{\lambda}\widehat{\mathbf{v}}\| = \delta, \quad \widehat{\lambda} = \frac{\widehat{\mathbf{v}}^* A \widehat{\mathbf{v}}}{\widehat{\mathbf{v}}^* \widehat{\mathbf{v}}}.$$

Prove that there is an eigenvalue λ of A such that

$$|\widehat{\lambda} - \lambda| \leq \delta.$$

Hint: Consider first the case where A is diagonal.

⚙ **Exercise 6.9** (Bauer–Fike theorem). Assume that $A \in \mathbf{C}^{n \times n}$ is diagonalizable: $AV = VD$. Show that, if $\widehat{\mathbf{v}}$ is a normalized approximation of an eigenvector which satisfies

$$\|\mathbf{r}\| := \|A\widehat{\mathbf{v}} - \widehat{\lambda}\widehat{\mathbf{v}}\| = \delta$$

for some $\widehat{\lambda} \in \mathbf{C}$, then there is an eigenvalue λ of A such that

$$|\widehat{\lambda} - \lambda| \leq \kappa_2(V)\delta.$$

Hint: Rewrite

$$\|\widehat{\mathbf{v}}\| = \|(A - \widehat{\lambda}I)^{-1}\mathbf{r}\| = \|V(D - \widehat{\lambda}I)^{-1}V^{-1}\mathbf{r}\|.$$

❁ **Exercise 6.10.** In *Exercise 6.8* and *Exercise 6.9*, we saw examples a posteriori error estimates which guarantee the existence of an eigenvalue of \mathbf{A} within a certain distance of the approximation $\hat{\lambda}$. In this exercise, our aim is to provide an answer to the following different question: given an approximate eigenpair $(\hat{\mathbf{v}}, \hat{\lambda})$, what is the smallest perturbation \mathbf{E} that we need to apply to \mathbf{A} in order to guarantee that $(\hat{\mathbf{v}}, \hat{\lambda})$ is an exact eigenpair, i.e. that

$$(\mathbf{A} + \mathbf{E})\hat{\mathbf{v}} = \hat{\lambda}\hat{\mathbf{v}}?$$

Assume that $\hat{\mathbf{v}}$ is normalized and let $\mathcal{E} = \{\mathbf{E} \in \mathbf{C}^{n \times n} : (\mathbf{A} + \mathbf{E})\hat{\mathbf{v}} = \hat{\lambda}\hat{\mathbf{v}}\}$. Prove that

$$\min_{\mathbf{E} \in \mathcal{E}} \|\mathbf{E}\|_2 = \|\mathbf{r}\|_2 := \|\mathbf{A}\hat{\mathbf{v}} - \hat{\lambda}\hat{\mathbf{v}}\|. \quad (6.18)$$

To this end, you may proceed as follows:

- Show first that any $\mathbf{E} \in \mathcal{E}$ satisfies $\mathbf{E}\hat{\mathbf{v}} = -\mathbf{r}$.
- Deduce from the first item that

$$\inf_{\mathbf{E} \in \mathcal{E}} \|\mathbf{E}\|_2 \geq \|\mathbf{r}\|_2.$$

- Find a rank one matrix \mathbf{E}_* such that $\|\mathbf{E}_*\|_2 = \|\mathbf{r}\|_2$, and then conclude. Recall that any rank 1 matrix can be written in the form $\mathbf{E}_* = \mathbf{u}\mathbf{w}^*$, with norm $\|\mathbf{u}\|_2\|\mathbf{w}\|_2$.

Equation (6.18) is a simplified version of the Kahan–Parlett–Jiang theorem and is an example of a backward error estimate. Whereas forward error estimates quantify the distance between an approximation and the exact solution, backward error estimates give the size of the perturbation that must be applied to a problem so that an approximation is exact.

❁ **Exercise 6.11.** Assume that $(\mathbf{x}_k)_{k \geq 0}$ is a sequence of normalized vectors in \mathbf{C}^n . Show that the following statements are equivalent:

- $(\mathbf{x}_k)_{k \geq 0}$ converges essentially to \mathbf{x}_∞ in the limit as $k \rightarrow \infty$.
- The angle $\angle(\mathbf{x}_k, \mathbf{x}_\infty)$ converges to zero in the limit as $k \rightarrow \infty$.
- The projector $\mathbf{P}_{\mathbf{x}_k}$ converges to $\mathbf{P}_{\mathbf{x}_\infty}$ in the limit as $k \rightarrow \infty$.

❁ **Exercise 6.12.** Assume that $\mathbf{A} \in \mathbf{C}^{n \times n}$ is skew-Hermitian. Derive a Lanczos-like algorithm for constructing an orthonormal basis of $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$ and calculating the reduced matrix

$$\mathbf{U}^* \mathbf{A} \mathbf{U},$$

where $\mathbf{U} \in \mathbf{C}^{n \times p}$ contains the vectors of the basis as columns. Implement your algorithm with $p = 20$ in order to approximate the dominant eigenvalue of the matrix \mathbf{A} constructed by the following piece of code:

```
n = 5000
A = rand(n, n) + im * randn(n, n)
A = A - A' # A is now skew-Hermitian
```

⚙ **Exercise 6.13.** Assume that $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ and $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$ are orthonormal bases of the same subspace $\mathcal{U} \subset \mathbf{C}^n$. Show that the projectors $\mathbf{U}\mathbf{U}^*$ and $\mathbf{W}\mathbf{W}^*$ are equal.

⚙ **Exercise 6.14.** Assume that $\mathbf{A} \in \mathbf{C}^{n \times n}$ is a Hermitian matrix with distinct eigenvalues, and suppose that we know the dominant eigenpair $(\mathbf{v}_1, \lambda_1)$, with \mathbf{v}_1 normalized. Let

$$\mathbf{B} = \mathbf{A} - \lambda_1 \mathbf{v}_1 \mathbf{v}_1^*.$$

If we apply the power iteration to this matrix, what convergence can we expect?

⚙ **Exercise 6.15.** Assume that $\widehat{\mathbf{v}}_1$ and $\widehat{\mathbf{v}}_2$ are two Ritz vectors of a Hermitian matrix \mathbf{A} relative to a subspace $\mathcal{U} \subset \mathbf{C}^n$. Show that, if the associated Ritz values are distinct, then $\widehat{\mathbf{v}}_1^* \widehat{\mathbf{v}}_2 = 0$.

6.6 Discussion and bibliography

The content of this chapter is inspired mainly from [14] and also from [11]. The latter volume contains a detailed coverage of the standard methods for eigenvalue problems. Some of the exercises are taken from [16], and others are inspired from [11].