

# 1 Iterative methods for large sparse eigenvalue problems

The eigenvalue problem is one of the fundamental problems in science. We wish to compute pairs  $\lambda \in \mathbb{C}$  and  $x \in \mathbb{C}^n$  such that

$$Ax = \lambda x, \quad (1.1)$$

and  $A \in \mathbb{R}^{n \times n}$ . The eigenvector is in this block assumed to be normalized as  $\|x\| = 1$ , with Euclidean norm such that  $\|x\|^2 = x^H x = 1$ .

## 1.1 Basic methods

### 1.1.1 Computing eigenvalues from eigenvectors

Before diving into the algorithms for eigenvalue problems, we will treat an easier problem.

**Problem:** Suppose  $x \in \mathbb{C}^n$  is an approximation of an eigenvector, compute an associated eigenvalue.

Assume for the moment an idealized situation where  $x$  is exactly an eigenvector. This means that (1.1) is satisfied, and we can multiply the equation from the left with  $x^H$ :

$$x^H Ax = \lambda x^H x,$$

such that

$$\lambda = \frac{x^H Ax}{x^H x}$$

This quotient can be used also if  $x$  is not an eigenvector and is usually referred to as the Rayleigh quotient.

**Definition 1.1.1** (Rayleigh quotient). *The quotient defined by*

$$r(x) := \frac{x^H Ax}{x^H x}$$

*is referred to as the Rayleigh quotient.*

For symmetric matrices, there are additional interpretations of the Rayleigh quotient. Given an approximate eigenvector  $x$ , it minimizes  $Ax - \mu x$  in the Euclidean norm: One can show that

$$\operatorname{argmin}_{\mu \in \mathbb{R}} \|Ax - \mu x\| = \frac{x^H Ax}{x^H x}$$

The Rayleigh quotient will in general only give you an approximation of the eigenvalue. The propagation of the approximation error can also be precisely described if  $x$  is sufficiently close to an eigenvector. More precisely, if we suppose  $x \in \mathbb{C}^n$  is an eigenvector corresponding to an eigenvalue  $\lambda$ , we have that

$$r(x + \varepsilon y) = \lambda + \mathcal{O}(\varepsilon). \quad (1.2)$$

In words, the error in the eigenvalue from the Rayleigh quotient is essentially of the order of magnitude of the error in the eigenvector. In the following this is made more concrete with an example and a theorem describing the accuracy. If  $A$  is symmetric (or hermitian) we have

$$r(x + \varepsilon y) = \lambda + \mathcal{O}(\varepsilon^2). \quad (1.3)$$

Note that  $\mathcal{O}(\varepsilon^2)$  is better than  $\mathcal{O}(\varepsilon)$  since we consider small  $\varepsilon$ .

### Rayleigh quotient

Consider the two matrices

$$A_1 = \begin{bmatrix} 2 & 5 \\ 0 & 3 \end{bmatrix} \text{ and } A_2 = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}.$$

Both matrices have an eigenvector  $x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  with eigenvalue  $\lambda = 2$ , but  $A_2$  is symmetric. The example code below illustrates that the Rayleigh quotient is much closer to the eigenvalue for the symmetric matrix  $A_2$ .

```
>> A1=[2 0;0 3];
>> A2=[2 5;0 3];
>> x=[1;0];
>> y=[1;1]; e=1e-4 % small perturbation
>> z=x+e*y;
>> z'*A1*z/(z'*z)
ans =
    2.000499959998002
>> z'*A2*z/(z'*z)
ans =
    2.000000009998000
```



**Theorem 1.1.2** (Accuracy of the Rayleigh quotient). Suppose  $(\lambda, x)$  is an eigenpair of  $A$  with  $\|x\| = 1$ . Let  $v = x + \varepsilon \Delta$  where  $\varepsilon \in \mathbb{R}$  and  $\|\Delta\| = 1$ . Then, for sufficiently small  $\varepsilon$

$$r(v) - \lambda = \begin{cases} \mathcal{O}(\varepsilon^2) & \text{if } x^T A = \lambda x^T \\ \mathcal{O}(\varepsilon) & \text{otherwise.} \end{cases}$$

Note that if  $(\lambda, x)$  is an eigenpair of  $A$  and  $A$  is symmetric we have that  $A^T x - \lambda x = 0$  whose transpose is  $x^T A - \lambda x^T = 0$ . Therefore, by Theorem 1.1.2 the accuracy is quadratic in  $\varepsilon$  for symmetric matrices.

*Proof.* First expand the Rayleigh quotient with the approximation

$$r(v) = \frac{(x + \varepsilon\Delta)^T A(x + \varepsilon\Delta)}{(x + \varepsilon\Delta)^T (x + \varepsilon\Delta)} = \frac{1}{x^T x + \varepsilon\alpha} (x^T A x + \varepsilon(\beta + \varepsilon\gamma))$$

where  $\alpha = x^T \Delta + \Delta^T x + \varepsilon \Delta^T \Delta$ ,  $\beta = x^T A \Delta + \Delta^T A x$  and  $\gamma = \Delta^T A \Delta$ . We will now use the Taylor expansion of functions of the form

$$\frac{1}{1+z} = 1 - z + z^2 - \dots$$

By selection  $z = \varepsilon\alpha$  and noting that  $x^T x = 1$  by assumption and using that  $Ax - \lambda x = 0$ , we conclude that

$$r(v) = \lambda + \varepsilon(x^T A - \lambda x^T) + \mathcal{O}(\varepsilon^2)$$

which reduces to the statement of the theorem.  $\square$

### 1.1.2 Basic eigenvalue methods

The Rayleigh quotient provides a procedure to numerically compute an eigenvalue approximation given an eigenvector approximation. Computing the eigenvector can be done in many ways. We first consider three basic algorithms.

- Power method (power iteration) summarized in Algorithm 1
- Inverse iteration summarized in Algorithm 2
- Rayleigh quotient iteration summarized in Algorithm 3

Read about these methods in TB pages 202-209.

### 1.1.3 Power method

The power method (or sometimes power iteration), is our first eigenvalue method. It consists of starting vector a vector  $v_0$ , we multiply this vector with  $A$ , scale the resulting vector and repeat the process:

$$v_{k+1} = \alpha_k A v_k, \quad k = 0, \dots$$

The scaling factor  $\alpha_k$  is used to prevent the iteration values  $v_k$  to become very small or very large which makes them more difficult to represent/store. (More precisely, we want to avoid overflow or underflow in the IEEE floating point arithmetic.) Typically the scaling is selected such that  $\|v_{k+1}\| = 1$ , which can be achieved by setting

$$\alpha_k = \frac{1}{\|A v_k\|}.$$

The operations can be re-ordered such it only requires one matrix vector product per iteration as in Algorithm 1.

An important property of the power method is that the only way we need to access the matrix  $A$  is in combination with a multiplication with a vector  $Ax$ : a so-called *matrix-vector product*. In many scientific applications, the matrix  $A$  may be so large that it is not possible to store it explicitly, but the matrix-vector product may still be available.

If we consider  $v_k$  as our eigenvector approximation, we can use the Rayleigh quotient to extract an eigenvalue approximation. Since  $\|v_k\|_2^2 = v_k^T v_k = 1$ , the Rayleigh quotient reduces to

$$\tilde{\lambda}_k = \frac{v_k^T A v_k}{v_k^T v_k} = v_k^T A v_k.$$

**Input:** A starting vector  $v$  with  $\|v\| = 1$   
**Output:** Eigenpair approximation  $(w, \tilde{\lambda})$   
**for**  $n = 1, 2, \dots$  **do**  
     $w = Av$   
     $v = w / \|w\|$   
     $\tilde{\lambda} = v^T A v$   
**end**

**Algorithm 1:** Power method (Power iteration).

### 1.1.4 Convergence of the power method

It turns out that the iterates  $v_k$  generated by the power method do indeed in general converge to an eigenvector. Under certain (not very restrictive) conditions one can show that

$$\tilde{\lambda}_k = \lambda_1 + \mathcal{O}\left(\frac{|\lambda_2|^k}{|\lambda_1|^k}\right).$$

where we have ordered the eigenvalues as  $|\lambda_1| \geq |\lambda_2| \geq \dots$ .

**Theorem 1.1.3.** Consider a matrix  $A \in \mathbb{C}^{n \times n}$ , and assume its largest eigenvalue is distinct in modulus such that

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \dots \geq |\lambda_n|.$$

If  $A = XDX^{-1}$  is a Jordan decomposition of  $A$  with  $D(1,1) = \lambda$ . Suppose the power method is initiated such that the first element of  $X^{-1}v_0$  is non-zero. Then,

$$|\tilde{\lambda}_k - \lambda_1| = \mathcal{O}\left(\frac{|\lambda_2|^k}{|\lambda_1|^k}\right).$$

See proof during lecture.

### 1.1.5 Inverse iteration

The next algorithm we consider is essentially a combination of what we know for the power method, and the observation that the eigenvalues of the matrix  $A$  and the matrix

$$B = (A - \mu I)^{-1}$$

are related by a simple relation.

Unlike the power method, inverse iteration does not involve a matrix vector product with  $A$  per iteration, but one solution to the linear system,  $(A - \mu I)^{-1}v$ . This operation is normally called a *linear solve*. A linear solve is in general much more computationally expensive than one matrix vector product.

If we denote  $\lambda_i(A)$ ,  $\lambda_i(B)$  the eigenvalues of  $A$  and  $B$  respectively, the eigenvalues are related by

$$\lambda_i(B) = \frac{1}{\lambda_i(A) - \mu} \quad (1.4)$$

The eigenvectors remain unchanged.

The transformation (1.4) has the useful property that the eigenvalues close to  $\mu$  will be transformed to large eigenvalues. Since inverse iteration converges to the eigenvector corresponding to the largest eigenvalue in general, we obtain with the application of the power method to  $B$ , which converges to the eigenvector corresponding to an eigenvalue of  $A$ , closest to  $\mu$ . The eigenvector extraction can be done with the Rayleigh quotient of  $A$ , rather than  $B$ , as shown in Algorithm 2.

The convergence follows from the fact that the method is equivalent to the power method applied to the matrix  $B$ :

$$\tilde{\lambda}_k = \lambda_1 + \mathcal{O}\left(\frac{|\lambda_J - \mu|^k}{|\lambda_K - \mu|^k}\right) \quad (1.5)$$

where  $\lambda_J$  is the eigenvalue of  $A$  that is closest to  $\mu$  and  $\lambda_K$  is the eigenvalue of  $A$  second closest to  $\mu$ .

**Input:** A starting vector  $v$  with  $\|v\| = 1$  and shift  $\mu$   
**Output:** Eigenpair approximation  $(w, \tilde{\lambda})$   
**for**  $n = 1, 2, \dots$  **do**  
    Solve linear system  $(A - \mu I)w = v$   
     $v = w / \|w\|$   
     $\tilde{\lambda} = v^T A v$   
**end**

**Algorithm 2:** Inverse iteration

An interpretation of (1.5): The convergence factor of inverse iteration is proportional to the distance between the shift and the closest eigenvalue. In formulas convergence to the eigenvector  $v$  is

$$\|v_{k+1} - v\| = \mathcal{O}(|\lambda_J - \mu| \|v_k - v\|)$$

### 1.1.6 Rayleigh quotient iteration

For the final basic algorithm, we use a combination of previous ideas. We can set  $\mu$  in inverse iteration as the Rayleigh quotient. This implies that when the eigenvector is a good approximation, the corresponding eigenvalue of  $B = (A - \mu I)^{-1}$  will be a very large value and therefore converge faster than constant  $\mu$ .

The theoretical convergence of Rayleigh quotient iteration can also be determined by combining results above. In this setting, we will simplify the inverse iteration convergence theory by noting that one step essentially multiplies the previous error with the convergence factor which is  $\lambda - \mu$ :

$$v_{k+1} - v = \mathcal{O}(\|v_k - v\| |\lambda - \mu|). \quad (1.6)$$

The relationship between (1.5) and (1.6) is consequence of linear convergence. A method which converges linearly with convergence factor  $\alpha$  can be described in two equivalent ways

- $v_k - v = \mathcal{O}(\alpha^k)$
- $v_{k+1} - v = \mathcal{O}(\alpha \|v_k - v\|)$

Formally, this can be derived from (1.5). Suppose now that we initiate Rayleigh quotient iteration with error  $\varepsilon$ :

$$\|v_k - v\| = \varepsilon$$

We compute the eigenvalue approximation with the Rayleigh quotient. The error of the Rayleigh quotient is given by Theorem 1.1.2. Therefore here we have

$$\lambda_k - \lambda = O(\|v_k - v\|^p) = O(\varepsilon^p).$$

Subsequently, the next eigenvector approximation is computed with inverse iteration whose error is propagated as (1.6):

$$v_{k+1} - v = O(\|v_k - v\| |\lambda - \lambda_k|) = O(\varepsilon^{p+1}) = \begin{cases} O(\varepsilon^3) & \text{if } x^T A = \lambda x^T \\ O(\varepsilon^2) & \text{otherwise.} \end{cases}$$

In the symmetric case, the error has reduced from  $\varepsilon$  to  $\varepsilon^3$ , and the method has cubic convergence for symmetric matrices.

However, in contrast to inverse iteration and the power method, the convergence theory does not determine to which eigenvalue the method converges; it highly depends on the starting vector.

Note that the Rayleigh quotient iteration can also be used for non-symmetric matrices, although it is sometimes presented only as a method for symmetric matrices.

**Input:** A starting vector  $v$  with  $\|v\| = 1$  and starting eigenvalue  $\mu$   
**Output:** Eigenpair approximation  $(w, \tilde{\lambda})$   
 Set  $\tilde{\lambda} = \mu$   
**for**  $n = 1, 2, \dots$  **do**  
     Solve linear system  $(A - \tilde{\lambda}I)w = v$   
      $v = w / \|w\|$   
      $\tilde{\lambda} = v^T A v$   
**end**

**Algorithm 3:** Rayleigh Quotient Iteration

## 1.2 Orthogonal matrices and orthogonalizing vectors

In basic linear algebra, we learn that two vectors  $x, y \in \mathbb{R}^n$  are orthogonal when  $y^T x = 0$ . The concept of orthogonality, and its generalization to matrices is very important in this course. We will use it mostly in different factorizations and decompositions of matrices.

The use of decompositions has been selected as one of the most influential concepts in algorithms in the 20th century: <https://www.siam.org/pdf/news/637.pdf> In this course we also cover other algorithms in the list of important algorithms.

### 1.2.1 Gram-Schmidt procedures

The Gram-Schmidt procedure is often explained as a procedure to orthogonalize vectors, meaning that given vectors stored in a matrix  $F = [f_1, \dots, f_m] \in \mathbb{R}^{n \times m}$  with  $n \geq m$  we try to determine  $q_1, \dots, q_m$  such that  $q_1, \dots, q_m$  are orthonormal and

$$\text{span}(f_1, \dots, f_m) = \text{span}(q_1, \dots, q_m).$$

Such vectors  $q_1, \dots, q_m$  exist if  $f_1, \dots, f_m$  are linearly independent vectors. Note that the matrix  $Q = [q_1, \dots, q_m] \in \mathbb{R}^{n \times m}$  is orthogonal in the sense of definition of orthogonal matrices (see background.pdf).

The Gram-Schmidt procedure can be directly derived by inductively applying the following result.

**Lemma 1.2.1.** Suppose  $Q = [q_1, \dots, q_m] \in \mathbb{R}^{n \times m}$  is an orthogonal matrix and suppose  $b \notin \text{span}(q_1, \dots, q_m)$ . Let

$$h = Q^T b$$

and

$$z = b - Qh = (I - QQ^T)b. \quad (1.7)$$

Let  $\beta = \|z\|$  and define

$$q_{m+1} := \frac{z}{\beta} \quad (1.8)$$

Then,

(a) the matrix  $[q_1, \dots, q_{m+1}]$  is an orthogonal matrix;

(b)  $b = h_1 q_1 + \dots + h_m q_m + \beta q_{m+1}$ ; and

(c)  $\text{span}(q_1, \dots, q_{m+1}) = \text{span}(q_1, \dots, q_m, b)$ .

*Proof.* Proof of (b): This is a direct consequence of (1.7) and (1.8). Proof of (a): Note that

$$[q_1, \dots, q_{m+1}]^T [q_1, \dots, q_{m+1}] = [Q, q_{m+1}]^T [Q, q_{m+1}] = \begin{bmatrix} Q^T Q & Q^T q_{m+1} \\ q_{m+1}^T Q & q_{m+1}^T q_{m+1} \end{bmatrix}$$

The conclusion (a) follows from the fact that  $Q^T Q = I$ ,

$$Q^T q_{m+1} = Q^T (I - QQ^T)b = 0$$

and  $q_{m+1}^T q_{m+1} = 1$ .

Proof of (c): In this course we will several times use the general property that if two rectangular matrices  $W \in \mathbb{R}^{n \times m}$  and  $V \in \mathbb{R}^{n \times m}$  are related by

$$W = VP \quad (1.9)$$

You have normally learned about the Gram-Schmidt procedure in basic linear algebra courses. We repeat it in a slightly different notation than normal (using orthogonal matrices). It turns out that the classical Gram-Schmidt is not always satisfactory.

In numerical linear algebra, the Gram-Schmidt procedure directly derived from Lemma 1.2.1 is typically called the *classical* Gram-Schmidt procedure in order to distinguish it from variants we discuss later.

The vector  $h \in \mathbb{R}^n$  is typically referred to as the Gram-Schmidt coefficients.

for some non-singular matrix  $P \in \mathbb{R}^{m \times m}$ , then  $\text{span}(W) = \text{span}(V)$ .  
 If we select  $P$  as

$$P = \begin{bmatrix} I & h \\ 0 & \|z\| \end{bmatrix}$$

then (1.9) is satisfied with  $V = [Q, q_{m+1}]$  and  $W = [Q, b]$ .  $\square$

### Classical Gram-Schmidt example

```
>> Q=(1/sqrt(2))*[1 -1; 1 1; 0 0; 0 0];
>> Q'*Q % Check if Q is orthogonal
ans =
    1.0000    0
         0    1.0000
>> b=randn(4,1);
>> h=Q'*b; % Compute Gram-Schmidt coefficients
>> z=b-Q*h; % Compute "orthogonal complement"
>> beta=norm(z);
>> q_new=z/beta;
>> Q_new=[Q,q_new]; % Construct new Q-matrix
>> Q_new'*Q_new % Check that Q_new is orthogonal
ans =
    1.0000    0    0
         0    1.0000    0
         0    0    1.0000
>> norm(Q_new*[eye(2), h; zeros(1,2), norm(z)]-[Q,b])
>> P=[eye(2), h; zeros(1,2), beta];
>> norm(Q_new*P-[Q,b]) % Check that span(Q_new)=span([Q,b])
ans =
    1.1444e-16
```

Although the above example suggests that classical Gram-Schmidt works, it will in general not be satisfactory in our context. It turns out that the classical Gram-Schmidt is very sensitive to round-off errors in certain situations.

A detailed analysis of the influence of the round-off errors can be found in (appendix) Section 1.7 from which extract one conclusion. If the vector to be orthogonalized is almost in the subspace, we are likely to obtain a large error. Suppose  $b = q + \delta e$  where  $q \in \text{span}(Q)$  (meaning there  $q = Qd$ ) and  $e \perp Q$  and  $\|e\| = 1$ , for a small  $\delta$ . Then, the round-off error is

$$\frac{|\varepsilon|}{|\delta|} \|Qd\| + \mathcal{O}(\varepsilon^2).$$

In practice, we have round-off errors in every floating point operation and a complete round-off error analysis is quite cumbersome. In our simplified analysis we assume that no error is introduced in the computation of  $z$  and  $\tilde{q}_{m+1}$ . In particular, no additional round-off error is introduced in (1.19) and (1.20).



which suggests that the round-off error is proportional to  $|\varepsilon|/|\delta|$ , and can be very large if  $\delta$  is very small.

---

**Conclusion of error analysis of classical Gram-Schmidt method.** The Gram-Schmidt procedure is likely to have a large round-off error if the vector  $b$  almost lies in the subspace  $\text{span}(Q)$ .

---

### Modified Gram-Schmidt

In this course we consider two variations of Gram-Schmidt which aim to improve the floating-point arithmetic problems described above.

We now derive the algorithm called *the modified Gram-Schmidt procedure* from the classical Gram-Schmidt procedure. For theoretical purposes we express the classical Gram-Schmidt in for-loops:

```
for i=1:m
    h(i)=Q(:,i)'*b;
end
z=b;
for i=1:m
    z=z-h(i)*Q(:,i)
end
```

Note that at iteration  $i$  of the second loop, we only need  $h(i)$  computed at the  $i$ th iteration the first loop such that we can merge the two loops:

```
z=b;
for i=1:m
    h(i)=Q(:,i)'*b;
    z=z-h(i)*Q(:,i);
end
beta=norm(z);
```

In the first step inside the for-loop, the vector  $z$  can be explicitly expressed as:

- Iteration  $i = 1$ :  $z = b$
- Iteration  $i = 2$ :  $z = b - h_1 q_1$
- $\vdots$
- Iteration  $i = m$ :  $z = b - h_1 q_1 - \dots - h_m q_{m-1}$

Now recall that the vectors  $q_1, \dots, q_m$  are assumed to be orthogonal. The following identities can be directly identified.

- Iteration  $i = 1$ :  $q_i^T z = q_i^T b$

The modified Gram-Schmidt procedure is equivalent to the classical Gram-Schmidt procedure in exact arithmetic, but different floating-point arithmetic.

Although modified Gram-Schmidt yields a different result in floating point arithmetic, it is not always clear that the result is better. In fact, theoretical understanding for this is still disputed by some scientists. You will investigate this in practice by for a specific situation in the homeworks.

Caution regarding terminology: In this course we consider  $Q \in \mathbb{R}^{n \times m}$  as an orthogonal matrix and want to orthogonalize  $b$  which result in algorithms above. In some literature (such as TB) Gram-Schmidt procedures are described for orthogonalizing an entire matrix  $A \in \mathbb{R}^{n \times (m+1)}$ .

- Iteration  $i = 2$ :  $q_i^T z = q_2^T (b - h_1 q_1) = q_2^T b - h_1 q_2^T q_1 = q_i^T b$
- $\vdots$
- Iteration  $i = m$ :  $q_i^T z = q_m^T b - h_1 q_m^T q_1 - \dots - h_m q_m^T q_{m-1} = q_m^T b = q_i^T b$

Note that for every iteration we have  $q_i^T z = q_i^T b$ . Therefore, we can replace  $Q(:, i)' * b$  with  $Q(:, i)' * z$  in the for-loop. This is what we call the modified Gram-Schmidt method.

```
z=b;
for i=1:m
    h(i)=Q(:,i)'\*z;
    z=z-h(i)*Q(:,i)
end
```

### Double Gram-Schmidt

The next approach to improve the classical Gram-Schmidt procedure is very naive. Since we know that round-off errors will make the vector  $z = b - Qh$  to not be orthogonal in practice, we can try to make it orthogonal by applying classical Gram-Schmidt again. This is what is called repeated Gram-Schmidt, or the special case double Gram-Schmidt.

```
>> h=Q'\*b;
>> z=b-Q*h;
>> g=Q'\*z;
>> z=z-Q*g;
>> h=h+g;
>> beta=norm(z);
>> z=z/norm(z);
```

## 1.3 Krylov methods

The power method was the basis of both inverse iteration and Rayleigh quotient iteration. These algorithms can be used to compute one eigenvector given an initial guess. In order to compute several eigenvalues we now extend the power method in a different way. We consider the space spanned by the iterates of the power method.

**Definition 1.3.1** (Krylov subspace). *The span of the iterates of the power method is called a Krylov subspace*

$$\mathcal{K}_m(A, b) := \text{span}(b, Ab, A^2b, \dots, A^{m-1}b).$$

Due to rounding error issues, the Krylov subspace is usually not computed from  $[b, Ab, A^2b, \dots, A^{m-1}b]$ , but rather represented with an orthogonal basis of  $\mathcal{K}_m(A, b)$ . The Arnoldi method can be seen as method to compute an orthogonal basis of a Krylov subspace. More precisely, the Arnoldi method is a method which generates an orthogonal matrix  $Q_m \in \mathbb{C}^{n \times m}$  such that

$$AQ_m = Q_{m+1}\underline{H}_m$$

where  $\underline{H}_m \in \mathbb{R}^{(m+1) \times m}$  and  $Q_{m+1} = [Q_m, q_{m+1}]$ . The matrix  $\underline{H}_m$  is a so-called Hessenberg matrix, which means that the elements below the first lower off-diagonal are zero:

$$\underline{H}_m = \begin{bmatrix} \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times \end{bmatrix}$$

The Arnoldi method can be used to compute many quantities. In the context of eigenvalue computations, we take the eigenvalues of  $H_m \in \mathbb{C}^{m \times m}$  as eigenvalue approximations.

Arnoldi's method for eigenvalue problems is also discussed in TB pages 251–264.

### 1.3.1 Derivation of the Arnoldi method

The Arnoldi method will be seen as method to compute the Arnoldi factorization, by expanding  $Q_m$  and  $\underline{H}_m$  to form  $Q_{m+1}$  and  $\underline{H}_{m+1}$ . The algorithm can be derived by induction. Suppose an Arnoldi factorization for  $m = 2$  is given

$$AQ_2 = Q_3\underline{H}_2 \quad (1.10)$$

and we wish to expand the matrices such that they satisfy

$$AQ_3 = Q_4\underline{H}_3. \quad (1.11)$$

This is a matrix equality and if we consider column 1, 2 of this equality we obtain exactly (1.10). Column 3 is given by multiplication with  $e_3$ :

$$AQ_3e_3 = Q_4\underline{H}_3e_3.$$

We simplify this equation to

$$Aq_3 = q_1h_{1,3} + q_2h_{2,3} + q_3h_{3,3} + q_4h_{4,3}. \quad (1.12)$$

Note that  $q_1, q_2, q_3$  are known since they form  $Q_3$ . It remains to determine  $h_{1,3}, \dots, h_{4,3}$  and  $q_4$ . If we denote the left-hand side of (1.12)

by  $b$  we see that the problem to determine the coefficients is exactly the problem we solved with Gram-Schmidt in Lemma 1.2.1. Therefore, the Arnoldi method essentially consists of applying matrix vector products and carrying out a Gram-Schmidt procedure.

### 1.3.2 The Arnoldi method

If we combine the (Gram-Schmidt) orthogonalization process with the Krylov subspace, we obtain the algorithm called the Arnoldi method.

**Input:** A starting vector  $b$   
**Output:** Eigenpair approximation  
 Set  $q_1 = b/\|b\|$ ,  $H_0$  =empty matrix  
**for**  $m = 1, 2, \dots$  **do**  
     Compute  $x = Aq_m$   
     Orthogonalize  $x$  against  $q_1, \dots, q_m$  by computing  $h \in \mathbb{C}^m$  and  $x_\perp \in \mathbb{C}^n$  such that  $Q^T x_\perp = 0$  and

$$x_\perp = x - Qh.$$

    Let  $\beta = \|x_\perp\|$   
     Let  $q_{m+1} = x_\perp/\beta$   
     Expand  $\underline{H}_{m-1}$  with one column:

$$\underline{H}_m := \begin{bmatrix} \underline{H}_{m-1} & h \\ 0 & \beta \end{bmatrix}$$

**end**

**Algorithm 4:** Arnoldi's method for eigenvalue problems.

### 1.3.3 The Lanczos method

There are various ways to improve and specialize the Arnoldi method for specific matrix structures. The most prominent specialization is for symmetric matrices and called the Lanczos method. First observe that  $H_m$  can be expressed as

$$H_m = Q_m^T A Q_m.$$

By transposing the left-hand side and the right-hand side we obtain

$$H_m^T = (Q_m^T A Q_m)^T = Q_m^T A^T Q_m = Q_m^T A = Q_m = H_m.$$

Hence,  $H_m$  is also symmetric. Since  $H_m$  is both symmetric and has the Hessenberg structure, it is a tridiagonal matrix. This forms the

**History:** The Arnoldi method was invented by Walter Edwin Arnoldi in 1951 and the Lanczos method is named after the work of Cornelius Lanczos in 1950. In those days, most eigenvalue problems arose in acoustics and vibrations that lead to symmetric matrices, which is one reason the symmetric specialization was invented first. The Krylov method is named after Aleksey Krylov (or Алексѣй Крылов) who presented some ideas for Krylov subspaces in the context of naval engineering in 1931. This method class was also deemed the one of the most important algorithms in the 20th century by SIAM - Society of industrial and applied mathematics.

**Current research:** Large parts of the current international numerical linear algebra research community focus on Krylov methods, with challenges ranging from modern hardware implementations to generalizations for structured eigenproblems in new emerging fields. Search the web for "book of abstracts" and "numerical linear algebra" and "Krylov".

The Lanczos iteration is also described in TB pages 276-278.

foundation of the derivation of the Lanczos method.

Output: Eigenpair approximations  
 Input: The matrix  $A$  and vector  $b$ .  
 $b$  =arbitrary,  $q_1 = b/\|b\|$ ,  $H_0$  =empty matrix  
**for**  $m = 1, 2, \dots$  **do**  
      $v = Aq_m$   
      $\alpha_m = q_m^T v$   
      $v = v - \beta_{m-1}q_{m-1} - \alpha_m q_m$   
      $\beta_m = \|v\|$   
      $q_{m+1} = v/\beta_m$   
**end**  
 Construct the matrix

$$H = \begin{bmatrix} \alpha_1 & \beta_1 & & \\ \beta_1 & \ddots & \ddots & \\ & \ddots & \ddots & \beta_{m-1} \\ & & \beta_{m-1} & \alpha_m \end{bmatrix}$$

**Algorithm 5:** The Lanczos method

#### 1.4 Convergence of Arnoldi's method for eigenvalue problems

Recall that, unless it breaks down,  $k$  steps of the Arnoldi method generates an orthogonal basis of a Krylov subspace, represented by a matrix  $Q = [q_1, \dots, q_k] \in \mathbb{C}^{n \times k}$  such that  $Q^*Q = I$  and

$$\text{span}(q_1, \dots, q_k) = \mathcal{K}_k(A, b) := \text{span}(b, Ab, \dots, A^{k-1}b).$$

The eigenvalue approximations (called Ritz values) are subsequently found from the eigenvalues of

$$H = Q^*AQ.$$

The matrix  $H \in \mathbb{C}^{k \times k}$  is a Hessenberg matrix and can be generated as a by-product of the Arnoldi method. We call a pair  $(\mu, Qv)$  a Ritz pair and  $Qv$  a Ritz vector, if  $v$  and  $\mu$  satisfy

$$Hv = \mu v.$$

##### 1.4.1 Bound for subspace-eigenvector angle

As a first indicator of the convergence we will characterize the following quantity

$$\text{error in eigenvector } x_i \sim \|(I - QQ^*)x_i\| \quad (1.13)$$

where

$$Ax_i = \lambda_i x_i.$$

Recall:  $Q \in \mathbb{C}^{n \times k}$  is an orthogonal matrix which means that  $Q^*Q = I \in \mathbb{C}^{k \times k}$ . However,  $I \neq QQ^* \in \mathbb{C}^{n \times n}$ .

It is very natural to associate the accuracy of the eigenvector with this quantity from a geometric perspective. The indicator in the right-hand side of (1.13) is called (the norm of) the orthogonal complement of the projection of  $x_i$  onto the space spanned by  $Q$  and it can be interpreted as the sine of the canonical angle between the Krylov subspace and an eigenvector. For the moment, we will only justify this indicator with this geometric reasoning and the following observation:

**Lemma 1.4.1.** *Suppose  $(\lambda_i, x_i)$  is an eigenpair  $A$ . If the Krylov subspace contains the eigenvector  $(x_i \in \mathcal{K}_k(A, b))$ , then the indicator vanishes  $\|(I - QQ^*)x_i\| = 0$  and there is at least one Ritz value  $\mu$  such that  $\mu = \lambda_i$ .*

In words:

- Suppose the Krylov subspace contains the eigenvector  $(x_i \in \mathcal{K}_k(A, b))$ . Then, there exists a vector  $z \in \mathbb{C}^k$  such that  $x_i = Qz$ . Moreover, this is an eigenvector of  $H$  such that the Arnoldi method will generate an exact eigenvalue of  $A$ . Moreover, the indicator is  $\|(I - QQ^*)x_i\| = \|(I - QQ^*)Qz\| = 0$ .
- If, similar to above,  $x_i \approx x \in \mathcal{K}_k(A, b)$ , we expect the indicator to be small and an eigenvalue of  $H$  also to be close  $\lambda_i$ .

The indicator can be bounded as follows, where we assume diagonalizability of the matrix.

**Theorem 1.4.2.** *Suppose  $A \in \mathbb{C}^{n \times n}$  is diagonalizable and let the matrix  $X = (x_1, \dots, x_n) \in \mathbb{C}^{n \times n}$  and diagonal matrix  $\Lambda \in \mathbb{C}^{n \times n}$  be the Jordan decomposition such that*

$$A = X\Lambda X^{-1}.$$

Suppose  $\alpha_1, \dots, \alpha_n \in \mathbb{C} \setminus \{0\}$  are such that

$$b = \alpha_1 x_1 + \dots + \alpha_n x_n \quad (1.14)$$

and

$$\varepsilon_i^{(m)} := \min_{\substack{p \in P_{m-1} \\ p(\lambda_i) = 1}} \max(|p(\lambda_1)|, \dots, |p(\lambda_{i-1})|, |p(\lambda_{i+1})|, \dots, |p(\lambda_n)|)$$

where  $P_n$  denotes polynomials of degree  $n$ . Suppose the Arnoldi method does not break down when applied to  $A$  and started with  $b$ . Let  $Q \in \mathbb{C}^{n \times m}$  be the orthogonal basis generated after  $m$  iterations. Then,

$$\|(I - QQ^*)x_i\| \leq \tilde{\zeta}_i \varepsilon_i^{(m)}, \quad (1.15)$$

where

$$\tilde{\zeta}_i = \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|\alpha_j|}{|\alpha_i|}.$$

The Arnoldi method produces an exact approximation if the Krylov subspace contains an eigenvector, or equivalently the indicator is zero.

Recall: The eigenvectors of a diagonalizable matrix form a basis of  $\mathbb{C}^n$ .

The indicator can be bounded by a product consisting of two scalar values:  $\varepsilon_i^{(m)}$  which only depends on the eigenvalues and iteration number; and  $\tilde{\zeta}_i$  only depending on the starting vector and eigenvectors.

*Proof.* The proof consists of three steps.

1. Consider any vector  $u \in \mathbb{C}^n$ . Then

$$\min_{z \in \mathbb{C}^m} \|u - Qz\|_2$$

is a linear least squares problem with a solution given by the normal equations  $Q^*u = Q^*Qz$ . Hence,  $z = Q^*u$ . This implies that (for any vector  $u$ ) we have

$$\min_{z \in \mathbb{C}^m} \|u - Qz\|_2 = \|u - QQ^*u\| = \|(I - QQ^*)u\|$$

2. Although we ultimately want to bound the left-hand side of (1.15), the proof is simplified by considering a scaling the left-hand side of (1.15) with  $\alpha_i$  as follows:

Apply step 1 reversely with  $u = \alpha_i x_i$

$$\begin{aligned} \|(I - QQ^*)\alpha_i x_i\| &= \min_{z \in \mathbb{C}^m} \|\alpha_i x_i - Qz\| \\ &= \min_{y \in \mathcal{K}_m(A, b)} \|\alpha_i x_i - y\| \end{aligned}$$

Now note that the space  $\mathcal{K}_m(A, b)$  can be characterized with polynomials. It is easy to verify that  $y \in \mathcal{K}_m(A, b)$  is equivalent to the existence of a polynomial  $p \in P_{m-1}$  such that  $y = p(A)b$ . Consequently,

$$\|(I - QQ^*)\alpha_i x_i\| = \min_{p \in P_{m-1}} \|\alpha_i x_i - p(A)b\|.$$

3. The final step consists of inserting the expansion of  $b$  in terms of eigenvectors (1.14) and applying appropriate bounds:

$$\begin{aligned}
 \|(I - QQ^*)\alpha_i x_i\| &= \min_{p \in P_{m-1}} \left\| \alpha_i x_i - p(A) \sum_{j=1}^n \alpha_j x_j \right\| \\
 &= \min_{p \in P_{m-1}} \left\| \alpha_i x_i - \sum_{j=1}^n \alpha_j p(\lambda_j) x_j \right\| \\
 &\leq \min_{\substack{p \in P_{m-1} \\ p(\lambda_i)=1}} \left\| \alpha_i x_i - \sum_{j=1}^n \alpha_j p(\lambda_j) x_j \right\| \\
 &= \min_{\substack{p \in P_{m-1} \\ p(\lambda_i)=1}} \left\| \alpha_i x_i - \alpha_i x_i - \sum_{\substack{j=1 \\ j \neq i}}^n \alpha_j p(\lambda_j) x_j \right\| \\
 &= \min_{\substack{p \in P_{m-1} \\ p(\lambda_i)=1}} \left\| \sum_{\substack{j=1 \\ j \neq i}}^n \alpha_j p(\lambda_j) x_j \right\| \\
 &\leq \left( \sum_{\substack{j=1 \\ j \neq i}}^n |\alpha_j| \right) \cdot \min_{\substack{p \in P_{m-1} \\ p(\lambda_i)=1}} \max_{j \neq i} (|p(\lambda_j)|) \\
 &= \left( \sum_{\substack{j=1 \\ j \neq i}}^n |\alpha_j| \right) \cdot \varepsilon_i^{(m)}
 \end{aligned}$$

The conclusion of the theorem is established by dividing the equation by  $|\alpha_i|$ .

□

Note that  $\|b\| = 1$  and  $\|x_1\| = \dots = \|x_n\| = 1$ . Hence the coefficients  $\alpha_1, \dots, \alpha_n$  are balanced. In particular they satisfy

$$1 = \|\alpha_1 x_1 + \dots + \alpha_n x_n\| \leq |\alpha_1| + \dots + |\alpha_n|.$$

and

$$\zeta_i = \frac{1}{|\alpha_i|} \sum_{j=1}^n |\alpha_j| - 1 \geq \frac{1}{|\alpha_i|} - 1$$

From this we can easily identify a very good situation and a very bad situation.

- Suppose for all  $j \neq i$ ,  $\alpha_j = \delta$  and suppose  $\delta$  is small. We have that  $\zeta_i = \frac{(n-1)\delta}{\alpha_i}$ . Due to balancing  $\alpha_i$  cannot be small. Hence,  $\zeta_i$  is small, showing fast convergence for this eigenvalue.
- On the other hand, if  $\alpha_i$  (the component of the starting vector in the direction of the  $i$ th eigenvector) is very small, we have  $\zeta_i \gg 1$  which implies that the right-hand side of (1.15) is large and we have slow convergence.

This serves as a justification for a more general property.



---

**Rule-of-thumb. Starting vector dependency.** The Arnoldi method for eigenvalue problems will “favor” eigenvectors which have large components in the starting vector.

---

The word “favors” is purposely vague. It should be interpreted as the situation that one observes often in practice, but certainly not always. If we have a particular structure in the matrix or starting vector, we might observe convergence to other eigenvalues.

*Bounding  $\varepsilon_i^{(m)}$*

In the characterization of the indicator in Theorem 1.4.2 above we introduced the quantity  $\varepsilon_i^{(m)}$ . This quantity bounds (up to a constant) the error in eigenvector  $x_i$  at iteration  $m$ . Although  $\varepsilon_i^{(m)}$  is defined through a polynomial optimization problem, which is complicated to solve, it is surprisingly easy to use this to obtain bounds providing qualitative understanding of the convergence of the Arnoldi method for eigenvalue problems. We illustrate the power with a specific bound.

Think:  $\varepsilon_i^{(m)}$  measures how “difficult” it is to push down a polynomial in points  $\lambda_j$ , for all  $j \neq i$  and maintain  $p(\lambda_i) = 1$ .

**Corollary 1.4.3.** Suppose  $C(\rho, c) \subset \mathbb{C}$  is a disk centered at  $c \in \mathbb{C}$  with radius  $\rho$  such that it contains all eigenvalues but  $\lambda_1$ . That is,  $\lambda_2, \dots, \lambda_n \in C(\rho, c)$  and  $\lambda_1 \notin C(\rho, c)$ . Then,

$$\varepsilon_1^{(m)} \leq \left( \frac{\rho}{|\lambda_1 - c|} \right)^{m-1}.$$

*Proof.* The proof consists of selecting a particular polynomial in the polynomial optimization problem,

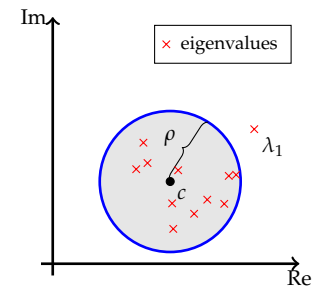
$$\begin{aligned} \varepsilon_1^{(m)} &:= \min_{\substack{p \in P_{m-1} \\ p(\lambda_1) = 1}} \max(|p(\lambda_1)|, \dots, |p(\lambda_{i-1})|, |p(\lambda_{i+1})|, \dots, |p(\lambda_n)|) \\ &= \max_{j \neq i} |q(\lambda_j)|, \end{aligned}$$

for any  $q \in P_{m-1}$  satisfying  $q(\lambda_1) = 1$ , in particular

$$q(z) = \frac{1}{(\lambda_1 - c)^{m-1}} (z - c)^{m-1}.$$

Hence, from the definition of  $\rho$  and  $c$  we have that

$$\begin{aligned} \varepsilon_1^{(m)} &\leq \max_{i > 1} \frac{|\lambda_i - c|^{m-1}}{|\lambda_1 - c|^{m-1}} \\ &\leq \frac{\rho^{m-1}}{|\lambda_1 - c|^{m-1}}. \end{aligned}$$



□

The result can be intuitively interpreted as follows. If we can construct a small disc that encloses all eigenvalues but one eigenvalue we expect fast (at least linear geometric) convergence for that eigenvalue. This can be achieved for an eigenvalue which is well separated from the rest of the eigenvalues and also in an outer part of the spectrum. We call this “extreme” isolated eigenvalues.

---

**Rule-of-thumb. Eigenvalue dependency.** Arnoldi’s method for eigenvalue problems favors convergence to “extreme” isolated eigenvalues.

---

Note the difference between an “extreme” eigenvalue and the eigenvalues which are largest in modulus (absolute value). The Arnoldi method will favor “extreme” whereas the power method will essentially always converge to the eigenvalue largest in modulus.

### 1.4.2 An a posteriori theorem

In the previous section we saw a characterization of the error involving the eigenvectors and eigenvalues of the matrix  $A$ . The following result provides an explicit characterization of  $\|Av - \mu v\|$  where  $(\mu, v)$  is an approximate eigenpair. It is expressed in terms of quantities computed during the iteration.

**Theorem 1.4.4.** Suppose  $Q_k$  and  $\underline{H}_k$  satisfy the Arnoldi relation

$$AQ_k = Q_{k+1}\underline{H}_k \quad (1.16)$$

where  $Q_k \in \mathbb{C}^{n \times k}$  and  $Q_{k+1} = [Q_k, q_{k+1}] \in \mathbb{C}^{n \times (k+1)}$  are orthogonal matrices. Moreover, suppose  $(\mu, v)$  is a Ritz pair such that  $H_k z = \mu z$  and  $v = Q_k z$ . Then,

$$\|Av - \mu v\|_2 = |h_{k+1,k}| |e_k^T z|. \quad (1.17)$$

*Proof.* From the fact that  $(\mu, v)$  is a Ritz pair, we have

$$\begin{aligned} Av - \mu v &= AQ_k z - \mu Q_k z \\ &= (AQ_k - Q_k H_k) z \\ &= h_{k+1,k} q_{k+1} e_k^T z \end{aligned}$$

The conclusion follows from the fact that  $e_k^T z$  is a scalar and  $q_{k+1}$  is normalized since  $Q_{k+1}$  is orthogonal. More precisely,  $\|Av - \mu v\|_2 = |h_{k+1,k}| \|q_{k+1} e_k^T z\| = |h_{k+1,k}| \|q_{k+1}\| |e_k^T z| = |h_{k+1,k}| |e_k^T z|$ .  $\square$

The result can be used to study break-down. Break-down corresponds to the situation where we cannot carry out that Gram-Schmidt orthogonalization process since the new vector is contained in the span

*A priori vs. a posteriori:* Error characterizations can be classified into two types. An *a priori* (latin for “from before”) error estimate involves quantities which are known before the algorithm is carried out. An *a posteriori* (latin for “from after”) error characterization involves quantities computed during the iteration. Theorem 1.4.2 is an *a priori* error bound. Theorem 1.4.4 is an (exact) *a posteriori* error characterization since the right-hand side involves  $H_k$  and  $z$  which are computed from the iteration.

Use  $v = Q_k z$ .

Use that since  $\underline{H}_k$  is a Hessenberg matrix, (1.16) can be written as  $AQ_k = Q_k H_k + h_{k+1,k} q_{k+1} e_k^T$ .

of previous iterations. It implies that the  $y_{\perp} = 0$  and  $\beta = 0$ . This implies in turn that  $h_{k+1,k} = 0$ . Hence, due to (1.17), if we have breakdown the error is already zero and the Ritz pairs are eigenpairs of the original problem.

## 1.5 Shift-and-invert Arnoldi method

We saw that the convergence of the Arnoldi method was given in terms of a polynomial optimization, which in turn gave bounds on the convergence factor; from which we conclude favorable convergence for the outer part of the spectrum. In an application, this may not necessarily be the eigenvalues of interest. This situation is similar to the power method. Similar to the construction of inverse iteration we can transform the problem and use the Arnoldi method on a matrix:

$$B = (A - \mu I)^{-1},$$

where  $\mu$  is called a shift (or target). This is called the shift-and-invert Arnoldi method.

Properties:

- The shift-and-invert Arnoldi method requires a linear solve per iteration, in contrast to the standard Arnoldi method which requires a matrix-vector multiplication.
- The convergence of shift-and-invert Arnoldi method is completely given by the convergence of the Arnoldi method with the transformed matrix  $B$ .
- In contrast to inverse iteration, which is essentially guaranteed to converge to the closest eigenvalue, the shift-and-invert Arnoldi method in practice often converges to eigenvalues close to  $\mu$ , but the precise relationship is more complicated, since the convergence of the Arnoldi method is more complicated than the convergence of the power method.
- The eigenvalue extraction in shift-and-invert Arnoldi method can be done in different ways. The standard approach to extract eigenvalue approximations is to use  $H_m$  such that

$$(A - \mu I)^{-1} Q_m = Q_{m+1} \underline{H}_m. \quad (1.18)$$

Another approach is to use  $G_m = Q_m^T A Q_m$ , where  $Q_m$  is generated from the Arnoldi method applied to  $(A - \mu I)$ .

## 1.6 Literature and further reading

The proof and reasoning above is inspired by [5]. Other convergence bounds involving Schur factorizations, that lead to similar qualitative understanding can be found in [6], where also complications of the non-generic cases are discussed. There are also further characterizations of convergence and the connection with potential theory [4]. In the above reasoning we characterized the angle between the subspace and the eigenvector. Although this serves as a very accurate prediction of the error in practice, it does not directly give a rigorous bound on the accuracy of Ritz pair. Several approaches to describe the convergence of Ritz values and Ritz vectors have been done in for instance [2, 3]. There is also considerable research on the effect of rounding errors in Krylov methods. Unlike many other numerical methods, the effect of finite arithmetic can improve the performance of the algorithm. See also the recent summary of the convergence of the Arnoldi method for eigenvalue problems [1]. The a posteriori error estimate in Theorem 1.4.4 is contained in some recent text-books in numerical linear algebra such as [7].

## 1.7 Appendix: Round-off error analysis of Gram-Schmidt

We now investigate what happens if we have an error in the computation of the Gram-Schmidt coefficients. In other words, we assume that  $h$  is approximated by

$$\tilde{h} = \begin{bmatrix} (1 + \varepsilon_1)h_1 \\ \vdots \\ (1 + \varepsilon_m)h_m \end{bmatrix} = \left( \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} + \underbrace{\begin{bmatrix} \varepsilon_1 & & \\ & \ddots & \\ & & \varepsilon_m \end{bmatrix}}_{\Lambda_\varepsilon} \right) Q^T b \quad (1.19)$$

where  $\varepsilon_1, \dots, \varepsilon_m$  are a small number introduced by the inexact evaluation of  $Q^T b$ , typically of order of the same order of magnitude  $\epsilon_{\text{mach}}$ . Our approximation of  $z$  satisfies

$$\tilde{z} = b - Q\tilde{h} = b - Q\Lambda_\varepsilon Q^T b(1 + \varepsilon) = z - Q\Lambda_\varepsilon Q^T b \quad (1.20)$$

such that

$$\begin{aligned} \tilde{q}_{m+1} &= \frac{1}{\|\tilde{z}\|} \tilde{z} = \frac{1}{\|z - Q\Lambda_\varepsilon Q^T b\|} \tilde{z} = \frac{1}{\sqrt{(z - Q\Lambda_\varepsilon Q^T b)^T (z - Q\Lambda_\varepsilon Q^T b)}} \tilde{z} = \\ &= \frac{1}{\sqrt{(z - Q\Lambda_\varepsilon Q^T b)^T (z - Q\Lambda_\varepsilon Q^T b)}} \tilde{z} = \frac{1}{\sqrt{\|z\|^2 + \|\Lambda_\varepsilon\|^2 \|Q Q^T b\|^2}} \tilde{z} = \\ &= \tilde{z} \left( \frac{1}{\|z\|} + \mathcal{O}(\varepsilon^2) \right), \quad (1.21) \end{aligned}$$

where  $\varepsilon = \|\Lambda_\varepsilon\|$ . The approximation of the new vector is

$$\tilde{q}_{m+1} = (z - Q\Lambda_\varepsilon Q^T b) \left( \frac{1}{\|z\|} + \mathcal{O}(\varepsilon^2) \right) = \frac{z}{\|z\|} - \frac{1}{\|z\|} Q\Lambda_\varepsilon Q^T b + \mathcal{O}(\varepsilon^2) \quad (1.22)$$

In this first-order estimation, we see that the error is small if

$$\frac{\|Q\Lambda_\varepsilon Q^T b\|}{\|z\|} = \frac{\|\Lambda_\varepsilon Q^T b\|}{\|z\|} \leq \varepsilon \frac{\|Q^T b\|}{\|z\|}$$

is small.

A bad situation can easily be identified, since we can construct a situation where  $\|z\|$  is small but  $Q^T b$  is not: Suppose  $b = q + \delta e$  where  $q = Qd$  and  $e \perp Q$  and  $\|e\| = 1$ . A direct computation leads to

$$\|\tilde{q}_{m+1} - \frac{z}{\|z\|}\| \leq \frac{|\varepsilon|}{|\delta|} \|Qd\| + \mathcal{O}(\varepsilon^2).$$

which suggests that the round-off error is proportional to  $|\varepsilon|/|\delta|$ .

## 1.8 References

- [1] M. Bellalij, Y. Saad, and H. Sadok. Further analysis of the Arnoldi process for eigenvalue problems. *SIAM J. Numer. Anal.*, 48(2):393–407, 2010.
- [2] Z. Jia. The convergence of generalized Lanczos methods for large unsymmetric eigenproblems. *SIAM J. Matrix Anal. Appl.*, 16(3):843–862, 1995.
- [3] Z. Jia and G. W. Stewart. On the convergence of ritz values, ritz vectors, and refined ritz vectors. Technical report, 1999.
- [4] A. B. Kuijlaars. Convergence analysis of Krylov subspace iterations with methods from potential theory. *SIAM Rev.*, 48(1):3–40, 2006.
- [5] Y. Saad. *Numerical methods for large eigenvalue problems*. SIAM, 2011.
- [6] G.W. Stewart. *Matrix Algorithms volume 2: eigensystems*. SIAM publications, 2001.
- [7] D. S. Watkins. *Fundamentals of matrix computations*. 3rd ed. Wiley, 2010.