

## 10 The Rayleigh Quotient and Inverse Iteration

From now on we will restrict the discussion to *real symmetric* matrices  $A \in \mathbb{R}^{m \times m}$  whose *eigenvalues*  $\lambda_1, \dots, \lambda_m$  are guaranteed to be *real* and whose *eigenvectors*  $\mathbf{q}_1, \dots, \mathbf{q}_m$  are *orthogonal*. Moreover, in this case Householder reduction will produce a *tridiagonal* matrix.

**Definition 10.1** Let  $A \in \mathbb{R}^{m \times m}$ ,  $\mathbf{x} \in \mathbb{R}^m$ . The quantity

$$r(\mathbf{x}) = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \in \mathbb{R}$$

is called a Rayleigh quotient.

**Remark** If  $\mathbf{x}$  is an eigenvector of  $A$ , then  $A\mathbf{x} = \lambda\mathbf{x}$  and

$$r(\mathbf{x}) = \frac{\lambda \mathbf{x}^T \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \lambda.$$

In general one can show that  $\lambda_{\min} \leq r(\mathbf{x}) \leq \lambda_{\max}$ .

Another fact about the Rayleigh quotient is that it forms the least squares best approximation to the eigenvalue corresponding to  $\mathbf{x}$ . Consider

$$\mathbf{x}\alpha \approx A\mathbf{x}.$$

This is similar to the standard least squares problem  $A\mathbf{x} \approx \mathbf{b}$ . However, now we can ask to find  $\alpha$  such that

$$\|A\mathbf{x} - \mathbf{x}\alpha\|_2$$

is minimized. The associated normal equations (cf.  $A^T A \mathbf{x} = A^T \mathbf{b}$ ) are given by

$$\mathbf{x}^T \mathbf{x} \alpha = \mathbf{x}^T A \mathbf{x} \quad \text{or} \quad \alpha = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = r(\mathbf{x}).$$

Thus, the Rayleigh quotient provides an *optimal estimate for the eigenvalue*. In fact,

$$r(\mathbf{x}) - \underbrace{r(\mathbf{q}_J)}_{=\lambda_J} = \mathcal{O}(\|\mathbf{x} - \mathbf{q}_J\|_2^2) \quad \text{as } \mathbf{x} \rightarrow \mathbf{q}_J.$$

Here  $\mathbf{q}_J$  is the eigenvector associated with  $\lambda_J$ , and the estimate shows that the convergence rate of  $r(\mathbf{x})$  to the eigenvalue is *quadratic*.

In order to prove this fact we need the first-order (multivariate) Taylor expansion of  $r$  about  $\mathbf{x} = \mathbf{q}_J$ :

$$r(\mathbf{x}) = r(\mathbf{q}_J) + (\mathbf{x} - \mathbf{q}_J)^T \nabla r(\mathbf{q}_J) + \mathcal{O}(\|\mathbf{x} - \mathbf{q}_J\|_2^2).$$

If we can show that the linear term is not present, then this identity will imply the quadratic convergence result we wish to prove. We will show that  $\nabla r(\mathbf{q}_J) = \mathbf{0}$ . First,

$$\nabla r(\mathbf{x}) = \left[ \frac{\partial r(\mathbf{x})}{\partial x_1}, \frac{\partial r(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial r(\mathbf{x})}{\partial x_m} \right].$$

One of these partials can be computed using the quotient rule as

$$\frac{\partial r(\mathbf{x})}{\partial x_j} = \frac{\partial}{\partial x_j} \left( \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \right) = \frac{\frac{\partial}{\partial x_j}(\mathbf{x}^T A \mathbf{x}) \mathbf{x}^T \mathbf{x} - \mathbf{x}^T A \mathbf{x} \frac{\partial}{\partial x_j} \mathbf{x}^T \mathbf{x}}{(\mathbf{x}^T \mathbf{x})^2}.$$

Here

$$\begin{aligned} \frac{\partial}{\partial x_j}(\mathbf{x}^T A \mathbf{x}) &= \frac{\partial}{\partial x_j}(\mathbf{x}^T) A \mathbf{x} + \mathbf{x}^T \frac{\partial}{\partial x_j}(A \mathbf{x}) \\ &= [0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0] A \mathbf{x} + \mathbf{x}^T A \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \end{aligned}$$

Here the ones are in the  $j$ -th position so that the fact that we are dealing with scalar quantities (and therefore can transpose taking advantage of the symmetry of  $A$ ) yields

$$\frac{\partial}{\partial x_j}(\mathbf{x}^T A \mathbf{x}) = 2(A \mathbf{x})_j.$$

Similarly, we can compute the other derivative as

$$\frac{\partial}{\partial x_j}(\mathbf{x}^T \mathbf{x}) = [0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0] \mathbf{x} + \mathbf{x}^T \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = 2x_j.$$

Putting the pieces back together we get

$$\begin{aligned} \frac{\partial r(\mathbf{x})}{\partial x_j} &= \frac{2(A \mathbf{x})_j}{\mathbf{x}^T \mathbf{x}} - \frac{(\mathbf{x}^T A \mathbf{x}) 2x_j}{(\mathbf{x}^T \mathbf{x})^2} \\ &= \frac{2}{\mathbf{x}^T \mathbf{x}} ((A \mathbf{x})_j - r(\mathbf{x}) x_j), \end{aligned}$$

and we see that the gradient is given by

$$\nabla r(\mathbf{x}) = \frac{2}{\mathbf{x}^T \mathbf{x}} ((A \mathbf{x}) - r(\mathbf{x}) \mathbf{x}).$$

Finally, letting  $\mathbf{x} = \mathbf{q}_J$  we have  $r(\mathbf{q}_J) = \lambda_J$  and indeed  $\nabla r(\mathbf{q}_J) = \mathbf{0}$  since  $A \mathbf{q}_J = \lambda_J \mathbf{q}_J$ .

## 10.1 Power Iteration

In order to take advantage of the Rayleigh quotient approximation of the eigenvalue we need a good approximation of an eigenvector. We begin with the basic form of the power iteration algorithm.

**Algorithm** (Power Iteration, simple form)

Initialize  $\mathbf{v}^{(0)}$  with an arbitrary nonzero vector

for  $k = 1, 2, \dots$

$$\mathbf{v}^{(k)} = A\mathbf{v}^{(k-1)}$$

end

This algorithm generates a sequence of vectors

$$\mathbf{v}^{(0)}, A\mathbf{v}^{(0)}, A^2\mathbf{v}^{(0)}, \dots$$

If we want to prove that this sequence converges to an eigenvector of  $A$ , the matrix needs to be such that it has a *unique largest eigenvalue*  $\lambda_1$ , i.e.,  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_m| \geq 0$ . There is another technical assumption. The initial vector  $\mathbf{v}^{(0)}$  needs to be chosen such that  $\mathbf{q}_1^T \mathbf{v}^{(0)} \neq 0$ . Otherwise, if  $\mathbf{v}^{(0)}$  is completely perpendicular to the eigenvector  $\mathbf{q}_1$ , the algorithm will not converge. Usually, we do not worry about this assumption since a small roundoff error will already ensure it is satisfied.

We now show that the sequence produced by the power iteration algorithm converges to a multiple of  $\mathbf{q}_1$ , the eigenvector corresponding to  $\lambda_1$ .

Write the initial vector as a linear combination of the eigenvectors of  $A$  (note that they form a basis for  $\mathbb{R}^m$  since  $A$  is assumed to be real symmetric):

$$\mathbf{v}^{(0)} = \sum_{j=1}^m a_j \mathbf{q}_j.$$

Then

$$\begin{aligned} \mathbf{v}^{(k)} &= A^k \mathbf{v}^{(0)} = A^k \sum_{j=1}^m a_j \mathbf{q}_j \\ &= \sum_{j=1}^m a_j A^k \mathbf{q}_j = \sum_{j=1}^m a_j \lambda_j^k \mathbf{q}_j \end{aligned}$$

since the  $\lambda^k$  are eigenvalues of  $A^k$ . Now we factor out the largest eigenvalue, i.e.,

$$\mathbf{v}^{(k)} = \lambda_1^k \left( a_1 \mathbf{q}_1 + \sum_{j=2}^m \left( \frac{\lambda_j}{\lambda_1} \right)^k a_j \mathbf{q}_j \right).$$

Since  $\lambda_j/\lambda_1 < 1$  the sum goes to zero as  $k \rightarrow \infty$ , and  $\mathbf{v}^{(k)}$  does indeed converge to a multiple of  $\mathbf{q}_1$ .

The previous algorithm is not particularly stable. We can improve it by ensure that  $\mathbf{v}^{(k)}$  is always of unit length. This leads to

**Algorithm** (Power Iteration, improved form)

Initialize  $\mathbf{v}^{(0)}$  with an arbitrary vector such that  $\|\mathbf{v}^{(0)}\|_2 = 1$

for  $k = 1, 2, \dots$

$$\begin{aligned}\mathbf{w} &= A\mathbf{v}^{(k-1)} \\ \mathbf{v}^{(k)} &= \mathbf{w} / \|\mathbf{w}\|_2 \\ \lambda^{(k)} &= [\mathbf{v}^{(k)}]^T A\mathbf{v}^{(k)}\end{aligned}$$

end

Note that we have also added the Rayleigh quotient estimate to get the largest eigenvalue of  $A$ .

It can be shown that convergence to the eigenvector is linear, while convergence to the eigenvalue is still quadratic. More precisely,

$$\begin{aligned}\|\mathbf{v}^{(k)} - (\pm \mathbf{q}_1)\|_2 &= \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \\ |\lambda^{(k)} - \lambda_1| &= \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right),\end{aligned}$$

which shows that the speed of convergence depends on the *gap* between the two largest eigenvalues of  $A$ . In particular, if the largest eigenvalue of  $A$  were complex (which it can't be for the real symmetric matrices we are considering), then  $\lambda_2 = \overline{\lambda_1}$  and the algorithm would not converge at all.

In order to find more than one eigenvalue we can perform *simultaneous power iteration* (see Chapter 28 in [Trefethen/Bau] for more details):

**Algorithm** (Simultaneous Power Iteration)

Initialize  $V^{(0)}$  with an arbitrary  $m \times n$  matrix of rank  $n$

for  $k = 1, 2, \dots$

$$V^{(k)} = AV^{(k-1)}$$

end

$$\hat{Q}^{(k)} \hat{R}^{(k)} = V^{(k)}$$

In the last step of the algorithm we perform a reduced QR factorization to obtain a well-behaved basis for the column space of  $V^{(k)}$ .

For this algorithm the speed of convergence will depend on the smallest gap among the first  $n + 1$  eigenvalues.

The problem with this algorithm is that *all* of the columns of  $V^{(k)}$  converge to a multiple of  $\mathbf{q}_1$ , and therefore the column of  $Q^{(k)}$  will form an extremely ill-conditioned basis for  $\text{range}(V^{(k)})$ . The fix is simple, but more expensive. We should orthonormalize at each step:

**Algorithm** (Orthogonal Simultaneous Iteration)

Initialize  $\hat{Q}^{(0)}$  with an arbitrary  $m \times n$  matrix with orthonormal columns

for  $k = 1, 2, \dots$

$$\begin{aligned} Z &= A\hat{Q}^{(k-1)} \\ \hat{Q}^{(k)}\hat{R}^{(k)} &= Z \end{aligned}$$

end

We will see later that this algorithm is very useful. In fact, it is equivalent to the QR iteration algorithm we will study soon (*not* the QR factorization algorithm).

We now return to basic power iteration, and consider a few more modifications. If we want to find the smallest eigenvalue instead of the largest one, then we perform power iteration for  $A^{-1}$  (since the eigenvalues of  $A^{-1}$  are the reciprocals of the eigenvalues of  $A$ ; see homework). Of course, we *do not want to compute*  $A^{-1}$ .

Instead we use the following

**Algorithm** (Inverse Iteration)

Initialize  $\mathbf{v}^{(0)}$  with an arbitrary vector such that  $\|\mathbf{v}^{(0)}\|_2 = 1$

for  $k = 1, 2, \dots$

Solve  $A\mathbf{w} = \mathbf{v}^{(k-1)}$  for  $\mathbf{w}$

$$\mathbf{v}^{(k)} = \mathbf{w} / \|\mathbf{w}\|_2$$

$$\lambda^{(k)} = [\mathbf{v}^{(k)}]^T A\mathbf{v}^{(k)}$$

end

**Remark** For this algorithm the matrix  $A$  needs to be factored *only once* (by Cholesky factorization for a symmetric  $A$ ).

Another modification that can be applied to either inverse iteration (or basic power iteration) is a *shift*  $\mu$ , i.e., we consider  $A - \mu I$  instead of  $A$ . Then

1. The eigenvalues of  $A - \mu I$  are  $\lambda_j - \mu$ , and
2. The eigenvectors of  $A - \mu I$  are still the same as those of  $A$ .

This is clear since

$$(A - \mu I)\mathbf{x} = \underbrace{A\mathbf{x}}_{=\lambda\mathbf{x}} - \underbrace{\mu I\mathbf{x}}_{=\mu\mathbf{x}} = (\lambda - \mu)\mathbf{x}.$$

The resulting algorithm (for inverse iteration is)

**Algorithm** (Inverse Iteration with Shift)

Initialize  $\mathbf{v}^{(0)}$  with an arbitrary vector such that  $\|\mathbf{v}^{(0)}\|_2 = 1$

for  $k = 1, 2, \dots$

Solve  $(A - \mu I)\mathbf{w} = \mathbf{v}^{(k-1)}$  for  $\mathbf{w}$

$$\mathbf{v}^{(k)} = \mathbf{w} / \|\mathbf{w}\|_2$$

$$\lambda^{(k)} = [\mathbf{v}^{(k)}]^T A \mathbf{v}^{(k)}$$

end

This algorithm will yield the eigenvalue *closest to*  $\mu$ . This means by picking appropriate shifts  $\mu$ , *any one* eigenvalue of  $A$  can be found. The rate of convergence to the eigenvector is still linear, and that to the eigenvalue is quadratic.

**Remark** If  $\mu = \lambda$ , i.e., one runs the algorithm with a known eigenvalue, then one step of inverse iteration will produce the associated eigenvector.

**Remark** Shifter power iteration — while theoretically possible — is not very useful since it converges to the eigenvalue *farthest away* from  $\mu$ .

## 10.2 A Cubically Convergent Improvement

If we update the estimate  $\mu$  for the eigenvalue with the Rayleigh quotient at each iteration we can get a cubically convergent algorithm:

**Algorithm** (Rayleigh Quotient Iteration)

Initialize  $\mathbf{v}^{(0)}$  with an arbitrary vector such that  $\|\mathbf{v}^{(0)}\|_2 = 1$

Initialize  $\lambda^{(0)} = [\mathbf{v}^{(0)}]^T A \mathbf{v}^{(0)}$

for  $k = 1, 2, \dots$

Solve  $(A - \lambda^{(k-1)}I)\mathbf{w} = \mathbf{v}^{(k-1)}$  for  $\mathbf{w}$

$\mathbf{v}^{(k)} = \mathbf{w} / \|\mathbf{w}\|_2$

$\lambda^{(k)} = [\mathbf{v}^{(k)}]^T A \mathbf{v}^{(k)}$

end

One can show that for almost all starting vectors we have

$$\begin{aligned} \|\mathbf{v}^{(k+1)} - (\pm \mathbf{q}_J)\|_2 &= \mathcal{O}\left(\|\mathbf{v}^{(k)} - (\pm \mathbf{q}_J)\|_2^3\right) \\ |\lambda^{(k+1)} - \lambda_J| &= \mathcal{O}\left(|\lambda^{(k)} - \lambda_J|^3\right), \end{aligned}$$

Since this is a cubically convergent algorithm one can expect the number of correct digits to *triple* in each iteration. This is illustrated in the MATLAB script `RayleighQuotient.m`. However, each iteration of the algorithm is fairly expensive since we need to solve a linear system with *different* system matrix in each iteration. Also, while the algorithm does usually converge, it need not converge to the eigenvalue closest to the initial shift  $\lambda^{(0)}$ .

**Remark** Cubically convergent algorithms are very rare in numerical algorithm (recall the the quadratically Newton iteration for finding roots of a nonlinear function is already considered a big deal).

**Remark** In order to find *all* eigenvalues and eigenvectors of  $A$  we can use *deflation* (which we will discuss a bit more in the next section in the context of QR iteration).