

On the Powers of a Matrix with
Perturbations*G. W. Stewart[†]

December 2001

Revised January 2002

ABSTRACT

Let A be a matrix of order n . The properties of the powers A^k of A have been extensively studied in the literature. This paper concerns the perturbed powers

$$P_k = (A + E_k)(A + E_{k-1}) \cdots (A + E_1),$$

where the E_k are perturbation matrices. We will treat three problems concerning the asymptotic behavior of the perturbed powers. First, determine conditions under which $P_k \rightarrow 0$. Second, determine the limiting structure of P_k . Third, investigate the convergence of the power method with error: that is, given u_1 , determine the behavior of $u_k = \nu_k P_k u_1$, where ν_k is a suitable scaling factor.

*This report is available by anonymous ftp from `thales.cs.umd.edu` in the directory `pub/reports` or on the web at `http://www.cs.umd.edu/~stewart/`.

[†]Department of Computer Science and Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742 (`stewart@cs.umd.edu`).

On the Powers of a Matrix with Perturbations

G. W. Stewart

ABSTRACT

Let A be a matrix of order n . The properties of the powers A^k of A have been extensively studied in the literature. This paper concerns the perturbed powers

$$P_k = (A + E_k)(A + E_{k-1}) \cdots (A + E_1),$$

where the E_k are perturbation matrices. We will treat three problems concerning the asymptotic behavior of the perturbed powers. First, determine conditions under which $P_k \rightarrow 0$. Second, determine the limiting structure of P_k . Third, investigate the convergence of the power method with error: that is, given u_1 , determine the behavior of $u_k = \nu_k P_k u_1$, where ν_k is a suitable scaling factor.

1. Introduction

Let A be a matrix of order n with eigenvalues $\lambda_1, \dots, \lambda_n$ ordered so that

$$|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_n|$$

and let $\rho(A) = |\lambda_1|$ denote the spectral radius of A . We will be concerned with extending the following three results about the behavior of the powers A^k of A . (These results are easily proved by exploiting the relations between norms and spectral radii; see [11, Sections I.2, II.1]. For earlier work on powers of a matrix see [3, 4, 6, 9].)

- The first result is classic. If $\rho(A) < 1$, then $\lim_{k \rightarrow \infty} A^k = 0$. Moreover, in any norm $\|A^k\|^{1/k} \rightarrow \rho(A)$, or equivalently the convergence of A^k to zero is faster than that of $[\rho(A) + \eta]^k$ for any $\eta > 0$. We say that the *root convergence index* of A^k is $\rho(A)$.
- The second result concerns the asymptotic form of A^k . Let $|\lambda_1| > |\lambda_2|$, and let the right and left eigenvectors corresponding to λ_1 be x and y , normalized so that $y^H x = 1$. Then

$$\lambda_1^{-k} A^k \rightarrow xy^H.$$

Moreover, the root convergence index is $|\lambda_2/\lambda_1|$.

- The third result concerns the convergence of the power method. Specifically, let u_1 be given and define

$$u_{k+1} = \nu_k A u_k, \quad k = 1, 2, \dots,$$

where ν_k is a normalizing factor (e.g., $\|Au_k\|^{-1}$). In the above notation, if $|\lambda_1| > |\lambda_2|$ and $y^H u_1 \neq 0$, then u_k , suitably scaled, converges to a multiple of x . The root index of convergence is $|\lambda_2/\lambda_1|$.

Now let E_1, E_2, \dots be a sequence of perturbation matrices and let

$$P_k = (A + E_k)(A + E_{k-1}) \cdots (A + E_1).$$

The purpose of this paper is to extend the three results above to the *perturbed powers* P_k .¹ Regarding the first result, we will show that if $\rho(A) < 1$ then for sufficiently small E_k , the P_k approach zero. Moreover, by making the E_k small enough we can bring the convergence ratio arbitrarily near $\rho(A)$. Regarding the second result, we will show that if

$$\sum_k \|E_k\| < \infty \tag{1.1}$$

in any norm then the P_k converge to xz^H for some z (which may be zero), and we will investigate the convergence rate. Finally, a finite-precision implementation of the power method results in perturbations E_k are of the order of the rounding unit and therefore do not satisfy (1.1). Thus, we cannot use the second result to analyze the convergence (actually nonconvergence) of the power method in the presence of rounding error. However, using other techniques we can show that in the presence of rounding error the power method will converge up to a point and then stagnate.

There is not a large literature on perturbed matrix powers. Ostrowski [9, Chapter 20] gives bounds on the size of the perturbed powers — results that imply that if $\rho(A) < 1$ and the powers are sufficiently small, then $P_k \rightarrow 0$. Higham and Knight [8] give a very detailed investigation of when $P_k \rightarrow 0$, along with some useful references. These results relate most naturally to our Theorem 3.1. However, the results Sections 4–5 do more than treat the convergence of P_k to zero — they reveal the structure of the matrices P_k , at the cost of additional hypotheses on the nature of the spectrum of A .

This paper is organized as follows. In establishing our extensions it will prove convenient to transform our matrices by certain similarity transformations, and any conditions placed on the transformed perturbations must be translated back to the original problem. Since the transformations can be ill conditioned, it is important to understand the source of the ill-conditioning. Accordingly, the next section is devoted to describing the two transformations we will use. In Section 3 we will establish our extension of the first result, and in Section 4 the extension of the second. In Section 5 we will give an analysis of the power method. It is worth noting that the last two sections end with a little hook: each provides a new result about the problem treated in the preceding section.

¹The phrase “perturbed powers” is, strictly speaking, a misnomer, since it is the factors, not the powers that are perturbed.

Throughout this paper the j th column of the identity matrix will be denoted by \mathbf{e}_j . In addition, $\|\cdot\|$ will denote a consistent family of norms such that that $\|A\|$ bounds the norm of any submatrix of A and $\|\text{diag}(d_1, \dots, d_n)\| = \max_i |d_i|$. This class includes the 1-, 2-, ∞ - norms but excludes the Frobenious norm [5, 10].

As above, we will use the *root convergence index* to measure speed of convergence. Specifically, if a_k is a sequence converging to zero, and $\rho = \limsup_k |a|^{1/k} < 1$ we say that a_k converges with root index ρ .

2. Two transformations

In deriving our results we will have to transform A into $\hat{A} = X^{-1}AX$, for some X appropriate to the problem at hand. In this case, we must also transform the perturbation matrices: $\hat{E}_k = X^{-1}E_kX$. Now $\|\hat{E}_k\| \leq \kappa(X)\|E_k\|$, where $\kappa(X) = \|X\|\|X^{-1}\|$ is the condition number of X . Hence in order to insure that a bound like $\|\hat{E}_k\| \leq \gamma$ holds we have to require that $\|E_k\| \leq \gamma/\kappa(X)$; that is, the bound on $\|E_k\|$ must be stronger by a factor of $\kappa^{-1}(X)$. Thus it is appropriate to examine the conditions under which $\kappa(X)$ is large — that is, under which the transformations are ill conditioned. There are two classes of transformations.

The first transformation is described in the following classic theorem (see, e.g., [11, Theorem I.2.8]).

Theorem 2.1. *For any $\eta > 0$ there is a matrix X such that*

$$\|X^{-1}AX\| \leq \rho(A) + \eta. \tag{2.1}$$

The theorem is proved by first transforming A to Schur form; i.e.,

$$U^H A U = T,$$

where U is unitary and T is upper triangular. If we then set $D_\alpha = \text{diag}(1, \alpha, \dots, \alpha^{n-1})$, the superdiagonal elements of $D^{-1}TD$ are $\alpha^{j-i}\tau_{ij}$. If we define $X_\alpha = UD_\alpha$, then as $\alpha \rightarrow 0$ we have $\|X_\alpha^{-1}AX_\alpha\| \rightarrow \|\text{diag}(\tau_{11}, \dots, \tau_{nn})\| = \rho(A)$. Consequently, by the continuity of norms we may set $X = X_\alpha$, where α is chosen so that (2.1) is satisfied.

As α decreases, X_α^{-1} becomes large. If the off-diagonal elements of T are nonzero, α becomes small along with ϵ , the more so in proportion as the off diagonal elements of T are large. Large diagonal elements in T are associated with Henrici's measure of nonnormality [6]. Hence, nonnormality in A may weaken our theorems. However, it should be stressed, that the above construction of X is designed to accommodate the worst possible case, and in particular situations there may be a better way to construct T . For example, if A has a complete, well-conditioned system of eigenvectors, then the matrix X of eigenvectors reduces A to diagonal form, so that (2.1) is satisfied for $\eta = 0$.

The following useful theorem describes second of our transformations.

Theorem 2.2. *Let λ_1 be a simple eigenvalue of A with right and left eigenvectors x_1 and y_1 normalized so that $\|x_1\|_2 = \|y_1\|_2 = 1$ and $\gamma = y_1^H x_1$ is positive. Let $\sigma = \sqrt{1 - \gamma^2}$. Then there is a matrix U with*

$$\kappa(U) = \frac{1 + \sigma}{\gamma} \quad (2.2)$$

in the 2-norm such that

$$U^{-1}AU = \begin{pmatrix} \lambda_1 & 0 \\ 0 & B \end{pmatrix}. \quad (2.3)$$

Proof. By appealing to the CS decomposition [5, 10], we can find orthonormal matrices $(x_1 \ x_2 \ X_3)$ and $(y_1 \ y_2 \ Y_3)$ such that

$$(x_1 \ x_2 \ X_3)^T (y_1 \ y_2 \ Y_3) = \begin{pmatrix} \gamma & \sigma & 0 \\ \sigma & \gamma & 0 \\ 0 & 0 & I \end{pmatrix}.$$

Let $U = (\gamma^{-\frac{1}{2}}x_1 \ \gamma^{-\frac{1}{2}}y_2 \ X_3)$ and $V = (\gamma^{-\frac{1}{2}}y_1 \ \gamma^{-\frac{1}{2}}x_2 \ Y_3)$. Then it is easily verified that $V^H = U^{-1}$. Moreover,

$$U^T U = \begin{pmatrix} 1/\gamma & \sigma/\gamma & 0 \\ \sigma/\gamma & 1/\gamma & 0 \\ 0 & 0 & I \end{pmatrix}$$

Thus $\|U\|_2^2 = \|U^T U\|_2 = (1 + \sigma)/\gamma$. Similarly $\|V\|_2^2 = (1 + \sigma)/\gamma$, which establishes (2.2). The fact that (2.3) is satisfied follows immediately from the fact that the first column of U is the right eigenvector x_1 and the first column of V is the left eigenvector y_1 . ■

The matrix U will be ill conditioned when γ^{-1} is small. But γ^{-1} — the secant of the angle between x_1 and y_1 — is a condition number for the eigenvalue λ_1 [5, Section 7.2.2], [11, Section I.3.2]. Thus U_1 will be ill conditioned precisely when λ_1 is.

3. Convergence to zero

We are now in the position to state and prove our extension of the first result. In fact, thanks to results of the last section, it is trivial to establish the following theorem. As we indicated in the introduction, this theorem is essentially due to Ostrowski [9, Chapter 20] and has been extended by Higham and Knight [8].

Theorem 3.1. *Let $\rho(A) < 1$ and consider the perturbed products*

$$P_k = (A + E_k)(A + E_{k-1}) \cdots (A + E_1). \quad (3.1)$$

For every $\eta > 0$ there is an $\epsilon > 0$ such that if $\|E_k\| \leq \epsilon$ then

$$\limsup_k \|P_k\|^{1/k} \leq \rho(A) + \eta. \quad (3.2)$$

Hence if $\rho(A) + \eta < 1$ then P_k converges to zero with root convergence index at greatest $\rho(A) + \eta$.

Proof. By Theorem 2.1 there is a nonsingular matrix X such that if $\hat{A} = X^{-1}AX$ then $\|\hat{A}\| \leq \rho(A) + \eta/2$. Let $\hat{E}_k = X^{-1}E_kX$, $\hat{P}_k = X^{-1}P_kX$, and $\hat{\epsilon} = \eta/2$. Then if $\|\hat{E}_k\| \leq \hat{\epsilon}$ we have $\|\hat{A} + \hat{E}_k\| \leq \rho(A) + \eta$, and $\|\hat{P}_k\| \leq [\rho(A) + \eta]^k$.

Transforming back to the original problem, we see that if $\|E_k\| \leq \epsilon \equiv \hat{\epsilon}/\kappa(X)$ then $\|P_k\| \leq \kappa(X)[\rho(A) + \eta]^k$. The inequality (3.2) now follows on taking k th roots. ■

There is little to add to this theorem. The price we pay for the perturbations is that to make the root convergence index approach $\rho(A)$ we must increasingly restrict the size of the perturbations. This is unavoidable. For if we fix the size of the error at ϵ we can always find E such that the largest eigenvalue of $A + E$ has magnitude $\rho(A) + \epsilon$. If we set $E_k = E$, then $P_k = (A + E)^k$, and the best root convergence index we can hope for is $\rho(A) + \epsilon$.

4. Convergence with a simple dominant eigenvalue

In this section we will treat the behavior of the perturbed powers when A has a single, simple dominant eigenvalue λ ; i.e., when $|\lambda_1| > |\lambda_2|$. By dividing by A by λ_1 , we may assume that $\lambda_1 = 1$. The basic result is given in the following theorem.

Theorem 4.1. *Let $1 = \lambda_1 > |\lambda_2|$ and let the right eigenvector corresponding to λ_1 be x . Let P_k be defined as in (3.1). If*

$$\sum_{i=1}^{\infty} \|E_k\| < \infty, \tag{4.1}$$

then for some (possibly zero) vector z we have

$$\lim_{k \rightarrow \infty} P_k = xz^H.$$

The root convergence index is not greater than

$$\max\{\rho, \sigma\},$$

where ρ is the largest of the magnitudes of the subdominant eigenvalues of A and

$$\sigma = \limsup \|E_k\|^{1/k}. \tag{4.2}$$

Proof. By Theorem 2.2, we may transform A so that it has the form $\text{diag}(1, B)$, where $\rho(B) < 1$. By Theorem 2.1, we may assume that $\|B\| \leq \beta < 1$. Note that the right eigenvector of the transformed matrix is \mathbf{e}_1 .

The theorem is best established in a vectorized form. First write

$$P_{k+1} = (A + E_k)P_k = AP_k + E_kP_k$$

Let $u \neq 0$ be given and let $p_k = P_k u$. Then

$$p_{k+1} = Ap_k + E_k p_k. \quad (4.3)$$

We will use this recurrence to show that the p_k approach a multiple of \mathbf{e}_1 .

Our first job is to find a condition that insures that the p_k remain bounded. To this end, partition (4.3) in the form

$$\begin{pmatrix} p_1^{(k+1)} \\ p_2^{(k+1)} \end{pmatrix} = \begin{pmatrix} 1 + e_{11}^{(k)} & e_{12}^{(k)H} \\ e_{21}^{(k)} & B + E_{22}^{(k)} \end{pmatrix} \begin{pmatrix} p_1^{(k)} \\ p_2^{(k)} \end{pmatrix}. \quad (4.4)$$

Now let $\epsilon_k = \|E_k\|$ and let $\pi_1^{(k)}$ and $\pi_2^{(k)}$ be upper bounds on $\|p_1^{(k)}\|$ and $\|p_2^{(k)}\|$. Then by taking norms in (4.4) we see that the components of

$$\begin{pmatrix} \pi_1^{(k+1)} \\ \pi_2^{(k+1)} \end{pmatrix} = \begin{pmatrix} 1 + \epsilon_k & \epsilon_k \\ \epsilon_k & \beta + \epsilon_k \end{pmatrix} \begin{pmatrix} \pi_1^{(k)} \\ \pi_2^{(k)} \end{pmatrix} \quad (4.5)$$

are upper bounds on $\|p_1^{(k+1)}\|$ and $\|p_2^{(k+1)}\|$. Thus if we write

$$\begin{pmatrix} 1 + \epsilon_k & \epsilon_k \\ \epsilon_k & \beta + \epsilon_k \end{pmatrix} = \text{diag}(1, \beta) + \begin{pmatrix} \epsilon_k & \epsilon_k \\ \epsilon_k & \epsilon_k \end{pmatrix} \equiv \text{diag}(1, \beta) + H_k,$$

then the p_k will be bounded provided the product

$$\prod_{k=1}^{\infty} \|\text{diag}(1, \beta) + H_k\| < \infty.$$

Now

$$\prod_{k=1}^{\infty} \|\text{diag}(1, \beta) + H_k\| \leq \prod_{k=1}^{\infty} (\|\text{diag}(1, \beta)\| + \|H_k\|) \leq \prod_{k=1}^{\infty} (1 + 4\epsilon_k).$$

It is well known that the product on the right is finite if and only if the series $\sum_k \epsilon_k$ converges [1, Section 5.2.2]. Hence a sufficient condition for the p_k to remain bounded is for (4.1) to be satisfied.

The next step is to show that $p_2^{(k)}$ converges to zero. Let π be a uniform upper bound on $\|p_1^{(k)}\|$ and $\|p_2^{(k)}\|$. From (4.5) we have

$$\pi_2^{(k+1)} \leq (2\epsilon_k \pi + \beta \pi_2^{(k)}).$$

Hence if we set $\pi = \hat{\pi}_1$ and define $\hat{\pi}_k$ by the recurrence

$$\hat{\pi}_{k+1} = \beta \hat{\pi}_k + 2\epsilon_k \pi,$$

we have that $\|p_2^{(k)}\| \leq \hat{\pi}_k$. But it is easy to see that if we define $\epsilon_0 = \frac{1}{2}$ then

$$\hat{\pi}_{k+1} = 2\pi(\beta^k \epsilon_0 + \beta^{k-1} \epsilon_1 + \cdots + \beta^1 \epsilon_{k-1} + \epsilon_k). \quad (4.6)$$

It follows that

$$\hat{\pi}_1 + \hat{\pi}_2 + \cdots = 2\pi(1 + \beta + \beta^2 + \cdots)(\epsilon_0 + \epsilon_1 + \epsilon_2 + \cdots) \quad (4.7)$$

But the geometric series in β on the right is absolutely convergent, and by (4.1) the series in ϵ_k is also. Thus the series on the left is absolutely convergent, and its terms $\hat{\pi}_k$ must converge to zero.

We must next show that $p_1^{(k)}$ converges. From the first row of (4.4) we have

$$p_1^{(k+1)} - p_1^{(k)} = e_{11}^{(k)} p_1^{(k)} + e_{12}^{(k)} p_2^{(k)}, \quad (4.8)$$

whence

$$|p_1^{(k+1)} - p_1^{(k)}| \leq 2\epsilon_k \pi. \quad (4.9)$$

Since $\sum_k \epsilon_k$ converges, if we set $p_1^{(0)} = 0$, the telescoping series $\sum_{j=0}^k (p_1^{(j+1)} - p_1^{(j)}) = p_1^{(k+1)}$ converges.

By taking $u = \mathbf{e}_i$, we find that the i th column of P_k converges to $\bar{w}_i \mathbf{e}_1$ for some w_i . Consequently, if we set $w^H = (\bar{w}_1 \cdots \bar{w}_n)$, then P_k converges to $\mathbf{e}_1 w^H$.

Finally, in assuming that $A = \text{diag}(1, B)$ we transformed the original matrix A by a similarity transformation X whose first column was the dominant eigenvector of A . It follows that the original P_k converge to

$$X \mathbf{e}_1 w^H X^{-1} = x z^H,$$

where $z^H = X^{-1} w^H$.

We now turn to the rates of convergence. The inequality (4.9) shows that $p_1^{(k)}$ converges as fast as $\|E_k\|$ approaches zero; i.e., its root convergence index is not greater than σ defined by (4.2).

To analyze the convergence of $p_2^{(k)}$ we make the observation that the reciprocal of the radius of convergence of any function $f(z)$ that is analytic at the origin is $\alpha = \limsup_k |a_k|^{1/k}$, where a_k is the k th coefficient in the power series of f [1, Section II.2.4]. We also note that in the expression (4.6) for $\hat{\pi}_{k+1}$ we can replace β by $\|B^k\|$ and still have an upper bound on $\|p_2^{(k+1)}\|$. Now let $r(\zeta) = \sum_k \|B^k\| \zeta^k$ and $s(\zeta) = \sum_k \|E_k\| \zeta^k$. Since $\|B^k\|^{1/k} \rightarrow \rho(B) = \rho$, we know that the radius of convergence of r is ρ^{-1} . By definition the radius of convergence of s is σ^{-1} . But by (4.7), $\limsup \hat{\pi}_k^{1/k}$ is the reciprocal of the radius of convergence of the function $p(\zeta) = r(\zeta)s(\zeta)$. Since the radius of convergence of p is at least as great as the smaller of ρ^{-1} and σ^{-1} , the root index of convergence of $p_2^{(k)}$ is not greater than $\max\{\rho, \sigma\}$. ■

There are four comments to be made about this theorem.

- By the equivalence of norms, if the condition (4.1) on the E_k holds for one norm, it holds for any norm. Thus, the condition on the errors does not depend on the similarity transformation we used to bring A into the form $\text{diag}(1, B)$. But this happy state of affairs obtains only because (4.1) is an asymptotic statement. In practice, the sizes of the initial errors, which do depend on the transformation, may be important.
- Since P_k converges to xz^H , if $z \neq 0$, at least one column of P_k contains an increasingly accurate approximation to x . In the error free case, z is equal to the left eigenvector of A , which is by definition nonzero. In general, however, we cannot guarantee that $z \neq 0$, and indeed it is easy to contrive examples for which z is zero.

However, it follows from (4.8) that

$$|p_1^{(k+1)}| \geq |p_1^{(k)}| - 2\pi\epsilon_k \geq |p_1^{(1)}| - 2\pi(\epsilon_k + \dots + \epsilon_1).$$

Hence if $2\pi \sum_k \epsilon_k < \|p_1^{(1)}\|$, then $\lim_k p_1^{(k)} \neq 0$, and hence $\lim_k P_k \neq 0$.

- The proof can be extended to the case where A has more than one dominant eigenvalue, provided they are all simple. The key is to use a generalization of Theorem 2.2 that uses bases for the left and right dominant eigenspaces of A , to reduce A to the form $\text{diag}(D, B)$, where $|D| = I$. The quantities $p_1^{(k)}$ and $p_1^{(k+1)}$ in (4.4) are no longer scalars, but the recursion (4.5) for upper bounds remains the same, as does the subsequent analysis.
- We have been interested in the case where A has a simple dominant eigenvalue of one. However, the proof of the theorem can easily be adapted to the case where $\rho(A) < 1$ with no hypothesis of simplicity (it is essentially the analysis of $p_2^{(k)}$ without the contributions from $p_1^{(k)}$). The result is the following corollary.

Corollary 4.2. *Let $\rho(A) < 1$ and let E_k satisfy (4.1). Then $P_k \rightarrow 0$ and the root convergence index is not greater than $\max\{\rho, \sigma\}$.*

5. The power method

The power method starts with a vector u_1 and generates a sequence of vectors according to the formula

$$u_{k+1} = \nu_k A u_k,$$

where ν_k is a normalizing factor. If A has a simple dominant eigenvalue (which we may assume to be one), under mild restrictions on u_1 , the u_k converge to the dominant eigenvector of A .

A backward rounding-error analysis shows that in the presence of rounding error we actually compute

$$u_{k+1} = \nu_k (A + E_k) u_k = (\nu_k \cdots \nu_1) P_k u_1.$$

where $\|E_k\|/\|A_k\|$ is of the order of the rounding unit [7, 10]. Theorem 4.1 is not well suited to analyzing this method for two reasons. First the E_k will all be roughly the same size, so that the condition (4.1) is not satisfied. But even if it were, it is possible for the P_k to approach zero while at the same time the normalized vectors u_k converge to a nonzero limit, in which case Theorem 4.1 says nothing useful. Accordingly, in this section we give a different convergence analysis for the power method.

As in the last section we will assume that $A = \text{diag}(1, B)$, where $\|B\| = \beta < 1$. Let $\epsilon_k = \|E_k\|$. We will normalize the u_k so that the first component is one and write

$$u_k = \begin{pmatrix} 1 \\ h_k \end{pmatrix}.$$

It is important to have some appreciation of the magnitudes of the quantities involved. If the computations are being done in IEEE double precision, ϵ will be around $\sqrt{n} \cdot 10^{-16}$; e.g., 10^{-14} if $n = 10,000$. If u_1 is a random vector, we can expect $\|h_1\|$ to be of order \sqrt{n} ; e.g., 100, if $n = 10,000$. Finally, since the ratio of convergence of the power method is approximately β , β must not be too near one; e.g., 0.99 gives unacceptably slow convergence. Thus we may assume that $\epsilon \|h_1\|$ and $\epsilon/(1 - \beta)$ are small.

Let η_k be an upper bound for $\|h_k\|$. We will derive an upper bound η_{k+1} for $\|h_{k+1}\|$, in the form of the quotient of a lower bound on the first component of $(A + E_k)u_k$ and an upper bound on the rest of the vector. We have

$$(A + E_k)u_k = \begin{pmatrix} 1 + e_{11}^{(k)} & e_{12}^{(k)H} \\ e_{21}^{(k)} & B + E_{22}^{(k)} \end{pmatrix} \begin{pmatrix} 1 \\ h_k \end{pmatrix}.$$

A lower bound on the first component of this vector is

$$1 - (1 + \eta_k)\epsilon_k$$

and an upper bound on the lower part is

$$\beta\eta_k + \epsilon_k(1 + \eta_k).$$

Hence

$$\|h_{k+1}\| \leq \eta_{k+1} \equiv \frac{\beta\eta_k + \epsilon_k(1 + \eta_k)}{1 - (1 + \eta_k)\epsilon_k}$$

Let

$$\varphi_\epsilon(\eta) = \frac{\beta\eta + \epsilon(1 + \eta)}{1 - (1 + \eta)\epsilon}.$$

so that

$$\eta_{k+1} = \varphi_{\epsilon_k}(\eta_k). \quad (5.1)$$

It is easily verified that if

$$c = \frac{1 - \beta - 2\epsilon}{\epsilon} \geq 2 \quad (5.2)$$

then φ_ϵ has a minimal fixed point

$$\eta_* = \frac{2}{c + \sqrt{c^2 - 4}}. \quad (5.3)$$

Moreover,

$$\varphi'_\epsilon(\eta) = \frac{\beta}{1 - (1 + \eta)\epsilon} + \frac{\beta\eta + \epsilon(1 + \eta)}{[1 - (1 + \eta)\epsilon]^2}\epsilon. \quad (5.4)$$

The following theorem gives conditions under which we can iterate for the fixed point η_* .

Theorem 5.1. *If*

$$\epsilon < \frac{1}{4}(1 - \beta), \quad (5.5)$$

then φ_ϵ has a unique smallest fixed point η_ given by (5.3) with $0 < \varphi'_\epsilon(\eta_*) < 1$. Moreover, if*

$$\eta_* < \eta_1 < \min \left\{ \frac{1}{2\epsilon} - 1, \frac{1}{\epsilon} \cdot \frac{1 - \beta - \epsilon - 2\epsilon^2}{1 + 2\beta + 2\epsilon} \right\}, \quad (5.6)$$

then the iteration

$$\eta_{k+1} = \varphi_\epsilon(\eta_k), \quad k = 1, 2, \dots, \quad (5.7)$$

converges from above to η_ .*

Proof. The condition (5.5) insures that c defined by (5.2) is less than two. Hence η_* given by (5.3) is the required fixed point. We have $\varphi_\epsilon(0) = \epsilon/(1 - \epsilon) > 0$, and

$$\varphi'_\epsilon(0) = \frac{\beta + \epsilon}{1 - \epsilon} + \frac{\epsilon^2}{(1 - \epsilon)^2} > 0.$$

Moreover $\varphi'(\eta)$ is strictly increasing. Hence the curve $\varphi(\eta)$ cannot cross the line $y = x$ unless its derivative at the crossing point is positive and less than one.

From the theory of fixed point iteration [2, Theorem 6.5.1] we know that the iteration (5.7) converges provided we start with η_1 in an interval $[\eta_*, \tau)$ for which φ'_ϵ is less than one. To compute such an interval, restrict η so that

$$\eta < \frac{1}{2\epsilon} - 1,$$

whence

$$\varphi'_\epsilon(\eta_*) \leq \varphi'_\epsilon(\eta) < \frac{\beta}{1 - (1 + \eta)\epsilon} + 2\frac{\beta\eta + \epsilon(1 + \eta)}{1 - (1 + \eta)\epsilon}\epsilon.$$

Setting this bound equal to one and solving for η , we get a solution τ satisfying

$$\tau = \frac{1}{\epsilon} \cdot \frac{1 - \beta - \epsilon - 2\epsilon^2}{1 + 2\beta + 2\epsilon}. \tag{5.8}$$

Note that since $1 - \beta > 4\epsilon$ and $\epsilon < \frac{1}{4}$, the numerator in (5.8) is positive. It follows that for η in $[\eta_*, \tau)$, we have $\varphi'_\epsilon(\eta) < 1$. ■

This theorem does not apply directly to the power method, since the bounding iteration (5.1) involves the varying errors ϵ_k . But owing to the monotonicity of φ_ϵ , if ϵ is an upper bound on the ϵ_k , then the iteration (5.7) provides a upper bounds on the quantities $\|h_k\|$. To the extent that these upper bounds reflect reality, the theorem has important things to say about the power method.

We are chiefly interested in the behavior of the iteration for small ϵ . In this case, the nearness of β to one influences the iteration in four ways.

- **Tolerable error.** The condition (5.5) — $\epsilon < \frac{1}{4}(1 - \beta)$ — suggests that for the power method to be effective the size of the error in A must decrease with $1 - \beta$. For our typical values, this is no problem, since $\epsilon = 10^{-14} < .01 = 1 - \beta$.
- **Size of the fixed point.** When ϵ is small, $\eta_* \cong \epsilon/(1 - \beta)$. Thus as β approaches 1, the limiting accuracy of the power method will be degraded.
- **Rate of convergence.** For small ϵ , φ_ϵ essentially a straight line with slope β over a wide range above η_* . Thus a β near one implies slow convergence.

• **Size of the convergence region.** For small ϵ the two expressions for the upper end τ of the convergence region in (5.6) become essentially

$$\frac{1}{2\epsilon} \quad \text{and} \quad \frac{1}{\epsilon} \cdot \frac{1-\beta}{1+2\beta}$$

It is seen that for $\beta < \frac{1}{4}$, the first expression determines τ , while otherwise the second expression determines τ . In particular, as β approaches one, the size of the convergence region decreases. For our typical parameters this fact is unimportant, since $\tau \cong 3 \cdot 10^{11}$, which is far greater than our estimate of 100 for $\|h_1\|$ when $n = 10,000$.

It is important to keep in mind that we have analyzed the diagonalized problem whose matrix is $\text{diag}(1, B)$. As we pointed out in Section 2, the norms of the errors in A must be multiplied by the condition number of the diagonalizing transformation. In particular, ill-conditioning in λ_1 will limit the accuracy of the final solution.

Although we have naturally focused on errors whose norms are bounded away from zero, we can use our analysis to show that if $\epsilon_k = \|E_k\|$ converges monotonically to zero and ϵ_1 is suitably small, then the power method converges. Specifically, we have the following theorem.

Theorem 5.2. *In the above notation, let $0 < \beta < 1$. For any η_1 , there is an ϵ_1 such that if the sequence $\epsilon_1, \epsilon_2, \dots$ approaches zero monotonically then the sequence defined by*

$$\eta_{k+1} = \varphi_{\epsilon_k}(\eta_k), \quad k = 1, 2, \dots,$$

converges monotonically to zero.

Proof. From (5.4) it is clear that if ϵ_1 is sufficiently small then $\varphi'_\epsilon(\eta) \leq \alpha < 1$ for any $\epsilon < \epsilon_1$ and $\eta < \eta_1$. It then follows from the theory of fixed point iterations that the sequence η_1, η_2, \dots is monotonic decreasing. Let its limit be $\hat{\eta}$.

We must show that $\hat{\eta} = 0$. Let $\delta > 0$ be given. Now $\lim_{\epsilon \rightarrow 0} \varphi_\epsilon(\eta) = \beta\eta$ uniformly on $[0, \eta_1]$. Hence there is an integer $K > 0$ such that

$$k \geq K \implies |\varphi_{\epsilon_k}(\eta_k) - \beta\eta_k| < \frac{\delta}{2}.$$

We may also assume that K is so large that

$$k \geq K \implies |\beta\eta_k - \beta\hat{\eta}| < \frac{\delta}{2}.$$

Then for $k \geq K$

$$|\eta_{k+1} - \beta\hat{\eta}| = |\varphi_{\epsilon_k}(\eta_k) - \beta\hat{\eta}| \leq |\varphi_{\epsilon_k}(\eta_k) - \beta\eta_k| + |\beta\eta_k - \beta\hat{\eta}| < \delta.$$

It follows that $\eta_k \rightarrow \beta\hat{\eta}$. But since $\eta_k \rightarrow \hat{\eta}$ and $\beta \neq 0$, we must have $\hat{\eta} = 0$. ■

This theorem has an important implication for the behavior of the perturbed powers P_k , which was treated in the previous section. The j th column of P_k , suitably scaled, is just the result of applying the unscaled power method with error to \mathbf{e}_j . Now suppose that $y^H \mathbf{e}_j \neq 0$, where y is the dominant left eigenvector. Then if $\epsilon_1 \geq \epsilon_2 \geq \dots$ and ϵ_1 is sufficiently small, the j th column of P_k , suitably scaled, approximates the dominant eigenvector of A , even if P_k converges to zero. Thus if we are interested only in the behavior of the columns of P_k , we can relax the condition that $\sum_k \epsilon_k < \infty$. However, the price we pay is a less clean estimate of the asymptotic convergence rate.

Acknowledgements

I would like to thank Donald Estep and Sean Eastman for their comments on this paper, and especially Sean Eastman for the elegant proof of Theorem 5.2. I am indebted to the Mathematical and Computational Sciences Division of the National Institute of Standards and Technology for the use of their research facilities.

References

- [1] L. V. Ahlfors. *Complex Analysis*. McGraw–Hill, New York, 1966.
- [2] G. Dahlquist and Å. Björck. *Numerical Methods*. Prentice–Hall, Englewood Cliffs, New Jersey, 1974.
- [3] W. Gautschi. The asymptotic behavior of powers of a matrix. *Duke Mathematical Journal*, 20:127–140, 1953.
- [4] W. Gautschi. The asymptotic behavior of powers of a matrix. II. *Duke Mathematical Journal*, 20:275–279, 1953.
- [5] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, MD, third edition, 1996.
- [6] P. Henrici. Bounds for iterates, inverses, spectral variation and fields of values of nonnormal matrices. *Numerische Mathematik*, 4:24–39, 1962.
- [7] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, 1996.
- [8] N. J. Higham and P. A. Knight. Matrix powers in finite precision arithmetic. *SIAM Journal on Matrix Analysis and Applications*, 16:343–358, 1995.

- [9] A. M. Ostrowski. *Solution of Equations and Systems of Equations*. Academic Press, New York, second edition, 1966.
- [10] G. W. Stewart. *Matrix Algorithms I: Basic Decompositions*. SIAM, Philadelphia, 1998.
- [11] G. W. Stewart. *Matrix Algorithms II: Eigensystems*. SIAM, Philadelphia, 2001.