

# Computing $A^\alpha$ , $\log(A)$ and related matrix functions by contour integrals

Nicholas Hale

*Oxford University Computing Laboratory*

Nicholas J. Higham

*School of Mathematics, The University of Manchester*

*(This work was supported by a Royal Society-Wolfson Research Merit Award.)*

Lloyd N. Trefethen

*Oxford University Computing Laboratory*

New methods are proposed for the numerical evaluation of  $f(\mathbf{A})$  or  $f(\mathbf{A})b$ , where  $f(\mathbf{A})$  is a function such as  $\mathbf{A}^{1/2}$  or  $\log(\mathbf{A})$  with singularities in  $(-\infty, 0]$  and  $\mathbf{A}$  is a matrix with eigenvalues on or near  $(0, \infty)$ . The methods are based on combining contour integrals evaluated by the periodic trapezoid rule with conformal maps involving Jacobi elliptic functions. The convergence is geometric, so that the computation of  $f(\mathbf{A})b$  is typically reduced to one or two dozen linear system solves, which can be carried out in parallel.

*Subject classifications:* AMS(MOS): 65F30, 65D30

*Key words and phrases:* Cauchy integral, conformal map, contour integral, matrix function, quadrature, rational approximation, trapezoid rule

Oxford University Computing Laboratory  
Numerical Analysis Group  
Wolfson Building  
Parks Road  
Oxford, England OX1 3QD

August, 2007

# 1 Introduction

It is well known that an analytic function  $f$  of a square matrix  $\mathbf{A}$  can be represented as a contour integral,

$$f(\mathbf{A}) = \frac{1}{2\pi i} \int_{\Gamma} f(z) (z\mathbf{I} - \mathbf{A})^{-1} dz, \quad (1.1)$$

where  $\Gamma$  is a closed contour lying in the region of analyticity of  $f$  and winding once around the spectrum  $\sigma(\mathbf{A})$  in the counterclockwise direction. However, this idea has not often been exploited for numerical computation. Here we propose efficient ways to make use of (1.1) in the case where  $\mathbf{A}$  has eigenvalues on or near the positive real axis  $(0, \infty)$  and  $f(z)$  is a function such as  $z^\alpha$  or  $\log z$  that is analytic apart from singularities or a branch cut on or near the negative real axis  $(-\infty, 0]$ . To be precise, it is not the matrix  $f(\mathbf{A})$  that we usually compute but the vector  $f(\mathbf{A})b$  for a given vector  $b$ . The first method we propose is to transplant the problem by a conformal map to an annulus and then apply the trapezoid rule, which converges geometrically (§2). This procedure will be especially effective in cases where  $\mathbf{A}$  is what might be called a “backslash matrix”: a large sparse matrix for which systems of equations  $(z\mathbf{I} - \mathbf{A})x = b$  can be solved efficiently by sparse direct methods but Krylov subspace iterations and Schur reduction to triangular form are impractical. An example is the matrix corresponding to the standard 5-point finite difference discretization of the Laplacian in two dimensions.

We then consider two modifications to this technique to take advantage of particular structure of  $f$ : first, if  $f$  has a branch cut but no singularities on  $(-\infty, 0)$  (§3); second, if  $f(z) = z^{1/2}$  (§4). (Here and throughout this paper, the notations  $z^{1/2}$ ,  $z^\alpha$ , and  $\log z$  refer to the standard branches of these functions, and  $\mathbf{A}^{1/2}$  and so on to the corresponding principal branches of the matrix functions.) Significant improvements can be achieved in such cases, including the avoidance of complex arithmetic.

All our methods based on quadrature formulas can be interpreted as making implicit use of rational approximations  $f(z) \approx r(z)$  in regions of the complex plane. In §5 we examine how these implicitly constructed rational functions compare with those obtained by solving a rational approximation problem directly. For the case  $f(z) = z^{1/2}$ , we find that the combination of our conformal map with the trapezoid rule reproduces a best rational approximation discovered by Zolotarev in 1877.

In §6 we consider the effect of complex eigenvalues on the convergence of our methods, and §7 examines three further numerical examples.

Before beginning, we shall slightly change the starting point. We have found that for the problems considered here a better way to obtain  $f(\mathbf{A})$  is often by computing  $\mathbf{A} \cdot \mathbf{A}^{-1}f(\mathbf{A})$ , replacing (1.1) with the formula

$$f(\mathbf{A}) = \frac{\mathbf{A}}{2\pi i} \int_{\Gamma} z^{-1} f(z) (z\mathbf{I} - \mathbf{A})^{-1} dz. \quad (1.2)$$

The reason is that the  $dz$  factor becomes large as  $\Gamma$  swings out into the complex plane, and the factor  $z^{-1}$  counters this effect. For the first two methods derived in §§ 2 and 3, the use of (1.2) instead of (1.1) typically improves accuracy by a digit or two. For the method of §4, its use is essential since the contour passes through  $z = \infty$ .

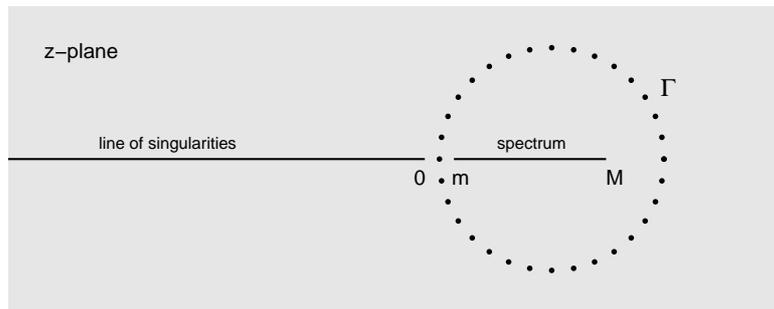


Figure 1: Depiction of the trapezoid rule applied directly to (1.2), with the contour  $\Gamma$  taken as a circle. This method is highly inefficient if  $M \gg m$ .

For general information about the theory and computation of functions of matrices, see [11].

## 2 Method 1: conformal map from an annulus

Let  $\mathbf{A}$  be a real matrix whose eigenvalues lie in the interval  $(0, \infty)$ . In many applications  $\mathbf{A}$  will be symmetric, but we do not require this. The assumption that  $\mathbf{A}$  is real will save a factor of two by allowing trapezoid sums to exploit symmetry. The assumption that the eigenvalues are real can be relaxed; our methods do not degrade very much if the eigenvalues move slightly off the real axis, as is illustrated in §6.

Let  $m$  and  $M$  be the minimum and maximum eigenvalues of  $\mathbf{A}$ , respectively. If  $\mathbf{A}$  is symmetric, or more generally normal, then  $M/m$  is its 2-norm condition number, whereas for general  $\mathbf{A}$  the condition number may be larger. We assume at the outset that  $m$  and  $M$  are known, although in practice, it would often be necessary first to estimate them. We also suppose that  $\sigma(\mathbf{A})$  more or less fills  $[m, M]$  in the sense that we do not attempt to exploit any gaps in the spectrum.

The assumption we make initially about  $f$  is that it is analytic in the slit complex plane  $\mathbb{C} \setminus (-\infty, 0]$ . This is true, for example, if  $f(z) = \log z$  or  $f(z) = z^\alpha = e^{\alpha \log z}$  (for any  $\alpha \in \mathbb{R}$ ) for the standard branches. It is also true of more complicated functions such as  $f(z) = \Gamma(z)$  or  $f(z) = \tanh(z^{1/2})$ .

One numerical approach to (1.2) is to surround  $[m, M]$  by a circle in the right half-plane and apply the trapezoid rule on this circle to approximate the integral, as sketched in Figure 1. However, this method will be terribly inefficient if  $\mathbf{A}$  is ill-conditioned, requiring  $O(M/m)$  or more linear system solves to get any accuracy at all [4], because  $\Gamma$  is contained in only a very narrow annulus of analyticity. Our strategy will be to improve the problem by a change of variables in the complex plane, which will bring the count down to  $O(\log(M/m))$ . The idea of combining the trapezoid rule with a change of variables was put forward for real integrals by Sag and Szekeres in 1964 [20] and Schwartz in 1969 [22] and developed more fully in the 1970s by Iri, Moriguti and Takasawa [12], Takahashi and Mori [17, 28, 29], and Stenger [25, 26, 27], among others. For inverse Laplace

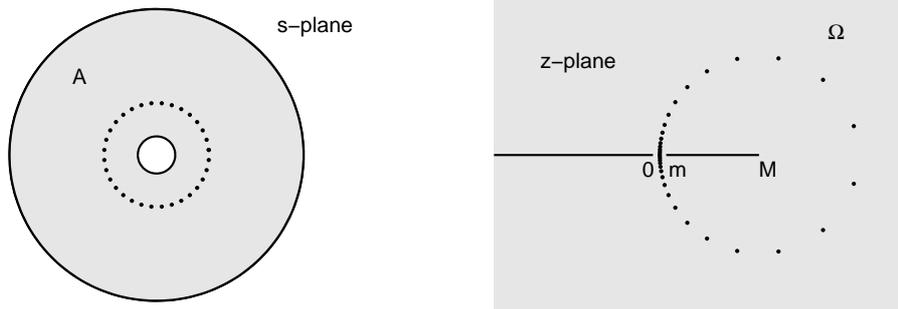


Figure 2: Conformal map for Method 1. First, transplant the entire doubly connected region of analyticity  $\Omega$  conformally to an annulus  $\mathcal{A}$ , with the line of singularities mapping to the outer boundary circle and the interval containing the spectrum to the inner one. Then apply the trapezoid rule over a circle in the annulus. This will correspond to a nonuniform distribution of quadrature points in the  $z$ -plane. In this figure  $M/m = 8$  and the annulus is quite thick. As  $M/m$  increases, it becomes thinner only logarithmically. The dots show quadrature nodes for  $N = 16$ .

transform contour integrals, such ideas were developed by Talbot in 1979 [30] and have been exploited by a variety of authors [10, 16, 23, 31], and the parameters have recently been optimized in various senses by Weideman and Trefethen [33, 34, 36, 37]. However, we are not aware of previous papers on such methods for the matrix function problems considered here.

Our aim is to transplant the region of analyticity of  $f$  and  $(z\mathbf{I} - \mathbf{A})^{-1}$  conformally to an annulus

$$\mathcal{A} = \{z \in \mathbb{C} : r < |z| < R\}.$$

According to results going back in part to Poisson in the 1820s and first fully worked out by Davis in the 1950s [5], [6, §4.6.5], the trapezoid rule applied over a circle within  $\mathcal{A}$  will converge geometrically as  $N \rightarrow \infty$ , where  $N$  is the number of sample points. Moreover, this will be a thicker annulus than was available before the conformal map (Figure 1), which will make the convergence constant quite favorable.

This strategy is sketched in Figure 2. The region in question is the doubly connected set

$$\Omega = \mathbb{C} \setminus ((-\infty, 0] \cup [m, M]).$$

We thus face a conformal mapping problem: how to map this domain  $\Omega$  onto an annulus  $\mathcal{A}$ ? More precisely it is the map from  $\mathcal{A}$  to  $\Omega$  that we shall need to compute. The radii  $R$  and  $r$  are not entirely at our disposal: the formulas below imply that  $R/r$  is determined by  $M/m$ .

The map can be carried out in three steps, as shown in Figure 3. First, the function

$$t = \frac{2Ki}{\pi} \log(-is/r)$$

carries the upper half of  $\mathcal{A}$  to the rectangle with vertices  $\pm K$  and  $\pm K + iK'$ . The numbers  $K$  and  $K'$  are complete elliptic integrals whose values are fixed by the next

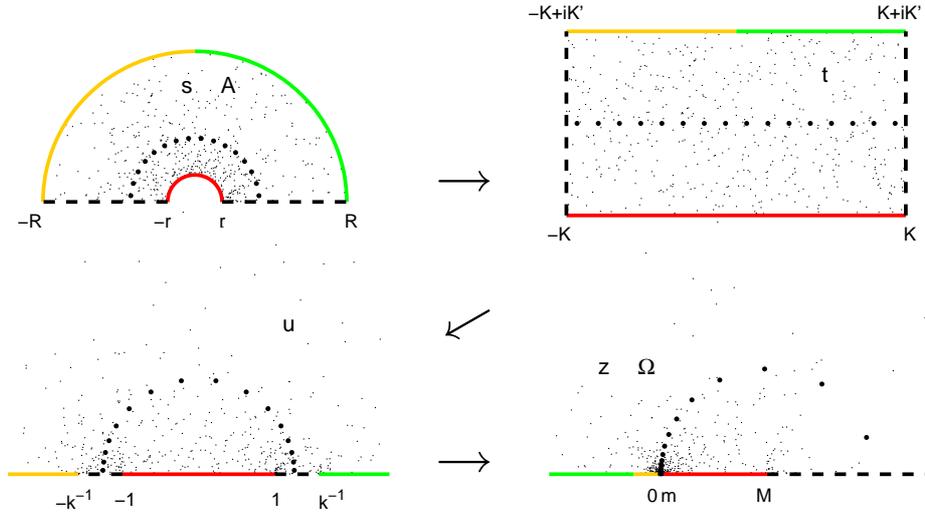


Figure 3: The map of Figure 2, constructed in three steps  $s \rightarrow t \rightarrow u \rightarrow z$ . Here  $m = 1$ ,  $M = 20$ . The boundaries are shown in various colors and line styles together with 500 random interior points in the rectangle and a typical set of  $N = 16$  quadrature points for the trapezoid rule in the  $t$ -plane. Notice how the narrow gap between 0 and  $m$  in the  $z$ -plane broadens to a comfortable ratio  $R/r \approx 5.8$  in the annulus.

step of the conformal map. In this second step the Jacobi elliptic function

$$u = \operatorname{sn}(t) = \operatorname{sn}(t|k^2), \quad k = \frac{\sqrt{M/m} - 1}{\sqrt{M/m} + 1} \quad (2.1)$$

maps the rectangle to the upper half-plane, with the ends mapping to  $[-k^{-1}, -1]$  and  $[1, k^{-1}]$  (see [1, 9]). Finally, the Möbius transformation

$$z = \sqrt{mM} \left( \frac{k^{-1} + u}{k^{-1} - u} \right) \quad (2.2)$$

carries the upper half-plane to itself in such a way that  $[-k^{-1}, -1]$  and  $[1, k^{-1}]$  are sent to  $[0, m]$  and  $[M, \infty]$ , respectively. We have now mapped the upper half of  $\mathcal{A}$  to the upper half of  $\Omega$ . Reflection across the interval  $(r, R)$  (using the Schwarz reflection principle) extends this map to all of  $\mathcal{A}$  and all of  $\Omega$ .

Conceptually, this map from  $\mathcal{A}$  to  $\Omega$  shows the essence of our method most clearly. To apply the trapezoid rule and derive corresponding theorems, however, we can bypass  $\mathcal{A}$  and work with the map from the  $K, K'$  rectangle in the  $t$ -plane to  $\Omega$ . By composing (2.2) and (2.1) we get

$$z = \sqrt{mM} \left( \frac{k^{-1} + \operatorname{sn}(t)}{k^{-1} - \operatorname{sn}(t)} \right), \quad k = \frac{\sqrt{M/m} - 1}{\sqrt{M/m} + 1}. \quad (2.3)$$

The integral (1.2) can accordingly be written

$$f(\mathbf{A}) = \frac{-\mathbf{A}}{2\pi i} \int_{-K+iK'/2}^{3K+iK'/2} z^{-1} f(z(t)) (z(t) - \mathbf{A})^{-1} \frac{dz}{du} \frac{du}{dt} dt;$$

the limits  $-K + iK'/2$  to  $K + iK'/2$  would correspond to the part of  $\Gamma$  in the upper half-plane (the dots in the fourth panel of Figure 3), and extending this to  $3K + iK'/2$  brings in the part in the lower half-plane. With the aid of the identities

$$\frac{dz}{du} = \frac{2k^{-1}\sqrt{mM}}{(k^{-1} - u)^2}, \quad \frac{du}{dt} = \operatorname{sn}'(t) = \sqrt{1 - k^2 u^2} \sqrt{1 - u^2} = \operatorname{cn}(t) \operatorname{dn}(t),$$

where  $\operatorname{cn}$  and  $\operatorname{dn}$  are further Jacobi elliptic functions in standard notation [1], this becomes

$$f(\mathbf{A}) = \frac{-\mathbf{A}\sqrt{mM}}{\pi i k} \int_{-K+iK'/2}^{3K+iK'/2} \frac{z^{-1} f(z(t)) (z(t) - \mathbf{A})^{-1} \operatorname{cn}(t) \operatorname{dn}(t)}{(k^{-1} - u)^2} dt. \quad (2.4)$$

Now since  $\mathbf{A}$  is real, the integrand is real-symmetric, i.e., the values it takes at two points  $z$  and  $\bar{z}$  are complex conjugate. It follows that  $f(\mathbf{A})$  is twice the real part of the value obtained by integrating over the first half of the contour, or if we cancel the  $i$  in the denominator,

$$f(\mathbf{A}) = \frac{-2\mathbf{A}\sqrt{mM}}{\pi k} \operatorname{Im} \int_{-K+iK'/2}^{K+iK'/2} \frac{z^{-1} f(z(t)) (z(t) - \mathbf{A})^{-1} \operatorname{cn}(t) \operatorname{dn}(t)}{(k^{-1} - u)^2} dt, \quad (2.5)$$

with  $z(t)$  given by (2.3).

We now apply the trapezoid rule with  $N$  equally spaced points in  $(-K + iK'/2, K + iK'/2)$ ,

$$t_j = -K + \frac{iK'}{2} + 2\frac{(j - \frac{1}{2})K}{N}, \quad 1 \leq j \leq N. \quad (2.6)$$

(Actually this choice should more precisely be called a midpoint rule. One could make it a true trapezoid rule by shifting the sample points to include the endpoints with weights  $1/2$ .) The result becomes  $f(\mathbf{A}) \approx f_N(\mathbf{A})$  with

$$f_N(\mathbf{A}) = \frac{-4K\mathbf{A}\sqrt{mM}}{\pi N k} \operatorname{Im} \sum_{j=1}^N \frac{f(z(t_j)) (z(t_j) \mathbf{I} - \mathbf{A})^{-1} \operatorname{cn}(t_j) \operatorname{dn}(t_j)}{z(t_j) (k^{-1} - u(t_j))^2}. \quad (2.7)$$

(If  $A$  is not real, extend the limits in (2.6) and (2.7) from  $N$  to  $2N$ , and in (2.7), replace  $4K$  by  $2K$  and multiply by  $-i$  instead of extracting the imaginary part.)

The convergence of the method we have described is geometric, with the constant of convergence worsening only logarithmically as  $M/m \rightarrow \infty$ .

**Theorem 1** *Let  $\mathbf{A}$  be a real matrix with eigenvalues in  $[m, M]$  and let  $f$  be a function analytic in  $\mathbb{C} \setminus (-\infty, 0]$ . Then the  $N$ -point conformally transplanted quadrature formula (2.7) converges at the rate*

$$\|f(\mathbf{A}) - f_N(\mathbf{A})\| = O(e^{\varepsilon - \pi K' N / (2K)}) \quad (2.8)$$

for any  $\varepsilon > 0$  as  $N \rightarrow \infty$ , and the constant in the exponent is asymptotically  $\pi K'/(2K) \sim \pi^2/\log(M/m)$  as  $M/m \rightarrow \infty$ . For any  $M$  and  $m$  we have

$$\|f(\mathbf{A}) - f_N(\mathbf{A})\| = O(e^{-\pi^2 N/(\log(M/m)+3)}). \quad (2.9)$$

**Proof** The formula (2.7) for  $f_N(\mathbf{A})$  comes from applying the trapezoid rule to a function in the  $t$ -plane that is analytic in a strip of half-width  $a = K'/2$ , with grid spacing  $\Delta t = 2K/N$ . Standard results on the trapezoid rule for periodic integrands imply convergence in such a context at the rate  $O(e^{\varepsilon-2\pi a/\Delta t})$  for any  $\varepsilon > 0$  [5, 6], and combining these numbers gives (2.8). The claim about asymptotics follows from the relationship  $K'/K \sim \pi/(\log(16) - \log(1-k^2)) \sim -\pi/\log(1-k)$  as  $k \rightarrow 1$  for Jacobi elliptic functions [1, eq. (17.3.26)] together with the relationship  $\log(M/m) \sim -2\log(1-k)$  implied by (2.1). The bound  $\pi^2/(\log(M/m) + 3) < (\pi/2)(K'/K)$  needed to derive (2.9) from (2.8) was established numerically. ■

We can illustrate Method 1 by the following MATLAB script. Standard MATLAB includes functions to evaluate  $\text{sn}(z)$  for real arguments only. This program needs complex arguments, and for this purpose it calls the functions `ellipkcp` and `ellipjc` from Driscoll's Schwarz–Christoffel Toolbox [7], [8], which is freely available online. The mathematical basis of these codes is outlined in [9].

This test code `method1`, like `method2` and `method3` to follow, computes the whole matrix function  $f(\mathbf{A})$  rather than just a vector  $f(\mathbf{A})b$ . Our purpose is to illustrate convergence rates as a function of the number of sample points  $N$ , which we can do just as well for a small matrix that can be fully inverted. The matrix illustrated here is the  $5 \times 5$  Pascal triangle matrix `pascal(5)`, with  $M/m \approx 10^4$ . The code begins by computing the minimal and maximal eigenvalues  $m$  and  $M$ , although in practice, some kind of estimation would normally be used.

```
% method1.m - evaluate f(A) by contour integral. The functions
%           ellipkcp and ellipjc are from Driscoll's SC Toolbox.

f = @sqrt;                % change this for another function f
A = pascal(5);            % change this for another matrix A
X = sqrtm(A);            % change this if f is not sqrt
I = eye(size(A));
e = eig(A); m = min(e); M = max(e); % use only for toy problems!
k = (sqrt(M/m)-1)/(sqrt(M/m)+1);
L = -log(k)/pi;
[K,Kp] = ellipkcp(L);
for N = 5:5:40
    t = .5i*Kp - K + (.5:N)*2*K/N;
    [u cn dn] = ellipjc(t,L);
    z = sqrt(m*M)*((1/k+u)./(1/k-u));
    dzdt = cn.*dn./(1/k-u).^2;
    S = zeros(size(A));
```

```

for j = 1:N
    S = S + (f(z(j))/z(j))*inv(z(j)*I-A)*dzdt(j);
end
S = -4*K*sqrt(m*M)*imag(S)*A/(k*pi*N);
error = norm(S-X)/norm(X);
fprintf('%4d %10.2e\n', N, error)
end

```

Here are the relative errors  $\|f(\mathbf{A}) - f_N(\mathbf{A})\|/\|f(\mathbf{A})\|$  printed by this program ( $\|\cdot\|$  is the 2-norm), showing 15-digit precision for  $N \approx 40$ . If  $M/m$  were only  $10^2$  (assuming  $\mathbf{A}$  is symmetric or more generally normal),  $N \approx 20$  would be enough to achieve the same accuracy. If  $z^{1/2}$  is replaced by  $\log z$ , the performance is about the same.

```

>> method1
    5   3.03e-02
   10   4.74e-04
   15   7.29e-06
   20   1.12e-07
   25   1.73e-09
   30   2.66e-11
   35   4.11e-13
   40   7.07e-15

```

In a sequence of computations like this, one could reuse certain data points. For example, the resolvents evaluated for  $N = 10$  are evaluated again for  $N = 30$ . We have not attempted to exploit this redundancy, and if one were doing so, it would be advantageous to switch from the midpoint to the true trapezoid formulation, as discussed above, to make the overlaps occur at multiples of 2 rather than 3.

We conclude this section with an example involving a more challenging function  $f$ , with singularities all along  $(-\infty, 0]$ , and a nonnormal matrix with eigenvalues at  $(3 \pm \sqrt{5})/2$ :

$$\mathbf{A} = \begin{pmatrix} 1 & \frac{1}{2} \\ 2 & 2 \end{pmatrix}, \quad f(z) = \Gamma(z).$$

The gamma function has poles at all the nonpositive integers, but that is no problem for Method 1. Taking  $N = 42$  is enough to give a result accurate to ten digits,

$$\Gamma(\mathbf{A}) = \begin{pmatrix} 2.0835578979 & -0.1960182234 \\ -0.7840728935 & 1.6915214512 \end{pmatrix},$$

as is readily checked by a diagonalization of  $\mathbf{A}$ . Of course the point of our method is that it is applicable also in cases where the dimension of  $\mathbf{A}$  is such that diagonalization is impractical, or indeed, if  $\mathbf{A}$  is not diagonalizable at all. Another approach to the computation of  $\Gamma(\mathbf{A})$  based on a Hankel integral has been discussed in [21].

### 3 Method 2: when $(-\infty, 0)$ is just a branch cut

The assumption of the last section was that  $f$  might have arbitrary singularities on  $(-\infty, 0]$ . In practice, many functions of interest, including  $z^\alpha$  and  $\log z$ , have a singularity at  $z = 0$  but just a branch cut on  $(-\infty, 0)$ . This means that  $f(z)$  can be analytically continued along any path that avoids  $z = 0$ ; the difficulty for  $z \neq 0$  is only that the value of  $f(z)$  will change if the path winds around the origin. Thus we may regard  $f$  either as a single-valued function on the slit plane  $\mathbb{C} \setminus (-\infty, 0]$ , or as a multivalued function on a Riemann surface associated with the punctured plane  $\mathbb{C} \setminus \{0\}$ .

For functions of this kind, Method 1 is not as efficient as it might be, since it employs a contour that avoids  $[m, M]$  and all of  $(-\infty, 0]$  rather than just  $[m, M]$  and  $\{0\}$ . One way to improve matters is sketched in Figure 4. If we introduce a new variable  $w = z^{1/2}$ ,  $dz = 2w dw$ , then (1.2) becomes

$$f(\mathbf{A}) = \frac{\mathbf{A}}{\pi i} \int_{\Gamma_w} w^{-1} f(w^2) (w^2 - \mathbf{A})^{-1} dw. \quad (3.1)$$

Under the square root, the branch cut of  $f$  along  $(-\infty, 0]$  in the  $z$  plane unfolds to the imaginary axis in the  $w$ -plane. The assumption on  $f(z)$  then implies that  $f(w^2)$  can be analytically continued to an analytic function throughout the slit  $w$ -plane  $\mathbb{C} \setminus (-\infty, 0]$ . Thus we have a new contour integral problem, where the contour should lie in the slit plane and enclose  $[m^{1/2}, M^{1/2}]$ . The geometry is the same as in the last section, but with  $[m, M]$  improved to  $[m^{1/2}, M^{1/2}]$ . Accordingly we may use the same method as before. The formula for the new method applied to a real matrix  $\mathbf{A}$  is

$$f_N(\mathbf{A}) = \frac{-8K\mathbf{A}(mM)^{1/4}}{\pi Nk} \operatorname{Im} \sum_{j=1}^N \frac{f(w(t_j)^2) (w(t_j)^2 \mathbf{I} - \mathbf{A})^{-1} \operatorname{cn}(t_j) \operatorname{dn}(t_j)}{w(t_j) (k^{-1} - u(t_j))^2}, \quad (3.2)$$

where the  $t_j$  are again given by (2.6) but now (2.3) is modified to

$$w = (mM)^{1/4} \left( \frac{k^{-1} + \operatorname{sn}(t)}{k^{-1} - \operatorname{sn}(t)} \right), \quad k = \frac{(M/m)^{1/4} - 1}{(M/m)^{1/4} + 1} \quad (3.3)$$

and the values of  $K$  and  $K'$  are the complete elliptic integrals associated with this new parameter  $k$ . For a complex matrix  $\mathbf{A}$  the modifications are as for Method 1.

Here is a MATLAB code.

```
% method2.m - Variant of method1.m for evaluating f(A) by contour
%                integration, assuming only a branch cut along (-inf,0).

f = @sqrt;           % change this for another function f
A = pascal(5);      % change this for another matrix A
X = sqrtm(A);       % change this if f is not sqrt
I = eye(size(A));
```

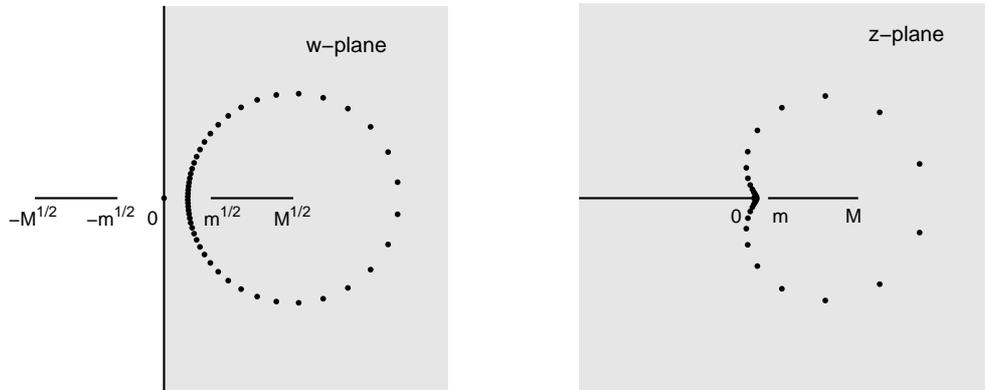


Figure 4: In Method 2, for problems where  $f(z)$  has only a branch cut on  $(-\infty, 0)$ , we apply the same quadrature points seen in the fourth panel of Figure 3 to the function  $f(w^2)$  in the  $w$ -plane (left). This square-roots the condition number and approximately doubles the rate of convergence. In the  $z$ -plane, the new method corresponds to a quadrature rule with the nodes shown on the right.

```

e = eig(A); m = min(e); M = max(e);    % only for toy problems
k = ((M/m)^(1/4)-1)/((M/m)^(1/4)+1);
L = -log(k)/pi;
[K,Kp] = ellipkcp(L);
for N = 5:5:25
    t = .5i*Kp - K + (.5:N)*2*K/N;
    [u cn dn] = ellipjc(t,L);
    w = (m*M)^(1/4)*((1/k+u)./(1/k-u));
    dzdt = cn.*dn./(1/k-u).^2;
    S = zeros(size(A));
    for j = 1:N
        S = S + (f(w(j)^2)/w(j))*inv(w(j)^2*I-A)*dzdt(j);
    end
    S = -8*K*(m*M)^(1/4)*imag(S)*A/(k*pi*N);
    error = norm(S-X)/norm(X);
    fprintf('%4d %10.2e\n', N, error)
end

```

The new code converges a good deal faster than the previous one, as expected, getting close to full precision for  $N \approx 25$  for this matrix with  $M/m \approx 10^4$ . Again the results are similar if  $z^{1/2}$  is replaced by  $\log z$ .

```

>> method2
5    2.97e-03

```

10	5.51e-07
15	7.03e-10
20	4.88e-12
25	7.29e-15

The same arguments as before give us a convergence theorem. Asymptotically as  $M/m \rightarrow \infty$ , Method 2 is twice as fast as Method 1, though the improvement is less than a factor of 2 for finite  $M/m$ .

**Theorem 2** *Let  $\mathbf{A}$  be a real matrix with eigenvalues in  $[m, M]$  and let  $f$  be a function analytic in  $\mathbb{C} \setminus (-\infty, 0]$  that can be continued analytically across  $(-\infty, 0)$  from the upper half-plane to the lower half-plane. Then the  $N$ -point improved formula (3.2) converges at the rate*

$$\|f(\mathbf{A}) - f_N(\mathbf{A})\| = O(e^{\varepsilon - \pi K' N / (2K)}) \quad (3.4)$$

for any  $\varepsilon > 0$ , where  $K$  and  $K'$  are the complete elliptic integrals associated with the parameter  $k$  of (3.3), and the constant in the exponent is asymptotically  $\pi K' / (2K) \sim 2\pi^2 / \log(M/m)$  as  $M/m \rightarrow \infty$ . For any  $M$  and  $m$  we have

$$\|f(\mathbf{A}) - f_N(\mathbf{A})\| = O(e^{-2\pi^2 N / (\log(M/m) + 6)}). \quad (3.5)$$

## 4 Method 3: special treatment for $f(z) = z^{1/2}$

The function  $f(w^2)$  in the  $w$ -plane in the last section, as sketched in Figure 4, has a singularity in general at the origin, and that is why the quadrature points must pass between 0 and  $m^{1/2}$ . For the particular case  $f(z) = z^{1/2} = w$ , however, there is no singularity, because the change of variables from  $z$  to  $w$  has eliminated it. Equation (3.1) becomes simply

$$\mathbf{A}^{1/2} = \frac{\mathbf{A}}{\pi i} \int_{\Gamma_w} (w^2 - \mathbf{A})^{-1} dw, \quad (4.1)$$

where  $\Gamma_w$  is a closed contour surrounding  $[m^{1/2}, M^{1/2}]$  but not  $[-M^{1/2}, -m^{1/2}]$ . This makes possible some further improvements. First of all, the quadrature points can now safely approach or even pass through 0. Secondly, we can use the symmetry to put them on the imaginary axis, corresponding to the negative real axis in the  $z$ -plane, thereby eliminating the complex arithmetic in the linear algebra problems. Indeed, since the integrand is of size  $O(w^{-2})$ , (4.1) is equivalent to an integral over the imaginary axis alone,

$$\mathbf{A}^{1/2} = \frac{i\mathbf{A}}{\pi} \int_{-i\infty}^{i\infty} (w^2 - \mathbf{A})^{-1} dw. \quad (4.2)$$

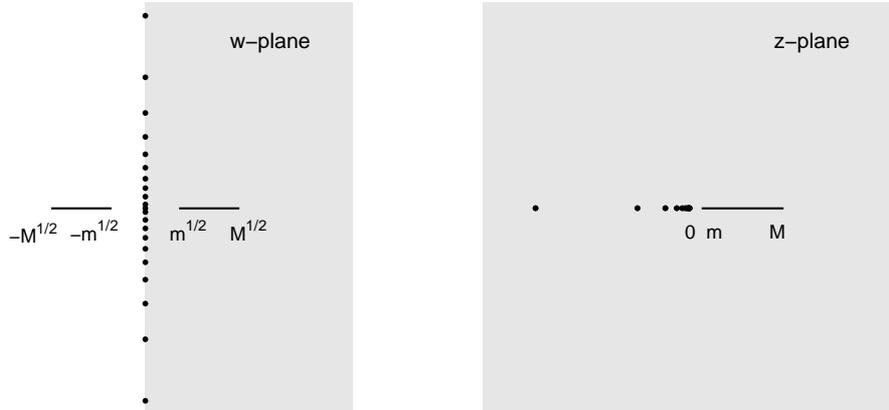


Figure 5: Quadrature points for Method 3. With  $f(z) = z^{1/2}$ , we can now fully exploit the symmetry in the  $w$ -plane. The points in the  $z$ -plane lie in  $(-\infty, 0]$ , so all the matrix arithmetic is real.

(Essentially the same formula appears as equation (6.1) in [11].) Thirdly, we can save a factor of two in discretizing this integral since  $(w^2 - \mathbf{A})^{-1} = ((-w)^2 - \mathbf{A})^{-1}$ , regardless of the properties of  $\mathbf{A}$ , so it will no longer matter whether or not  $\mathbf{A}$  is real.

Figure 5 shows how Figure 4 changes in this special case. Figure 6 shows the conformal map we can use to compute this transformation, which is essentially the middle step of Figure 3, in which a rectangle is mapped to the upper half-plane. The rectangle in the  $t$ -plane is the same as before, and the function  $w = m^{1/2} \operatorname{sn}(t | k^2)$  with

$$k = m^{1/2}/M^{1/2} \quad (4.3)$$

maps it onto the upper half  $w$ -plane. This time, however, it is the vertical midline of the rectangle, the imaginary interval  $[0, iK']$ , where the trapezoid rule is applied. This gives us the approximation

$$f_N(\mathbf{A}) = \frac{-2K'm^{1/2}\mathbf{A}}{\pi N} \sum_{j=1}^N (w(t_j)^2 \mathbf{I} - \mathbf{A})^{-1} \operatorname{cn}(t_j) \operatorname{dn}(t_j) \quad (4.4)$$

with

$$t_j = i(j - \frac{1}{2})K'/N, \quad 1 \leq j \leq N. \quad (4.5)$$

Equation (4.4) works well, but like the algorithms of the past two sections, it has the disadvantage of requiring evaluation of elliptic functions for complex arguments. In this case, however, the arguments (4.5) are purely imaginary. This allows us to take advantage of equations (16.20.1)–(16.20.3) of [1], and the code below implements (4.4) requiring only MATLAB's standard elliptic routines `ellipke` and `ellipj`. For  $M/m > 10^6$ , however, this approach can be unstable, and better numerical stability is achieved by returning to `ellipkqp` and `ellipjc` from the SC Toolbox. The code has been constructed in such a way that this alteration can be carried out by replacing the line beginning

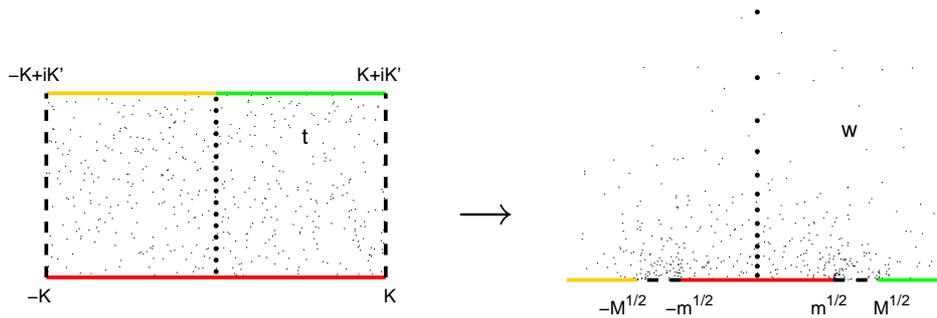


Figure 6: The conformal map for the quadrature points of Figure 5 is the scaled Jacobi sine function  $w = m^{1/2} \operatorname{sn}(t|k^2)$  with  $k = m^{1/2}/M^{1/2}$ . Compare Figure 3.

$K_p = \dots$  by the two lines  $L = -.5*\log(k^2)/\pi$ ;  $[K, K_p] = \operatorname{ellipkkp}(L)$ ; and the two lines beginning  $[\operatorname{sn} \operatorname{cn} \operatorname{dn}] = \dots$  and  $\operatorname{cn} = \dots$  by the single line  $[\operatorname{sn} \operatorname{cn} \operatorname{dn}] = \operatorname{ellipjc}(t, L)$ ; .

```
% method3.m - Variant of method1.m and method2.m for evaluating sqrt(A)

A = pascal(5); % change this for another matrix A
X = sqrtm(A); % this cannot be changed
I = eye(size(A)); e = eig(A);
m = min(e); M = max(e); % only for toy problems
k2 = m/M; % elliptic functions parameter k^2
Kp = ellipke(1-k2);
for N = 5:5:20
    t = 1i*(.5:N)*Kp/N;
    [sn cn dn] = ellipj(imag(t), 1-k2);
    cn = 1./cn; dn = dn.*cn; sn = 1i*sn.*cn;
    w = sqrt(m)*sn;
    dzdt = cn.*dn;
    S = zeros(size(A));
    for j = 1:N
        S = S - inv(A-w(j)^2*I)*dzdt(j);
    end
    S = (-2*Kp*sqrt(m)/(pi*N))*A*S;
    error(N) = norm(S-X)/norm(X);
    fprintf('%4d %10.2e\n', N, error(N))
end
```

Method 3 converges very quickly. For our matrix with condition number  $M/m \approx 10^4$ , nearly full precision is achieved with  $N \approx 20$ , and we get 6-digit precision with just 10 backslashes.

```
>> method3
      5  9.47e-04
     10  2.24e-07
     15  5.30e-11
     20  1.10e-14
```

The map involved in Method 3 is essentially the same as that of Method 1, except with  $z$  improved to  $w = z^{1/2}$ . Accordingly the same arguments as before lead to the following analogue of Theorem 1 with convergence rates twice as fast.

**Theorem 3** *Let  $\mathbf{A}$  be a real or complex matrix with eigenvalues in  $[m, M]$ . Then the formulas (4.3)–(4.5) converge to  $\mathbf{A}^{1/2}$  at the rate*

$$\|\mathbf{A}^{1/2} - f_N(\mathbf{A})\| = O(e^{\varepsilon - 2\pi KN/K'}) \quad (4.6)$$

for any  $\varepsilon > 0$  for the values of  $K$  and  $K'$  associated with the parameter (4.3), and the constant in the exponent is asymptotically  $\pi K'/(2K) \sim 2\pi^2/\log(M/m)$  as  $M/m \rightarrow \infty$ . For any  $M$  and  $m$  we have

$$\|\mathbf{A}^{1/2} - f_N(\mathbf{A})\| = O(e^{-2\pi^2 N/(\log(M/m)+3)}). \quad (4.7)$$

For  $\mathbf{A} = \text{pascal}(5)$ , for example, the theorem guarantees convergence at the rate  $O(5^{-N})$ , and for a matrix with  $M/m = 10$  this speeds up to  $O(41^{-N})$ . Even with  $M/m = 10^{10}$  the error is  $O(2.1^{-N})$ .

Figure 7 compares Methods 1, 2, and 3 for the function for which they are all applicable,  $\mathbf{A}^{1/2}$ , for two matrices with  $M/m \approx 62.0$  and  $2.06 \times 10^7$ . Method 3 converges twice as fast as Method 1 in both cases, and it does not lose a factor of 2 if  $\mathbf{A}$  is complex. (Its difficulties with rounding errors for the ill-conditioned matrix can be fixed by the use of the alternative formulation described above.) Method 2 has essentially the same convergence rate for the ill-conditioned matrix but falls between the other two methods for the well-conditioned one.

## 5 Connection with rational approximation

Since the time of Gauss it has been known that quadrature formulas are associated with rational functions: the nodes are the poles, and the weights are the residues. Thus the methods we have proposed can be interpreted as implicit descriptions of rational approximations

$$f(x) \approx r(x), \quad x \in [m, M]. \quad (5.1)$$

One might ask, how close are these approximations to optimal? Conversely, if one constructs the optimal rational approximations, do they lead to good methods for evaluating

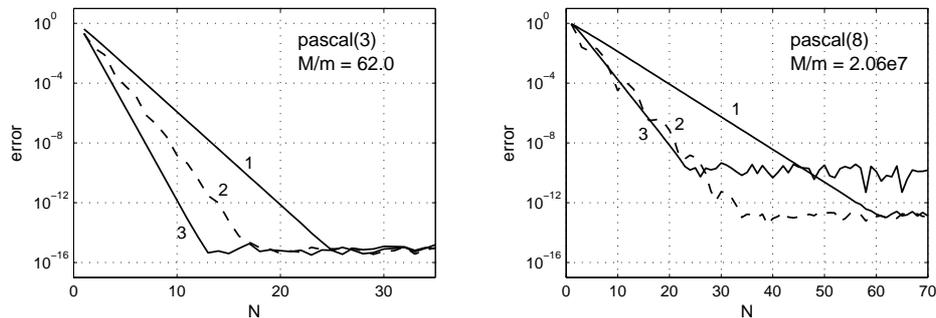


Figure 7: Convergence of Methods 1, 2, and 3 for computing  $\mathbf{A}^{1/2}$  with  $\mathbf{A} = \text{pascal}(3)$  and  $\text{pascal}(8)$ . Methods 1 and 2 are also applicable to more general functions than the square root. The loss of final accuracy in Method 3 can be fixed by the complex formulation described just above the listing of `method3.m`.

$f(\mathbf{A})$ ? The same questions are considered in the context of Talbot contours for inverse Laplace transforms in [34].

The link to rational functions goes as follows. Suppose that by some quadrature procedure we derive an approximation

$$f(\mathbf{A}) \approx \mathbf{A} \sum_{j=1}^N \gamma_j (z_j - \mathbf{A})^{-1} \quad (5.2)$$

as in (2.7), (3.2), or (4.4), where  $f$  is assumed to have spectrum in  $[m, M]$ . This approximation defines a rational function of type  $(N, N)$  (numerator and denominator of degree  $\leq N$ ),

$$r(z) = \sum_{j=1}^N \frac{\gamma_j z}{z_j - z}. \quad (5.3)$$

It is easily seen that if the difference of the two sides of (5.2) has matrix norm  $\varepsilon$ , then  $|f(\lambda) - r(\lambda)| \leq \varepsilon$  must hold for any  $\lambda$  which is an eigenvalue of  $\mathbf{A}$ . Thus (5.1) certainly holds for the eigenvalues  $\lambda$ , and presumably also for other  $x \in [m, M]$  if the approximation was derived without reference to particular eigenvalues of  $\mathbf{A}$ .

Figure 8 examines our three approximations from this point of view for the special case  $f(z) = z^{1/2}$ ,  $[m, M] = [1, 100]$ ,  $N = 6$ . The error curves  $f(x) - r(x)$  are plotted for Methods 1, 2, and 3, and the fourth panel of the figure shows the error curve for the best supremum-norm approximation, which equioscillates 14 times. (Computing the best approximation for such problems is not a trivial matter. We used the Carathéodory–Fejér method [32], as outlined in [34].) An examination of the axis scales for these curves reveals that the approximations are very good, with that of Method 3 in particular quite close to optimal.

A deeper analysis of Figure 8 points to some interesting mathematics. The methods we have put forward in this paper are general ones, combining conformal mapping and trapezoid rule quadrature to handle general functions  $f(z)$ . The third panel of the

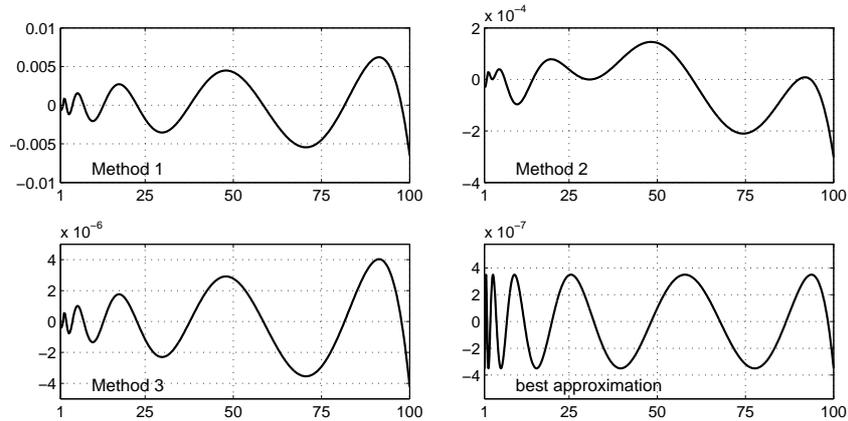


Figure 8: Errors  $z^{1/2} - r(z)$  for various type  $(N, N)$  rational approximations to  $z^{1/2}$  on  $[1, 100]$  with  $N = 6$ .

figure, however, reveals an unexpected optimality property of the Method 3 variant for  $f(z) = z^{1/2}$ . If one scales this error curve by  $z^{-1/2}$ , it is found to perfectly equioscillate. From this equioscillation one can infer that (5.3) is in this case the optimal type  $(N, N)$  approximation to  $z^{1/2}$  on  $[m, M]$  with respect to the maximum norm weighted by  $z^{-1/2}$ . Is this best rational approximation already known? Yes indeed, it was discovered by Zolotarev in 1877, who showed that several approximation problems of this kind, though not the unweighted problem of Figure 8, can be solved by the use of Jacobi elliptic functions.<sup>1</sup> The poles for these best approximations, he found, are evenly spaced with respect to the variable we have called  $t$  [2, App. E]. For more on the connection of Zolotarev's work to matrix functions, see [13, 14, 35].

In other words, the trapezoid rule has rediscovered for us one of Zolotarev's results from 130 years ago! For computing  $\mathbf{A}^{1/2}$  when  $\mathbf{A}$  has a positive real spectrum, it follows that there is probably no advantage in using the trapezoid rule rather than explicit rational approximations. Incidentally, the use of rational approximations to  $\mathbf{A}^{1/2}$  becomes particularly interesting when  $\mathbf{A}$  is singular, with eigenvalues in an interval  $[0, M]$ . It was a celebrated discovery of D. J. Newman in 1964 that whereas polynomial approximations to  $z^{1/2}$  on such an interval converge only algebraically, rational approximations converge at a rate  $O(\exp(-C\sqrt{N}))$  [18]. (Newman's paper deals with approximation of  $|z|$ , but the result can be carried over to  $z^{1/2}$  by a change of variables.) It would be interesting to see if this behavior can be recovered by a conformal map and the trapezoid rule, but we do not pursue this idea here as it is tied to the special case  $f(z) = z^{1/2}$ .

The advantage of the trapezoid rule technique is that it is general, applying to arbitrary functions  $f$  and to matrices whose spectra may not be all real, yet still readily yielding theorems about asymptotic rates of convergence. We now consider some of

<sup>1</sup>Egor Ivanovich Zolotarev was a student of Chebyshev, and he learned the theory of elliptic functions from Weierstrass during a visit to Berlin in 1872. Zolotarev made significant contributions to number theory as well as approximation theory before dying at age 31 as a result of being hit by a train at what is now called the Pushkin station outside St. Petersburg [24].

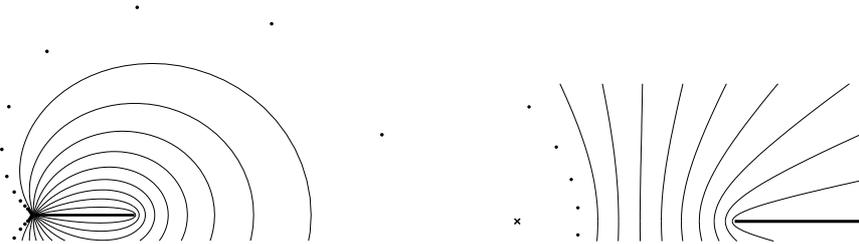


Figure 9: On the left, repetition of the right panel of Figure 4, now with curves shown inside the quadrature points. If the eigenvalues of  $\mathbf{A}$  are not confined to the real interval  $[m, M]$  (solid line) but lie within the innermost curve, for example, Method 2 will still converge exponentially with a rate 90% as fast as before. The next curve corresponds to 80%, then 70%, and so on down to 10%. On the right, a closeup of the same figure near the origin (marked by a cross).

these matters.

## 6 Complex eigenvalues and eigenvalue estimates

The methods we have described are flexible. Though we have assumed that the spectrum of  $\mathbf{A}$  lies in an interval  $[m, M]$  and that  $m$  and  $M$  are known, the convergence rate degrades only slowly as these conditions are relaxed.

Figure 9 gives an indication of the mathematical reason for this robustness. The general explanation is that in an integral like (1.1), there is no need for the spectrum to be real, provided it is enclosed by the contour  $\Gamma$ . For a quantitative analysis we note that each of our Theorems 1, 2, 3 is proved by reducing the original problem in the  $z$ -plane to a function analytic in a rectangle in the  $t$ -plane. The width of the rectangle determines the convergence rate, and if the spectrum of  $\mathbf{A}$  moves off the real axis, this narrows the rectangle in a predictable way. For example, in the case of Theorem 1, the rectangle has half-width  $K'/2$  and is shown in the second panel of Figure 3, where the bottom edge is the portion of the boundary that maps conformally to  $[m, M]$ . If the spectrum lies in a region in the  $z$ -plane larger than  $[m, M]$ , this will narrow the rectangle of analyticity. The curves in the figure accordingly correspond to images in the  $z$ -plane of horizontal lines in the  $t$ -plane at distances 90%, 80%,  $\dots$ , 10% of the original distance  $K'/2$  below the midline.

Perhaps it is instructive to focus on the middle of the nine curves in Figure 9. So long as the spectrum of  $\mathbf{A}$  lies within this quite generous region, the convergence of Method 2 will slow down by at most a factor of 2. Similarly ample regions govern the convergence of Methods 1 and 3. In the case of Method 3, it can be shown that the 50% curve is exactly the circle of radius  $M(1 - m/M)^{1/2}$  about the point  $z = M$ . For  $M \gg m$  this is approximately the circle about  $z = M$  that passes through the point  $z = m/2$ . Thus the loss if the spectrum widens from an interval to a disk, and a bigger disk at that, is at worst a factor of 2.

The arguments just made apply to our numerical methods with no adjustments of parameters, but if the spectrum of  $\mathbf{A}$  is complex, such adjustments may be advantageous. For example, in the second panel of Figure 3, if the lower half of the rectangle must be narrowed by an amount  $\alpha$  because  $\mathbf{A}$  has a complex spectrum, then one can halve the impact on the convergence rate by moving the line of quadrature points upwards a distance  $\alpha/2$ .

Here is an example of the success of these methods for matrices with complex spectra. In MATLAB, `gallery('parter',32)` constructs the  $32 \times 32$  nonsymmetric Toeplitz matrix with entries  $a_{ij} = (i - j + \frac{1}{2})^{-1}$ . The norm of  $\mathbf{A}$  matches  $\pi$  to 15 digits, and the eigenvalues are complex numbers of absolute value slightly less than  $\pi$  lying approximately on a semicircle in the right half-plane. If Method 2 is applied to compute  $\log(\mathbf{A})$  with parameters  $m = 0.25$  and  $M = 8$  (these parameters were obtained experimentally, not on the basis of analysis of the problem), and with the imaginary part of the quadrature points in (2.6) shifted from  $0.5K'$  to  $0.6K'$ , here are the results:

5	1.31e-02
10	3.99e-05
15	3.53e-07
20	1.58e-09
25	2.76e-12
30	2.08e-14

We see that although the spectrum is complex, thanks to the better condition number, the convergence is nearly as fast as for the example of §3. Figure 10 explains what is going on. If the imaginary part is fixed at  $0.5K'$ , the convergence is about half as fast, becoming about three-quarters as fast in that case if  $M$  is increased to 64. If it is fixed at  $0.7K'$ , the convergence can be improved further, but in this case the quadrature points in the  $z$ -plane cross the branch cut on the negative real axis, so for a successful computation certain values must be negated in accordance with the nonstandard branch of the square root.

This example illustrates that the methods we have put forward can be applied broadly, but also that parameter choices are involved. If these methods are to become the basis of general tools, it will be necessary to develop systematic ways to estimate eigenvalues or related information from which to derive such parameters. For example, our programs `method1`–`method3` choose  $m$  and  $M$  by computing the eigenvalues of  $A$  and then finding their minimum and maximum. In practice this would usually be impractical, and in any case the codes as written will give the wrong result if the eigenvalues are complex.

## 7 Three more examples

We conclude with three more numerical examples illustrating a highly nonnormal matrix of small dimension, a fractional matrix power of medium dimension, and the square root

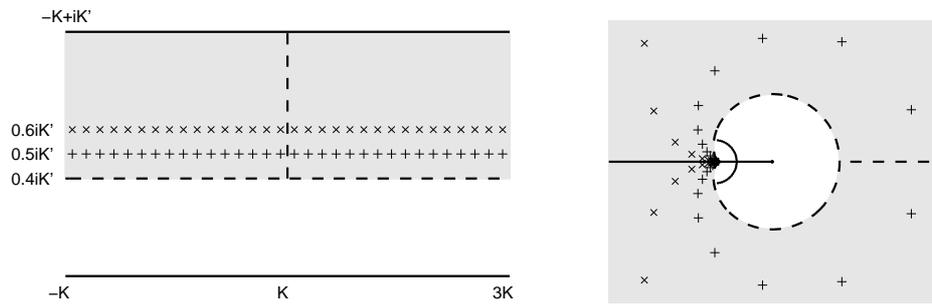


Figure 10: Example of the  $32 \times 32$  Parter matrix. The plot on the left shows the rectangle in the  $t$ -plane with 32 quadrature points at heights  $0.5K'$  (plusses) and  $0.6K'$  (crosses). The plot on the right shows their images in the  $z$ -plane under the Method 2 map with parameters  $m = 0.25$ ,  $M = 8$ . The dashed curve on the right encloses the eigenvalues of  $\mathbf{A}$ , which appear fused in a semicircle, and corresponds to the line at height  $0.4K'$  on the left. Above this line, in the shaded strip, the integrand will be analytic, so convergence is guaranteed at a rate corresponding to the lower half-width  $0.1K'$  for the plusses and  $0.2K'$  for the crosses, that is, 20% and 40% of the rate given by Theorem 2 for a matrix with eigenvalues in  $[0.25, 8]$ .

of a sparse matrix of large dimension.

By a highly nonnormal matrix we mean a matrix whose eigenvectors, though a complete set may exist, are far from orthogonal. For an example of this kind we take a Frank matrix, an upper-Hessenberg matrix with integer entries  $a_{ij} = n + 1 - j$  on and above the diagonal and  $a_{ij} = n + 1 - i$  on the subdiagonal: MATLAB's `gallery('frank', n)`. All the eigenvalues are real and positive, and it is well-known that the condition numbers of the smaller eigenvalues increase rapidly as  $n$  increases. Thus the matrix is highly non-normal. For the  $12 \times 12$  Frank matrix we computed the square root using a modified version of `method3.m`, with parameters set to their true values  $m \approx 0.031$ ,  $M \approx 32.2$ . For the exact square root we took the matrix computed in 64-digit arithmetic via diagonalization using MATLAB's Symbolic Math Toolbox. The convergence shown in Figure 11 is rapid, the error reaching  $1.7 \times 10^{-10}$  at  $N = 12$  and thereafter levelling off due to the effects of rounding errors. Convergence need not be monotonic, as is confirmed by the sharp increase in error at  $N = 9$ . It is interesting to note that the error in the square root computed by `sqrtm` is  $2 \times 10^{-9}$ , so contour integration is giving about one more digit of accuracy than the Schur method, which is the standard method of choice if a Schur decomposition can be computed [11].

The dashed curve in Figure 11, corresponding to a diagonal matrix with the same eigenvalues as the Frank matrix, shows that in this example the nonnormality has had a modest effect on the convergence and a great effect on the final accuracy. The reason for the latter is that Method 3 requires the computation of various matrix inverses  $(w(t_j)^2 \mathbf{I} - \mathbf{A})^{-1}$  (or solution of corresponding systems of equations, if computing  $\mathbf{A}^{1/2}b$  instead of  $\mathbf{A}^{1/2}$ ). The numbers  $w(t_j)^2$  are far from the spectrum, but not from the  $\varepsilon$ -pseudospectra for  $\varepsilon \ll 1$ , so these matrix problems are well-conditioned in the normal case but highly ill-conditioned in the nonnormal case, causing great amplification of

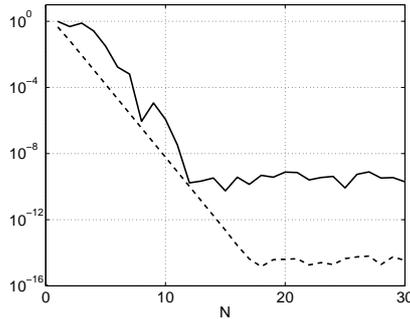


Figure 11: Convergence of Method 3 in computing  $\mathbf{A}^{1/2}$  for  $\mathbf{A} = \text{gallery}(\text{'frank'}, 12)$  (solid curve) and for a diagonal matrix with the same eigenvalues (dashed).

rounding errors. The right solution would be to employ an integration contour that avoids the pseudospectra as well as the spectrum.

For a second example, we use Method 2 to compute  $\mathbf{A}^{1/7}b$ , where  $b$  is a random vector and  $\mathbf{A}$  is a  $598 \times 598$  real nonsymmetric matrix with positive eigenvalues and  $m \approx 2.5$ ,  $M \approx 6.1 \times 10^9$ . To be precise,  $\mathbf{A}$  is the Chebyshev spectral differentiation matrix obtained by the MATLAB commands  $\mathbf{A} = \text{gallery}(\text{'chebspec'}, 600)^2$  followed by  $\mathbf{A} = \mathbf{A}(2:\text{end}-1, 2:\text{end}-1)$ . One way to compute  $\mathbf{A}^{1/7}b$  is by means of the identity  $\mathbf{A}^{1/7} = \exp(\log(\mathbf{A})/7)$ , or  $\text{expm}(\logm(\mathbf{A})/7)$  in MATLAB, and on a 2003 laptop, this takes 32 seconds. Another approach is first to reduce  $\mathbf{A}$  orthogonally to Hessenberg form, then solve  $N$  Hessenberg systems of equations in the algorithm of Method 2. With  $N = 40$  this gets comparable accuracy in about 4 seconds. By increasing the dimension, one could increase the disparity in timings.

Our final example (Table 1) applies Method 3 to compute  $\mathbf{A}^{1/2}b$ , where  $\mathbf{A}$  is the matrix generated by the MATLAB command  $\text{gallery}(\text{'poisson'}, n)$ , an  $n^2 \times n^2$  matrix corresponding to a standard 5-point finite difference discretization of the Laplacian, and  $b$  is the vector of dimension  $n^2$  of all ones. For  $m$  and  $M$  we take the estimates  $2\pi^2/(n+1)^2$  and 8, respectively. This example shows that the methods described here are capable of handling large sparse matrices. The final column of the table compares against MATLAB's built-in method of forming  $\text{sqrtm}(\mathbf{A})$  (the full matrix) and then multiplying by  $b$ , a method that is impractical for larger dimensions.

## References

- [1] M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions*, Dover, 1972 (originally published in 1964).
- [2] N. I. Achieser, *Theory of Approximation*, Frederick Ungar, 1956.
- [3] M. Benzi and G. H. Golub. Bounds for the entries of matrix functions with applications to preconditioning. *BIT*, 39(3):417–438, 1999.

Table 1: Computation of  $\mathbf{A}^{1/2}b$  by Method 3 for discrete Laplacians of various dimensions. In each case  $N$  is chosen large enough for 10 digits of relative accuracy, and time is measured in seconds on a 2003 laptop.

$n^2$	$N$	$M/m$	time	time (sqrtm)
16	8	10.1	0.001	0.0006
64	9	32.8	0.02	0.005
256	10	117.	0.04	0.2
1024	12	441.	0.2	21.
4096	14	1712.	1.0	26 minutes
16384	15	6744.	6.0	< 1 day?

- [4] P. I. Davies and N. J. Higham. Computing  $f(A)b$  for matrix functions  $f$ . In Artan Borici, Andreas Frommer, Báalint Joó, Anthony Kennedy, and Brian Pendleton, editors, *QCD and Numerical Analysis III*, volume 47 of *Lecture Notes in Computational Science and Engineering*, pages 15–24. Springer-Verlag, Berlin, 2005.
- [5] P. J. Davis, “On the numerical integration of periodic analytic functions”, in R. E. Langer, ed., *On Numerical Integration: Proceedings of a Symposium, Madison, April 21–23, 1958*, Math. Res. Ctr., U. of Wisconsin, 1959, pp. 45–59.
- [6] P. J. Davis and P. Rabinowitz, *Methods of Numerical Integration*, 2nd ed., Academic Press, New York, 1984.
- [7] T. A. Driscoll. The Schwarz–Christoffel toolbox. <http://www.math.udel.edu/~driscoll/software/SC/>.
- [8] T. A. Driscoll. Algorithm 843: Improvements to the Schwarz–Christoffel toolbox for MATLAB. *ACM Trans. Math. Software*, 31(2):239–251, 2005.
- [9] T. A. Driscoll and L. N. Trefethen, *Schwarz–Christoffel Mapping*, Cambridge U. Press, 2002.
- [10] I. P. Gavriljuk and V. L. Makarov, Exponentially convergent parallel discretization methods for the first order evolution equations, *Comp. Meth. Appl. Math.* 1 (2001), 333–355.
- [11] N. J. Higham, *Functions of Matrices: Theory and Computation*, book to appear.
- [12] M. Iri, S. Moriguti and Y. Takasawa, On a certain quadrature formula (in Japanese), Kokyuroku RIMS, Kyoto Univ., 91 (1970), 82–119, translated into English in *J. Comp. Appl. Math.* 17 (1987), 3–20.
- [13] A. D. Kennedy, Approximation theory for matrices, *Nucl. Phys. B* 128 (2004), 107–116.

- [14] C. S. Kenney and A. J. Laub, The matrix sign function, *IEEE Trans. Automat. Control* 40 (1995), 1330–1348.
- [15] E. Martensen, Zur numerischen Auswertung uneigentlicher Integrale, *Zeit. Angew. Math. Mech.* 48 (1968), T83–T85.
- [16] W. McLean and V. Thomée. Time discretization of an evolution equation via Laplace transforms. *IMA J. Numer. Anal.*, 24(3):439–463, 2004.
- [17] M. Mori, Discovery of the double exponential transformation and its developments, *Publ. RIMS, Kyoto U.* 41 (2005), 897–935.
- [18] D. J. Newman, Rational approximation to  $|x|$ , *Michigan Math. J.* 11 (1964), 11–14.
- [19] S.-D. Poisson, Sur le calcul numérique des intégrales définies, *Mémoires de L'Académie Royale des Sciences de L'Institut de France* 4 (1827), pp. 571–602 written in 1823.
- [20] T. W. Sag and G. Szekeres, Numerical evaluation of high-dimensional integrals, *Math. Comp.* 18 (1964), 245–253.
- [21] T. Schmelzer and L. N. Trefethen, Computing the gamma function, *SIAM J. Numer. Anal.* 45 (2007), 558–571.
- [22] C. Schwartz, Numerical integration of analytic functions, *J. Comp. Phys.* 4 (1969), 19–29.
- [23] D. Sheen, I. H. Sloan, and V. Thomée, A parallel method for time-discretization of parabolic problems based on contour integral representation and quadrature, *Math. Comp.* 69 (1999), 177–195.
- [24] K.-G. Steffens, *The History of Approximation Theory: From Euler to Bernstein*, Birkhäuser, 2006.
- [25] F. Stenger, Integration formulae based on the trapezoidal formula, *J. IMA* 12 (1973), 103–114.
- [26] F. Stenger, Numerical methods based on Whittaker cardinal, or sinc functions, *SIAM Review* 23 (1981), 165–224.
- [27] F. Stenger, *Numerical Methods Based on Sinc and Analytic Functions*, Springer-Verlag, New York, 1993.
- [28] H. Takahasi and M. Mori, Quadrature formulas obtained by variable transformation, *Numer. Math.* 12 (1973), 206–219.
- [29] H. Takahasi and M. Mori, Double exponential formulas for numerical integration, *Publications RIMS, Kyoto U.* 9 (1974), 721–741.

- [30] A. Talbot, The accurate numerical inversion of Laplace transforms, *J. Inst. Maths. Applics.* 23 (1979), 97–120.
- [31] N. M. Temme, *Special Functions*, Wiley, New York, 1996.
- [32] L. N. Trefethen and M. H. Gutknecht, The Carathéodory–Fejér method for real rational approximation, *SIAM J. Numer. Anal.* 20 (1983), pp. 420–436.
- [33] L. N. Trefethen and J. A. C. Weideman, The trapezoid rule in scientific computing, article for *SIAM Review* in preparation.
- [34] L. N. Trefethen, J. A. C. Weideman and T. Schmelzer, Talbot contours and rational approximation, *BIT Numer. Math.* 46 (2006), 653–670.
- [35] J. van den Eshof, A. Frommer, Th. Lippert, K. Schilling, and H. A. van der Vorst, Numerical methods for the QCD overlap operator. I. Sign-function and error bounds, *Computer Phys. Commun.* 146 (2002), 203–224.
- [36] J. A. C. Weideman, Optimizing Talbot’s contours for the inversion of the Laplace transform, *Math. Comp.*, to appear.
- [37] J. A. C. Weideman and L. N. Trefethen, Parabolic and hyperbolic contours for computing the Bromwich integral, *Math. Comp.*, to appear.