

## SOLVING BURGERS' EQUATION USING

approximation methods for solving partial differential and integral equations in higher dimensions, where the ability to construct near optimal rational (or exponential) approximations to functions of one variable is a key component.

Since the seminal result in [21], it has been known that functions with singularities may be efficiently approximated in the  $\mathbf{L}^\infty$  norm using proper rational functions. Indeed, the number of poles required to approximate a function with singularities is directly related to the sparsity of the function's wavelet coefficients (see [16, Theorem 11.1]). However, in contrast to more traditional  $\mathbf{L}$ -type methods (using e.g., wavelet bases as in [2]), the use of such optimal  $\mathbf{L}^\infty$ -type approximations in numerical analysis has been limited due to a lack of efficient and robust algorithms.

Given a proper rational function  $\mathbf{f}$ , we present an algorithm—which we refer to as the reduction algorithm—to compute, for a fixed number 20.8801030e

many times. For example, in the context of solving Burgers' equation with viscosity  $\nu = 10^{-6}$  and approximation tolerance  $\epsilon = 10^{-6}$ , on the order of a million applications of the reduction algorithm are performed.

For functions with  $n$  poles resulting from intermediate computations, the reduction algorithm requires only  $\mathcal{O}(m n)$  operations to find an optimal approximation with  $m$  poles. In our numerical experiments with the reduction algorithm, we find that an approximation error of  $10^{-6}$  may be reliably obtained within double precision arithmetic, even when the number of poles  $n$  is large and their spatial distribution is highly clustered.

There is a significant literature devoted to applications of the AAK approach in control theory (cf. [23]), signal processing (cf. [8]), and numerical analysis (cf. [25, 27, 29, 5]), to mention just a select few. The reformulation of the AAK theory given here could be related to the approaches taken in [28], [20], and [10]. However, as far as we know, all of the AAK-type algorithms discussed in the literature require  $\mathcal{O}(n^3)$  operations when applied to a rational function with  $n$  poles, and may require extended precision arithmetic if high accuracy of the result is desired. In contrast, our reduction algorithm requires only  $\mathcal{O}(m n)$  operations to find an optimal approximation with  $m$  poles and achieves high accuracy ( $\epsilon = 10^{-6}$ ) using only double precision arithmetic.

We show in this paper that solutions of Burgers' equation with viscosity  $\nu$  require only  $\mathcal{O}(\log \frac{1}{\nu}) + \mathcal{O}(\log \frac{1}{\epsilon})$  poles for its rational approximation with an  $L^\infty$  error of size  $\epsilon$ . Burgers' equation has been traditionally used to test the limits of new numerical methods since the solution develops sharp transition regions that need to be captured adaptively. Conceptually, the two closest adaptive methods are those in [24] and [2]. While in [2] adaptivity is achieved by adding wavelet scales when needed, the algorithm in [24] achieves spectral accuracy by adding Chebyshev nodes.

smaller number of poles. As mentioned earlier, our reduction algorithm is based on a theorem of Adamyán, Arov, and Krein ([1]), which concerns the approximation of a periodic function  $\mathbf{f}$ , essentially bounded on the unit circle  $\mathbb{D}$ , by a meromorphic function  $\mathbf{r}(\mathbf{z})$  ( $\mathbf{z} = e^{-i\mathbf{x}}$ ) containing a specified number of poles in the unit disk. We limit our presentation to rational functions  $\mathbf{f}$  taking real values on  $\mathbb{D}$ . This case turns out to be particularly important, as it allows us to develop a practical algorithm based on approximating the Fourier series coefficients of  $\mathbf{f}$  with positive index. More general functions  $\mathbf{f}$  may be dealt with by using the techniques in

The fact that there are exactly  $m$  zeros in the unit disk, corresponding to the index  $m$  of the con-eigenvalue  $\lambda_m$ , is a consequence of the AAK theory. As shown in Section 4.1 (see equations (4.8)), the key to the high accuracy of evaluating the function  $v(z)$  is the relationship

$$(2.4) \quad v(z) = \sum_{i=1}^m \frac{u_i}{1 - \bar{z} z_i}, \quad i = 1, \dots, m,$$

which, together with the  $m$  poles  $1/\bar{z}_i$ , uniquely determines  $v(z)$ .

**Step 3:** Find the coefficients  $u_i$  of  $g(z)$  by solving the  $m \times m$  linear system,

$$(2.5) \quad \sum_{i=1}^m \frac{1}{1 - \bar{z}_i z_j} u_i = \sum_{i=1}^m \frac{z_j}{1 - \bar{z}_i z_j}, \quad j = 1, \dots, m.$$

Denoting  $\|f - g\|_\infty = \sup_{x \in [0, 2\pi]} |f(e^{-ix}) - g(e^{-ix})|$ , the resulting rational approximation  $g(e^{-ix})$  satisfies  $\|f - g\|_\infty \leq \frac{1}{2m}$  and, thus, is close to the best  $L^\infty$ -error achievable by rational functions with no more than  $m$  poles in the unit disk (see also [25] for a discussion of optimal rational approximations).

*Remark 1.* In Step 3, we solve for the coefficients  $u_i$  in  $O(m^3)$  operations by exploiting the structure of Cauchy matrices (see [11, 7]). We note that such a solver may require quadruple precision if the overall desired approximation error is smaller than  $10^{-16}$ . However, since  $m = \log(\frac{1}{\epsilon})$  is small, Step 3 for finding coefficients  $u_i$  does not impact the overall speed of the algorithm even if performed in quadruple precision.

*Remark 2.* In applications where the function  $f(e^{-ix})$  has singularities or sharp transitions, the poles  $z_j$  in the rational representation of  $f(e^{-ix})$  may be located very close to the unit circle (and/or to each other). In such cases, it is advantageous to maintain the poles in the form  $z_j = \exp(-\beta_j)$ , since they are well separated on a logarithmic scale. Importantly, the reduction algorithm computes the new poles

( $\mathbf{r} < \mathbf{1}$ ). This estimate shows that, for accuracy  $\epsilon$ , we may reasonably expect  $\mathbf{O}(\log \frac{1}{\epsilon})$  terms in our approximation. In fact, we have observed this behavior in our numerical experiments.

Let us now briefly discuss the algorithmic aspects behind efficiency and accuracy of solving steps **1-3** above.

**2.2. Accurate computation of con-eigenvalues/eigenvectors.** For Step **1**, we use a recent algorithm developed and analyzed in [?] for computing con-eigenvalues of Cauchy matrices with high relative accuracy, which we briefly describe in this section.

It is well-known that standard eigenvalue algorithms compute an approximate con-eigenvalue  $\widehat{\lambda}$

using (2.4) to rewrite (2.3) as

$$\sum_{i=1}^n \frac{\overline{v_i}(z)}{1 - \overline{z}} = m v(z),$$

we see that

SOLVING



to the known eigenvalues  $\lambda_1, \dots, \lambda_{m-1}$ , one by one. We then orthogonalize these  $m-1$  vectors using the stabilized Gram-Schmidt procedure, thus yielding a basis  $\hat{q}_1, \dots, \hat{q}_{m-1}$  for the invariant subspace  $\text{span}\{q_1, \dots, q_{m-1}\} = \text{span}\{\hat{q}_1, \dots, \hat{q}_{m-1}\}$ . Finally, we use simultaneous inverse iteration applied to  $\hat{q}_1, \dots, \hat{q}_{m-1}, q$ , where  $q$  is chosen randomly. Notice that each step of this process

§ess

depend on the timestep  $\mathbf{t}$ , the number  $\mathbf{M}_t$  of quadrature nodes in time, and the number of quadrature nodes  $\mathbf{M}_x$  used in space to discretize the convolution kernels. From the rapid decay of the periodic heat kernel,

$$\mathbf{K}(\mathbf{x}, \mathbf{t}) = \frac{1}{4\sqrt{\mathbf{t}}} \sum_{\mathbf{k} \in \mathbb{Z}} e^{-(\mathbf{x}+\mathbf{k})^2/(4\mathbf{t})},$$

where  $\nu$  is the viscosity parameter in (3.1), it follows that  $\frac{1}{4\sqrt{\mathbf{t}}}$  and  $\frac{1}{4\sqrt{\mathbf{t}}}$  are localized to a  $\mathcal{O}(\sqrt{\mathbf{t}})$  neighborhood of  $\mathbf{x} = 0$  (see Section 4.2 for details).

We assume that the initial function  $\mathbf{u}(\mathbf{x}, 0) = \mathbf{u}(\mathbf{x})$  is given as a periodic rational function of the form

$$\mathbf{u}(\mathbf{x}) = \sum_{\mathbf{j}=\mathbf{1}}^{\mathbf{M}_0} \frac{\mathbf{j}}{e^{-i\mathbf{x}} - \mathbf{j}} + \sum_{\mathbf{i}=\mathbf{1}}^{\mathbf{M}_0} \frac{\overline{\mathbf{j}}}{e^{-i\mathbf{x}} - \overline{\mathbf{j}}} + \dots,$$

and that this representation is nearly optimal. We then solve the system of equations (3.2) by approximating each function  $\mathbf{u}_l$  using the reduction algorithm. We obtain, via fixed point iteration applied to (3.2) and the reduction algorithm, rational functions  $\mathbf{u}_l(\mathbf{x})$  of the form,

$$(3.3) \quad \mathbf{u}_l(\mathbf{x}) = \sum_{\mathbf{j}=\mathbf{1}}^{\mathbf{M}_l} \frac{\mathbf{j},l}{e^{-i\mathbf{x}} - \mathbf{j},l} + \sum_{\mathbf{j}=\mathbf{1}}^{\mathbf{M}_l} \frac{\overline{\mathbf{j},l}}{e^{-i\mathbf{x}} - \overline{\mathbf{j},l}} + \dots,$$

which solve (3.2) to a specified level of precision, and have a (near) optimal number of poles.

More specifically, given  $\mathbf{u}_j^{(m)} = \mathbf{u}_j(\mathbf{x})$ ,  $1 \leq j \leq \mathbf{M}_t$ , at iteration  $m$ , we use (3.2) to define the next iterates  $\mathbf{u}_l^{(m+1)}(\mathbf{x})$  for  $l = 1, \dots, \mathbf{M}_t$ ,

$$(3.4) \quad \mathbf{u}_l^{(m+1)}(\mathbf{x}) = \sum_{\mathbf{p}=\mathbf{1}}^{\mathbf{M}_x} \frac{1}{\mathbf{p}} \mathbf{u}_l^{(m)}\left(\mathbf{x} - \frac{\mathbf{i}}{\mathbf{p}}\right) + \sum_{\mathbf{p}=\mathbf{1}}^{\mathbf{M}_x} \frac{1}{\mathbf{p}} \sum_{\mathbf{j}=\mathbf{1}}^{\mathbf{M}_t} \left\{ \left( \mathbf{u}_j^{(m+1)}\left(\mathbf{x} - \frac{\mathbf{i}}{\mathbf{p}}\right) \right) \right\} + \sum_{\mathbf{j}=\mathbf{1}}^{\mathbf{M}_t}$$





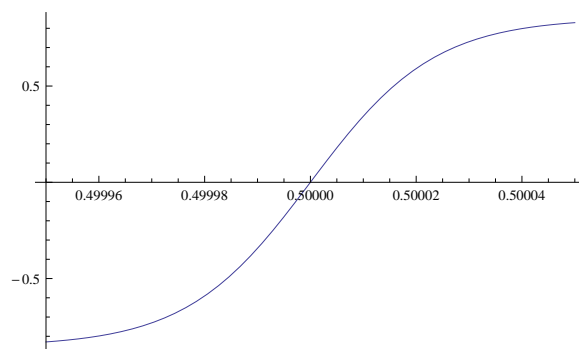
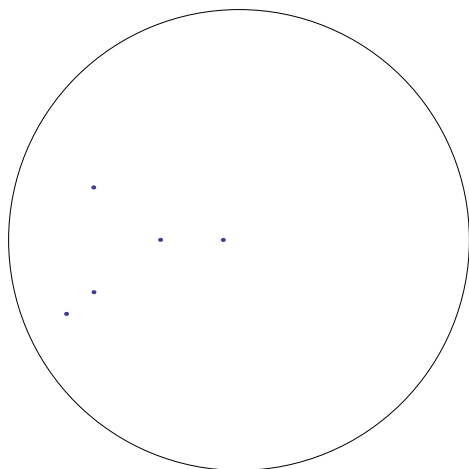


FIGURE 3.3. Solution  $\mathbf{u}(\mathbf{x}, \mathbf{t})$  at time  $\mathbf{t} = .4$ , localized about the transition region  $(1/2 - 10^{-4}, 1/2 + 10^{-4})$ . Note the absence of any Gibbs-type phenomena.



Suppose  $\mathbf{f} \in \mathbf{L}^\infty$  has the

we calculate from (4.4)

$$\begin{aligned} \sum_{j=0}^{\infty} \left( \sum_{m=0}^M m^{i+j} \right) v_j &= \sum_{m=0}^M m^i \sum_{j=0}^{\infty} m^j v_j \\ &= \sum_{m=0}^M m^i v(m) = w_i. \end{aligned}$$

Now multiplying both sides of the last equation by  $z^i$  and summing, we obtain

$$(4.6) \quad \sum_{m=0}^M \frac{m^i}{1 - mz} v(m) = z^i w(z).$$

Similarly, from (4.5), we have

$$\begin{aligned} \sum_{j=0}^{\infty} \left( \sum_{m=0}^M \frac{m^{i+j}}{m} \right) w_j &= \sum_{m=0}^M \frac{m^i}{m} \sum_{j=0}^{\infty} \frac{m^j}{m} w_j \\ &= \sum_{m=0}^M \frac{m^i}{m} \left( \frac{1}{m} w \left( \frac{1}{m} \right) \right) = v_i. \end{aligned}$$

Finally, multiplying by  $z^i$  and summing, we arrive at

$$(4.7) \quad \sum_{m=0}^M \frac{1}{1 - \frac{z}{m}} \frac{1}{m} w \left( \frac{1}{m} \right) = v(z).$$

Hence, for a function  $f$  of the form (4.3), the functions  $v$  and  $w$  in (4.2) turn out to be rational and fully determined by their values at the poles of  $f$ . Taking  $z = \frac{1}{n}$  and  $z = \frac{1}{n}$  in equations (4.6) and (4.7), respectively, we obtain

$$\sum_{m=0}^M m^i$$

Let us define the vectors  $\mathbf{p}$  and  $\mathbf{q}$  with entries  $p_m = \frac{1}{2}v(m)$ ,  $q_m = \frac{1}{m}w(m)$ , and the positive definite matrix  $\mathbf{C}$  with entries

$$C_{mn} = \frac{\frac{1}{2} \frac{1}{m} \frac{1}{n}}{1 - \frac{1}{m} \frac{1}{n}}.$$

Then the above equations are equivalent to

$$\begin{aligned} \mathbf{C} \mathbf{p} &= \mathbf{q}, \\ \mathbf{C} \mathbf{q} &= \mathbf{p}, \end{aligned}$$

which may be reduced to a con-eigenvalue problem for  $\lambda > 0$ , see [15, Section 4.6]. One simple way to see this and obtain an equation of the form (2.2) is by defining  $\mathbf{x} = \mathbf{p} + \mathbf{q}$ . If  $\mathbf{x} = \mathbf{0}$ , then  $\mathbf{iq} = \overline{\mathbf{ip}}$  and hence

$$\mathbf{C}(\mathbf{ip}) = \overline{\mathbf{ip}}.$$

If  $\mathbf{x} \neq \mathbf{0}$ , we have

$$\mathbf{C} \mathbf{x} = \lambda \mathbf{x}$$

and, in both cases, we obtain a con-eigenvalue problem for the matrix  $\mathbf{C}$ .

**4.2. Discretization of Burgers' equation.** We rewrite the equation (3.1) in semi-group form (see, e.g., [14, 17, 18, 3])

$$(4.9) \quad \mathbf{u}(t) = e^{t\mathbf{L}}\mathbf{u}(0) + \int_0^t e^{(t-\tau)\mathbf{L}}\mathbf{N}(\mathbf{u}(\tau))d\tau,$$

where  $\mathbf{u}(t)$  denotes the function  $\mathbf{u}(\cdot, t)$ . The operator  $\mathbf{L}$ ,  $\mathbf{L}\mathbf{u}(\mathbf{x}) = \mathbf{u}_{\mathbf{x}\mathbf{x}}$ , represents the linear part of (3.1) while the operator  $\mathbf{N}$ ,  $\mathbf{N}(\mathbf{u}) = \frac{1}{2}(\mathbf{u})_{\mathbf{x}}$ , represents the nonlinear part. The action of the operator  $e^{t\mathbf{L}}$  on a function  $\mathbf{f}$  is given by

$$(e^{t\mathbf{L}}\mathbf{f})(\mathbf{x}) = \int_{-\frac{1}{2}}^{\frac{1}{2}} \mathbf{K}(\mathbf{y}, t)\mathbf{f}(\mathbf{x} - \mathbf{y})d\mathbf{y}, \text{ with } \mathbf{K}(\mathbf{y}, t) = \frac{1}{4t} \sum_{\mathbf{k} \in \mathbb{Z}} e^{-(\mathbf{y}+\mathbf{k})^2/(4t)}.$$

To discretize equation (4.9) in time, we use the approximation

$$\mathbf{N}(\mathbf{u}(\tau)) \approx \sum_{j=1}^{M_t} \mathbf{R}_j(\tau)\mathbf{N}(\mathbf{u}(\tau_j)), \quad \tau \in [0, t]$$

where  $\{\tau_j\}_{j=1}^{M_t}$  denote the Gauss-Legendre nodes on the interval  $(0, t)$ , and  $\mathbf{R}_j(\tau)$  denote the Legendre interpolating polynomials for these nodes, i.e.,

$$\mathbf{R}_j(\tau_m) = \delta_{jm}, \quad \text{for } j, m = 1, \dots, M_t.$$

Taking  $\tau = \tau_l$  in (4.9), we obtain the semi-discrete system of equations

$$(4.10) \quad \mathbf{u}_l = e^{\tau_l \mathbf{L}} \mathbf{u} + \sum_{j=1}^{M_t} \left( \int_0^{\tau_l} e^{(\tau_l - \tau)\mathbf{L}} \mathbf{R}_j(\tau) d\tau \right) \mathbf{N}(\mathbf{u}_j), \quad 1 \leq l \leq M_t,$$

where  $\mathbf{u}_l = \mathbf{u}_l(\mathbf{x})$  denote the computed values of  $\mathbf{u}$  at time  $\tau = \tau_l$  and  $\mathbf{u} = \mathbf{u}(\mathbf{x}, 0)$ .



For the spatial discretization, using  $\mathbf{N}(\mathbf{u}) = 1/2 (\mathbf{u})_{\mathbf{x}}$  and

- [9] G. Dahlquist and Å. Björck. *Numerical methods in scientific computing. Vol. I.* Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2008.
- [10] Ph. Delsarte, Y. Genin, and Y. Kamp. On the role of the Nevanlinna-Pick problem in circuit and system theory. *International Journal of Circuit Theory and Applications*, 9 (2):177–187, 1981.
- [11] J. Demmel. Accurate singular value decompositions of structured matrices. *SIAM J. Matrix Anal. Appl.*, 21 (2):562–580, 1999.
- [12] J. Demmel, M. Gu, S. Eisenstat, I. Slapnicar, K. Veselic, and Z. Drmac. Computing the singular value decomposition with high relative accuracy. *LAPACK Working Note*, 119 (CS-97-48), 1997.
- [13] J. Demmel and K. Veselic. Jacobi's method is more accurate than QR. *SIAM J. Matrix Anal. and Appl.*, 1 (4):1204–1245, 1992.
- [14] E. Hille and R. S. Phillips. *Functional Analysis and Semi-groups.* American Mathematical Society, Providence, RI, 1957.