

# NEAR OPTIMAL RATIONAL APPROXIMATIONS OF LARGE DATA SETS

ANIL DAMLE, GREGORY BEYLKIN, TERRY HAUT AND LUCAS MONZÓN

ABSTRACT. We introduce a new computationally efficient algorithm for constructing near optimal rational approximations of large (one-dimensional) data sets. In contrast to wavelet-type approximations, these new approximations are effectively shift invariant. We note that the complexity of current algorithms for computing near optimal rational approximations prevents their use for large data sets.

In order to obtain a near optimal rational approximation of a large data set, we first construct its intermediate B-spline representation. Then, by using a new rational approximation of B-splines, we arrive at a suboptimal rational approximation of the data set. We then use a recently developed fast and accurate reduction algorithm for obtaining a near optimal rational approximation

in contrast to wavelet decompositions, rational functions are closed under translations and, thus, optimal rational approximations are shift invariant. Indeed, shifting an optimal rational approximation yields the optimal approximation of the shifted function or data.

Our rational representations are optimal in the sense that, for a given accuracy of approximation, the number of poles is minimal. We say that the approximation is “near optimal” if, instead of the desired accuracy  $\epsilon$ , our algorithms yield accuracy  $\epsilon'$ , where  $\epsilon'$  is slightly smaller than  $\epsilon$ . In such case the number of poles may not be minimal in the strict sense (we note that we have an *a posteriori* check to identify such situation, if needed). We use the term “suboptimal”, if we know that the number of poles definitely exceeds the optimal number (for a given accuracy).

For functions given analytically or for functions described by a relatively small number of samples, there are several methods for obtaining their near optimal rational approximations [5, 6, 7]. For a large data set these methods are impractical due to their computational complexity. On the other hand, computing a wavelet decomposition of a large data set does not present a difficulty since its computational cost is linear in the number of samples; we use these facts in our approach.

We first compute a B-spline representation of the data, which provides a simple and efficient method for a transition to a suboptimal rational representation. For this purpose, we construct a new rational approximation of B-splines, where the poles are arranged on a rectangular grid aligned with the location of spline knots. We then split the data into large segments, and compute suboptimal rational approximations for each segment. Finally, we compute a near optimal rational approximation using a recently developed, fast and accurate algorithm in [10].

Although the example provided here is compression of audio recordings, the algorithm may be used to compress and analyze other types of signals, e.g., signals obtained by continuous, global seismic monitoring. In particular, we view compression via near optimal rational approximations as the first step in signal analysis since the poles carry frequency and time information. As shown in [6], poles of near optimal rational approximations concentrate near the singularities of functions. For signals, this corresponds to locations of rapid change, such as occurring when a piano key is struck or at wave arrivals in seismic recordings. The location of the poles also carries information about local frequency content of the signal in a manner similar to wavelets, i.e., the logarithmic distance of these locations from the real axis corresponds to wavelet scales.

We start in Section 2 by providing the background information on the key existing algorithms that facilitate the development of our new approach. Next, in Section 3, we construct a rational approximation (with special properties) of a B-spline to be used in intermediate computations. Then, in Section 4, we describe in detail the algorithm for constructing near optimal

NEAR OPTIMAL RATIONAL

Following [6], from (1) we obtain the rational representation

$$\mathbf{f}(\mathbf{x}) = -2\mathcal{R}e \left( \sum_{\mathbf{j}=1}^{\mathbf{M}} \mathbf{w}_{\mathbf{j}} \right)$$

. In contrast to standard algorithms, the con-eigenvalues (and con-eigenvectors) are computed with high relative accuracy in  $\mathcal{O}(\mathbf{M}^2\mathbf{M}_0)$  operations.

Find all the roots inside the unit disk of the function

$$\mathbf{v}(\mathbf{z}) = \frac{1}{\mathbf{M}} \sum_{j=1}^{\mathbf{M}_0} \frac{\sqrt{d_j} u_j}{1 - \mu_j \mathbf{z}}.$$

Note that there are exactly  $\mathbf{M}$  roots  $\mu_j$  of  $\mathbf{v}(\mathbf{z})$  inside the unit disk based on results from [3].

Finally solve for the residuals  $\mathbf{r}_j$  of  $\mathbf{r}(\mathbf{z})$  by solving the  $\mathbf{M} \times \mathbf{M}$  linear system

$$\sum_{j=1}^{\mathbf{M}} \frac{\mathbf{r}_j}{1 - \mu_j \bar{\mathbf{k}}} = \sum_{j=1}^{\mathbf{M}_0} \frac{d_j}{1 - \mu_j \bar{\mathbf{k}}}.$$

Using this algorithm, we obtain  $\mathbf{f} - \mathbf{r} \approx \mathbf{m}$ , which provides a near optimal representation of  $\mathbf{f}(\mathbf{z})$  using only  $\mathbf{M}$  pairs of conjugate-reciprocal poles,  $\mu_j$  and  $\bar{\mu}_j^{-1}$ . The computational complexity of this algorithm is  $\mathcal{O}(\mathbf{M}^2\mathbf{M}_0)$ , where  $\mathbf{M}$  is the number of resulting poles and  $\mathbf{M}_0$  is the original number of poles. Since typically  $\mathbf{M} \ll \mathbf{M}_0$ , this algorithm is effectively linear in its practical use.

**2.3. Spline representations.** We use an intermediate representation via B-splines as the first step towards computing the (near) optimal rational approximation. Although theoretically we may use scaling functions of any wavelet-type basis, the choice of B-splines reduces the computational cost of this intermediate step.

We recall the definition of the  $\mathbf{m}^{\text{th}}$  degree B-spline as

$$\mathbf{m}(\mathbf{x}) = \mathbf{m}-1(\mathbf{x}) \circ \mathbf{o}(\mathbf{x}),$$

with

$$\mathbf{o}(\mathbf{x}) = \begin{cases} 1, & \mathbf{x} \in \left[-\frac{1}{2}, \frac{1}{2}\right] \\ 0, & \text{otherwise,} \end{cases}$$

(see e.g., [8]). For convenience, we only use B-splines of odd degree. It is easy to show that, in this case,  $\mathbf{m}$  is a piecewise polynomial of degree  $\mathbf{m}$  with knots on the integers and supported on  $[-(\mathbf{m} + 1)/2, (\mathbf{m} + 1)/2]$ . To represent periodic functions, we use periodized versions of B-splines. Let us introduce the 1-periodic function

$$\mathbf{a}_m(\mathbf{z}) = \sum_{\mathbf{j} \in \mathbb{Z}} \widehat{\mathbf{m}}(\mathbf{z} + \mathbf{j})^2 = \sum_{\mathbf{l} = -\frac{\mathbf{m}-1}{2}}^{\frac{\mathbf{m}-1}{2}} \widehat{\mathbf{m}}(\mathbf{z} + \mathbf{l})^2$$

Given a uniformly sampled 1-periodic function  $\mathbf{f}$ , we seek the coefficients  $\mathbf{j}$  such that

$$(7) \quad \mathbf{f}\left(\frac{\mathbf{k}}{2\mathbf{N}}\right) = \sum_{\mathbf{j}=0}^{2\mathbf{N}} \mathbf{j} \mathbf{m}(\mathbf{k}-\mathbf{j}), \quad \mathbf{k} = 0, \dots, 2\mathbf{N}.$$

The algorithm in [4, 12] rapidly computes the coefficients  $\mathbf{j}$  in (7) using the Fast Fourier Transform (FFT). It performs the following steps:

Set  $\hat{\mathbf{f}}_{\mathbf{k}} = \mathbf{f}\left(\frac{\mathbf{k}}{2\mathbf{N}}\right)$  and compute, for  $\mathbf{k} = 0, \dots, 2\mathbf{N}$ ,

$$\hat{\mathbf{f}}_{\mathbf{k}} = \sum_{\mathbf{n}=0}^{2\mathbf{N}} \mathbf{f}_{\mathbf{n}} e^{-\frac{i}{2\mathbf{N}+1} \mathbf{k} \mathbf{n}}$$

using the FFT.

Compute, for  $\mathbf{k} = 0, \dots, 2\mathbf{N}$ ,

$$\hat{\mathbf{g}}_{\mathbf{k}} = \frac{\hat{\mathbf{f}}_{\mathbf{k}}}{\mathbf{a}_{\mathbf{m}}\left(\frac{\mathbf{k}}{2\mathbf{N}+1}\right)}.$$

The B-spline coefficients are now obtained via the FFT as

$$\mathbf{j} = \frac{1}{2\mathbf{N}+1} \sum_{\mathbf{n}=0}^{2\mathbf{N}} \hat{\mathbf{g}}_{\mathbf{n}} e^{\frac{i}{2\mathbf{N}+1} \mathbf{j} \mathbf{n}}, \quad \mathbf{j} = 0, \dots, 2\mathbf{N}$$

This algorithm requires  $\mathcal{O}(\mathbf{N} \log \mathbf{N})$  operations. The details may be found in the appendix in [12].

### 3. RATIONAL REPRESENTATION OF B-SPLINES

In this section we construct rational approximations of B-splines. In our construction we force the real parts of the poles to be integers  $\mathbf{l} \in \mathbf{Z}$ , so that the poles are aligned with the knots of the B-spline. As we explain below, this reduces the cost of intermediate computations.

Specifically, we are looking for a suboptimal rational approximation of the form (5), with poles  $\mathbf{l} \pm i \mathbf{k}$ , so that

$$(8) \quad \left| \mathbf{m}(\mathbf{x}) + 2 \sum_{\mathbf{l}=-\mathbf{m}+}^{\mathbf{m}+} \sum_{\mathbf{k}=1}^{\mathbf{R}} \frac{\mathbf{u}_{\mathbf{k},\mathbf{l}}(\mathbf{x}-\mathbf{l}) - \mathbf{v}_{\mathbf{k},\mathbf{l}} \mathbf{k}}{(\mathbf{x}-\mathbf{l})^2 + \frac{2}{\mathbf{k}}} \right|,$$

where the number of rows of poles,  $\mathbf{R}$ , will be chosen later. We note that the constraint on the real part of the poles arranges them on a rectangular grid (see Figure 2).

We start by computing a near optimal rational approximation of a B-spline following the approach in [6]. For a given  $\mathbf{m}$ , we evaluate  $\hat{\mathbf{f}}$  at a sufficient number of samples; specifically for  $\mathbf{m} = 32$  we have

$$(9) \quad \mathbf{h}_{\mathbf{n}} = \hat{\mathbf{f}}_{\mathbf{m}}\left(\frac{\mathbf{n}}{32}\right), \quad \mathbf{n} = 0, 1, \dots, 60,$$

where

$$(10) \quad \hat{\mathbf{m}}(\cdot) = \left( \frac{\sin}{\cdot} \right)^{\mathbf{m}+1},$$

and use the algorithm in Section 2.1 to construct a near optimal rational approximation.

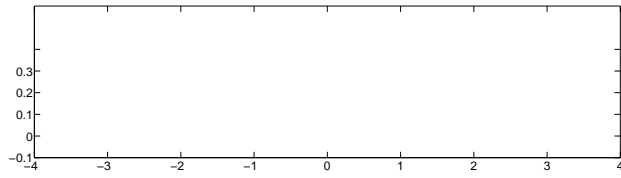
An example of a near optimal rational approximation of a B-spline of degree  $\mathbf{m} = 3$  may be found in [6]. As observed in that paper, the poles concentrate towards the locations of the knots of the B-spline since its third derivative is discontinuous at these points. In our application we would like to use a higher degree B-spline to lessen the impact of discontinuities and obtain fewer poles. In Figure 1 we present the results for a near optimal approximation of a  $\mathbf{m}$ th degree B-spline using the same approach as in [6]. Since the poles,  $\mathbf{t}_j \pm i\mathbf{s}_j$ , appear in complex conjugate pairs, in Figure 1 we display (on a  $\log_{10}$  scale) only those with negative imaginary part.

We then seek a suboptimal rational representation of  $\mathbf{x}(\cdot)$  with poles in the locations indicated in (8) and use the near optimal approximation to select the parameters  $\mathbf{k}$  in (8). Taking into account that the poles closer to the real line are responsible for the high frequency content of the representation, whereas those furthest away capture the lower frequency content, we limit the range for the imaginary parts of our suboptimal poles by using the corresponding maximum,  $\mathbf{s}^+$ , and minimum,  $\mathbf{s}^-$ , of the near optimal poles. We select three rows of poles, i.e.,  $\mathbf{R} = 3$  in (8), by choosing imaginary parts  $\mathbf{s}_1 = \mathbf{s}^+$ ,  $\mathbf{s}_3 = \mathbf{s}^-$ , and

$$\mathbf{s}_2 = e^{-(\log \mathbf{s}^+ + \log \mathbf{s}^-)}.$$

The real part for all of these poles are at locations  $\mathbf{l}$ , where  $\mathbf{l} = -\frac{\mathbf{m}+1}{2}, \dots, \frac{\mathbf{m}+1}{2}$  (recall that  $\mathbf{m}$  is odd). The choice of three rows of poles is based on the degree of the B-spline and our accuracy requirements (see Figure 2(b)) and may be different in other applications.

Once the location of poles is fixed, the weights in (8) are obtained by `solvd[()-M45.96020r7`





the optimization package CVX [9]. The resulting absolute error is shown in Figure 2(b).

#### 4. NEAR OPTIMAL RATIONAL APPROXIMATIONS

We now briefly describe the key steps

compute the B-spline coefficients for each section of the signal. Once the B-spline coefficients for each section are found, by adding com-

the support of both segments, we preserve the overall accuracy of the merged approximation. In our experiments, we reduce the set of poles

parts,

$$\begin{aligned}
\mathbf{f}_k &= \mathbf{f}_k^+ + \mathbf{f}_k^- + \mathbf{f}_k^{\text{local}} \\
(16) \quad &= \sum_{\mathbf{x}_k - \mathbf{t}_j} \frac{1}{s} \left( \frac{\mathbf{u}_j + i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j - i\mathbf{s}_j} + \frac{\mathbf{u}_j - i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j + i\mathbf{s}_j} \right) \\
&+ \sum_{\mathbf{t}_j - \mathbf{x}_k} \frac{1}{s} \left( \frac{\mathbf{u}_j + i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j - i\mathbf{s}_j} + \frac{\mathbf{u}_j - i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j + i\mathbf{s}_j} \right) \\
&+ \sum_{|\mathbf{x}_k - \mathbf{t}_j| < s} \left( \frac{\mathbf{u}_j + i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j - i\mathbf{s}_j} + \frac{\mathbf{u}_j - i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j + i\mathbf{s}_j} \right),
\end{aligned}$$

and evaluate  $\mathbf{f}_k^+$ ,  $\mathbf{f}_k^-$  and  $\mathbf{f}_k^{\text{local}}$  separately, where that of  $\mathbf{f}_k^{\text{local}}$  proceeds directly. The condition on the factor  $s$  is described below ( $s = 5$  is a typical choice). It remains to describe an algorithm for evaluating  $\mathbf{f}_k^+$  since  $\mathbf{f}_k^-$  is computed in a similar manner.

We have

$$\begin{aligned}
\mathbf{f}_k^+ &= \sum_{\mathbf{x}_k - \mathbf{t}_j} \frac{1}{s} \left( \frac{\mathbf{u}_j + i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j - i\mathbf{s}_j} + \frac{\mathbf{u}_j - i\mathbf{v}_j}{\mathbf{x}_k - \mathbf{t}_j + i\mathbf{s}_j} \right) \\
(17) \quad &= 2 \sum_{\mathbf{x}_k - \mathbf{t}_j} \frac{1}{s} \int_{-\infty}^{\infty} e^{-e^y(\mathbf{x}_k - \mathbf{t}_j) + y} (\mathbf{u}_j \cos(e^y \mathbf{s}_j) - \mathbf{v}_j \sin(e^y \mathbf{s}_j)) dy.
\end{aligned}$$

The effective range of integration in (17) is finite due to the exponential ( $e^{-e^y(\mathbf{x}_k - \mathbf{t}_j)}$ ) and super-exponential ( $e^y$ ) decay of the integrand. Our choice of the factor  $s$  prevents an excessive oscillatory behavior of the integrand within that range. In order to obtain an exponential approximation of the form

$$(18) \quad \mathbf{f}_k^+ = \sum_{\mathbf{t}_j} \sum_{\mathbf{x}_k} \sum_{l=1}^L \mathbf{l}_j e^{-\mu_l(\mathbf{x}_k - \mathbf{t}_j)}, \quad s \leq \mathbf{x}_k - \mathbf{t}_j \leq \mathbf{T}, \quad \mathcal{R}e(\mu_l) > 0,$$

(where  $\mathbf{T}$  is sufficiently large to accommodate a given segment of the signal), we may now proceed as in [5, 7]. Indeed, we discretize the integral in (17) to any desired precision and use an appropriate algorithm to reduce the number of terms.

In (18) we may switch the order of summation and, as a result, construct a recursion (see [13, 5]). Denoting

$$\mathbf{q}_{k,l} = \sum_{\mathbf{t}_j} \sum_{\mathbf{x}_k} \mathbf{l}_j e^{-\mu_l(\mathbf{x}_k - \mathbf{t}_j)},$$

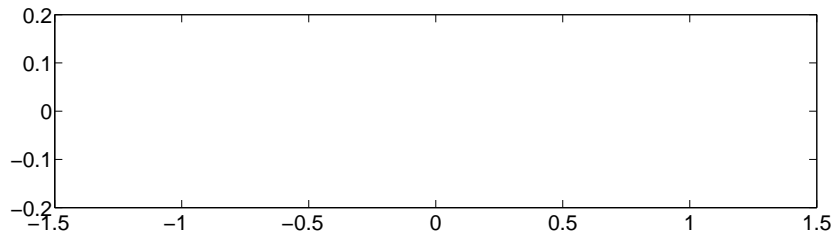
we obtain

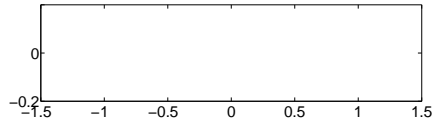
$$\begin{aligned}
\mathbf{q}_{k+1,l} &= \sum_{\mathbf{t}_j} \sum_{\mathbf{x}_{k+1}} \mathbf{l}_j e^{-\mu_l(\mathbf{x}_{k+1} - \mathbf{t}_j)} \\
&= e^{-\mu_l(\mathbf{x}_{k+1} - \mathbf{x}_k)} \mathbf{q}_{k,l} + \sum_{\mathbf{x}_k < \mathbf{t}_j} \sum_{\mathbf{x}_{k+1}} \mathbf{l}_j e^{-\mu_l(\mathbf{x}_{k+1} - \mathbf{t}_j)}.
\end{aligned}$$

This recursion leads to an  $\mathcal{O}(\mathbf{L} \cdot \mathbf{K}) + \mathcal{O}(\mathbf{L} \cdot \mathbf{M})$  algorithm for computing  $\mathbf{f}_k^+$ .

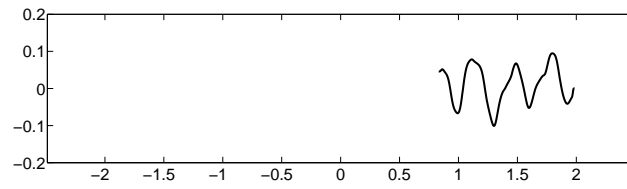
## 5. NUMERICAL EXAMPLES

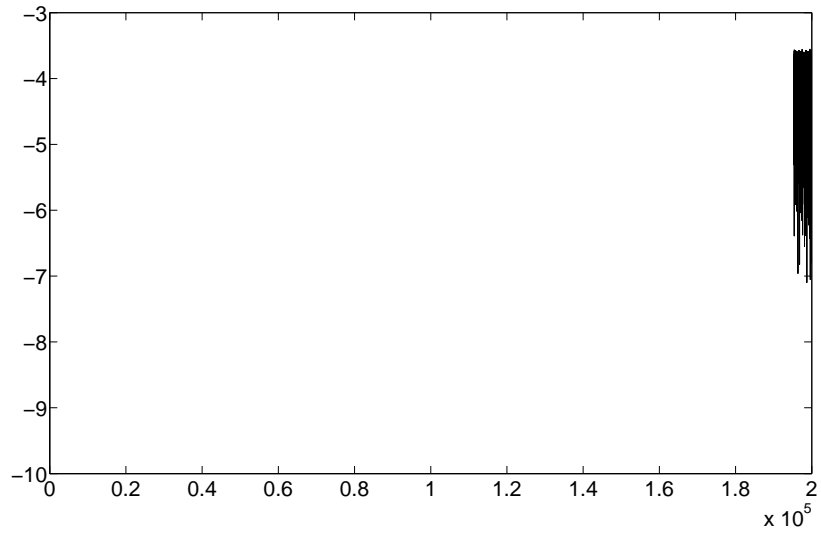
We have computed several approximations using the algorithm from Section 4. Since one of the potential applications for this method is a compression scheme, we illustrate our algorithm using a large data set from a high quality audio recording.











- [12] B. A. Jones, G. H. Born, and G. Beylkin. A Cubed Sphere Gravity Model for Fast Orbit Propagation. *AAS/AIAA Spaceflight Mechanics Meeting, Advances in the Astronautical Sciences*, 134:567–584, 2009.
- [13] N. Yarvin and V. Rokhlin. An improved fast multipole algorithm for potential fields on the line. *SIAM J. Numer. Anal.*, 36(2):629–666 (electronic), 1999.

DEPARTMENT OF APPLIED MATHEMATICS, UNIVERSITY OF COLORADO, BOULDER,  
CO 80309-0526