



Padé-Jacobi Approximants

Jan Hesthaven, Sidi Mahmoud Kaber

► To cite this version:

| Jan Hesthaven, Sidi Mahmoud Kaber. Padé-Jacobi Approximants. 2016. hal-01418450

HAL Id: hal-01418450

<https://hal.sorbonne-universite.fr/hal-01418450>

Preprint submitted on 16 Dec 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Padé-Jacobi Approximants

J.S. Hesthaven[†] and S.M. Kaber[‡]

December 14, 2016

Abstract

We discuss projection based Padé-Jacobi approximants in general and present in particular an exact rational approximation to the Sign function. This serves as vehicle to analyze the behavior of Padé-Jacobi approximants for discontinuous functions. The analysis shows that the Padé-Jacobi approximant is superior in several ways to classic polynomial approximations of discontinuous functions, provided the parameters in the approximations are chosen carefully. Guidelines for this is obtained through the analysis.

1 Introduction

The nonuniform pointwise convergence, known as the Gibbs phenomenon, of polynomial approximations to discontinuous function is a well known and much studied phenomenon, see e.g. [10] and references therein. Among the consequences of the Gibbs phenomenon is the lack of convergence at the jump with an overshoot of approximately 9% of the jump size, a global $\mathcal{O}(N^{-1})$ convergence rate in mean, and a steepness of the approximation right at the jump being proportional to the length, N , of the polynomial expansion.

The literature is rich with methods trying to reduce or even eliminate these problems. The perhaps simplest approach is that of modal filtering, essentially relying on forcing the expansion to converge more rapidly [20, 10, 13]. An alternative approach is physical space filtering using mollifiers [11, 19], yielding similar behavior. Both methods, however, do not overcome the lack of convergence at the point of discontinuity. To achieve this, information about the shock location is needed. With this, the Gibbs phenomenon can be completely resolved [10], albeit this approach has considerable practical problems.

In this work we shall discuss the use of rational functions, Padé-Jacobi approximants, for the representation of discontinuous functions. As rational functions are richer than simple polynomial expansions, one can hope that the impact of the discontinuity will be less severe and, further, that one could use this as a postprocessing tool to reduce the impact of the Gibbs phenomena in polynomial expansions.

To study the fundamental behavior of Padé-Jacobi approximants of discontinuous functions, we present a family of exact rational approximations to the Sign function, considered as a prototype of discontinuous

functions. This enables a complete analysis of the behavior of this approximation as characterized by the maximum size of the overshoot and the achievable steepness at the point of discontinuity. As we shall show, the use of a rational approximation allows one to dramatically reduce the overshoot and increase the steepness while recovering high order accuracy away from the jump.

There has been some recent activity in the exploration of Padé-forms for the reconstruction of Gibbs oscillations. In particular, work for the Fourier case can be found in [8, 6, 4], for the Chebyshev in [16], and for the Legendre approximations in [5, 12]. However, much of this has been of a qualitative character and for special polynomial families only.

In Sec. 2, we recall some properties of the Jacobi polynomials and the Padé-Jacobi problem. Section 3 is devoted to the derivation of an exact solution of the Padé-Jacobi approximation problem for the Sign function. In Sec. 4, we consider the optimization of the Padé-Jacobi solution by varying several of the free parameters. Section 5 contains a few remarks.

2 Jacobi Polynomials and Padé-Jacobi Approximations

In the following we shall recall various definitions and properties of Jacobi polynomials and expansions, as well as define exactly what we mean by Padé-Jacobi approximations in this work.

2.1 Jacobi Polynomials and Expansions

For $\alpha > -1$, the symmetric Jacobi polynomials, $P_n^{(\alpha)}(x)$, also known as the ultraspherical polynomials, are defined as the polynomial eigenfunctions to the singular Sturm-Liouville problem

$$\mathcal{A}_\alpha P_n^{(\alpha)}(x) = \lambda_n^\alpha P_n^{(\alpha)}(x) \quad , \quad x \in [-1, 1] \quad , \quad (1)$$

where

$$\mathcal{A}_\alpha \varphi = -\frac{1}{\omega_\alpha} (\omega_{\alpha+1} \varphi')' \quad ,$$

with the weight function

$$\omega_\alpha = (1 - x^2)^\alpha \quad ,$$

and the eigenvalue

$$\lambda_n^\alpha = n(n + 2\alpha + 1) \quad .$$

One easily proves that the Jacobi polynomials are the unique polynomial solution [18] to Eq.(1), once a normalization is chosen. The standard choice, also used here, is

$$P_n^{(\alpha)}(1) = \frac{\Gamma(n + \alpha + 1)}{\Gamma(n + 1)\Gamma(\alpha + 1)} \quad ,$$

where $\Gamma(x)$, $x \geq 0$ represents the classic Euler Gamma function. Recall that $\Gamma(n + 1) = n\Gamma(n) = n!$.

We introduce the Pochhammer symbol

$$(z)_n = \frac{\Gamma(z+n)}{\Gamma(z)} = (z+1)(z+2)\dots(z+n-1) ,$$

and note that $(1)_n = n!$ and $(z)_1 = z$. Recall also that for $k \in \mathbb{N}$

$$\forall n \geq k+1 : (-k)_n = 0 . \quad (2)$$

An important property of the Pochhammer symbol is expressed in the Saalchütz's formula [7]

$${}_3F_2(-n, a, b; d, 1+a+b-d-n; 1) = \frac{(d-a)_n(d-b)_n}{(d)_n(d-a-b)_n}, \quad (3)$$

with ${}_3F_2(-n, a, b; d, 1+a+b-d-n; 1)$ being the hypergeometric function defined as

$${}_3F_2(a, b, c; d, e; z) := \sum_{k=0}^{\infty} \frac{(a)_k(b)_k(c)_k}{(d)_k(e)_k} \frac{1}{k!} z^k .$$

Well known examples of ultraspherical polynomials are the Chebyshev polynomials ($\alpha = -1/2$) and the Legendre polynomials ($\alpha = 0$). The ultraspherical polynomials have a number of important properties which we shall exploit. In particular, all the polynomials are mutually orthogonal in the inner product

$$(u, v)_\alpha = \int_{-1}^1 u(x)v(x)\omega_\alpha dx, \quad (4)$$

with the associated weighted L_α^2 norm

$$\|u\|_\alpha^2 = (u, u)_\alpha .$$

The normalization is given by

$$\gamma_n^\alpha = \left(P_n^{(\alpha)}, P_n^{(\alpha)} \right)_\alpha = \frac{2^{2\alpha+1}}{2n+2\alpha+1} \frac{\Gamma(n+\alpha+1)^2}{\Gamma(n+1)\Gamma(n+2\alpha+1)} . \quad (5)$$

Another important property of the ultraspherical polynomials is their even-odd characteristics

$$P_n^{(\alpha)}(x) = (-1)^n P_n^{(\alpha)}(-x) . \quad (6)$$

We also recall the special value at $x = 0$ as

$$P_{2n}^{(\alpha)}(0) = (-1)^n 2^{-2n} \binom{2n+\alpha}{n} \quad (7)$$

and zero otherwise due to Eq.(6).

Finally we shall need the relation [18]

$$\frac{d}{dx} P_n^{(\alpha)}(x) = \frac{1}{2} (n+2\alpha+1) P_{n-1}^{(\alpha+1)}(x) . \quad (8)$$

If we now consider functions, $u(x) \in L_\alpha^2$, i.e., for which $\|u\|_\alpha < \infty$, we can seek polynomial approximations as

$$u(x) = \sum_{n=0}^{\infty} \hat{u}_n P_n^{(\alpha)}(x) \quad , \quad \hat{u}_n = \frac{1}{\gamma_n^\alpha} \left(u, P_n^{(\alpha)} \right)_\alpha \quad ,$$

by orthogonality.

Let us consider the truncated expansion

$$u_N^{(\alpha)}(x) = \sum_{n=0}^N \hat{u}_n P_n^{(\alpha)}(x), \quad (9)$$

i.e., $u_N^{(\alpha)} \in \mathbf{P}_N$ where \mathbf{P}_N is the space of algebraic polynomials of degree less than or equal to N . The orthogonality of the Jacobi polynomials implies

$$\forall p \in \mathbf{P}_N \quad : \quad \left(u - u_N^{(\alpha)}, p \right)_\alpha = 0.$$

It is well known that the polynomial expansion is convergent in the mean but not uniformly. In particular, if the smoothness is measured in the Sobolev space H_α^p of functions u and their derivatives up to order p in L_α^2 , there exists a constant c_p such that

$$\|u - u_N^{(\alpha)}\|_\alpha \leq c_p N^{-p} \|u\|_{H_\alpha^p}.$$

For a smooth function, i.e., p large, this provides an accurate approximation and the approximation error decreases rapidly to zero as N goes to infinity. This is one of the main motivations for using spectral methods for solving partial differential equations with regular solutions. We refer the reader to [1].

However, for problems with discontinuous solutions, the expansion exhibits non-uniform convergence and a phenomenon known as the Gibbs phenomenon [14] as illustrated in Fig. 1 where the truncated ($N = 20$ and $N = 100$) Legendre expansion, i.e., for $\alpha = 0$, of the Sign function is displayed. One observes the oscillations, especially near the discontinuity. If the parameter N is increasing, the size of the oscillations decrease everywhere except near the discontinuity where $\mathcal{O}(1)$ oscillations (overshoot/undershoot) remain. Furthermore, the global nature of the oscillations limits the pointwise accuracy to first order away from $x = 0$.

One of the objectives of this work is to consider Padé-Jacobi approximations of the Sign function and attempt to answer the question of which one among this family would be best suited to approximate the Sign function. As measures of success we shall consider

- the overshoot/undershoot of the approximation at the point of discontinuity as characterized by the Gibbs constant of the expansion.
- the ability to reproduce the discontinuity characterized by the steepness of the approximation.

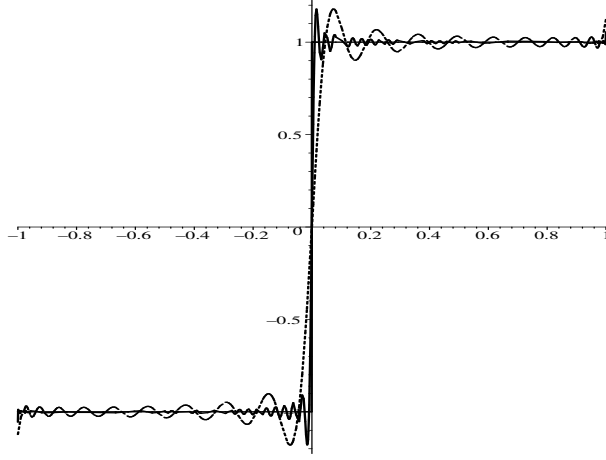


Figure 1: Legendre approximations of the Sign function: $N = 20$ (dashed) and $N = 100$ (solid).

2.2 Padé-Jacobi Approximants

We shall consider the Padé-Jacobi approximation to u in a Galerkin sense, i.e., find $\mathcal{P} \in \mathbf{P}_N$ and $\mathcal{Q} \in \mathbf{P}_M$ such that

$$(\mathcal{Q}u - \mathcal{P}, p)_\alpha = 0, \quad \forall p \in \mathbf{P}_K, \quad (10)$$

with $K \leq M + N$. This shall be used to define the linear Galerkin type Padé-Jacobi approximation of order (N, M) to u as the rational function

$$\mathcal{R}_{N,M}(x) = \frac{\mathcal{P}(x)}{\mathcal{Q}(x)},$$

where $(\mathcal{P}, \mathcal{Q})$ satisfies Eq.(10).

Remark 2.1 For $M = 0$ and $K = N$, the pair $(\mathcal{Q} \equiv 1, \mathcal{P} = u_N^{(\alpha)})$ defined in (9) is a solution of the Padé-Jacobi approximation problem.

It is important for practical purposes that the complete knowledge of u is not needed to solve the problem (10), only $u_{N+M}^{(\alpha)}$ is required. In Eq. (10) we take $p = P_k^{(\alpha)}$ with $k = N + 1, \dots, K$, to get

$$(\widehat{\mathcal{Q}u})_k^{(\alpha)} = 0, \quad \forall k = N + 1, \dots, K.$$

This is a linear system of $K - N$ equations and $M + 1$ unknowns (the coefficients of \mathcal{Q} in a basis of \mathbf{P}_M). Once a non trivial solution of this system is found (such a solution always exists if $K \leq N + M$), the numerator $\mathcal{P} \in \mathbf{P}_N$ is simply computed by

$$\widehat{\mathcal{P}}_k^{(\alpha)} = (\widehat{\mathcal{Q}u})_k^{(\alpha)}, \quad \forall k = 0, \dots, N.$$

Hence the main problem is the determination of the denominator.

3 Analysis of the Sign-Function

Let us now return to a more thorough analysis of the behavior of the Padé-Jacobi approximation for the most basic discontinuous function – the Sign function

$$u(x) = \begin{cases} -1 & x < 0 \\ 1 & x > 0 \end{cases} . \quad (11)$$

3.1 An Exact Solution

For the purpose of analysis, we shall seek the approximation to the Sign function using

$$\mathcal{P}(x) = \sum_{n=0}^N \hat{p}_{2n+1}^\alpha P_{2n+1}^{(\alpha)}(x) \in \mathbf{P}_{2N+1} ,$$

and

$$\mathcal{Q}(x) = \sum_{m=0}^M \hat{q}_{2m}^\alpha x^{2m} \in \mathbf{P}_{2M} ,$$

where we have used the parity of the problem and Eq.(6) to reduce the complexity of the problem.

We must now seek $\mathcal{P} \in \mathbf{P}_{2N+1}$ and $\mathcal{Q} \in \mathbf{P}_{2M}$ to satisfy Eq.(10), i.e.,

$$\left(\mathcal{Q}u - \mathcal{P}, P_{2k+1}^{(\alpha)} \right)_\alpha = 0 , \quad \forall k \leq K ,$$

where we have again utilized the parity of the problem to reduce the complexity.

Orthogonality of the Jacobi basis immediately yields

$$\left(\mathcal{P}u, P_{2k+1}^{(\alpha)} \right)_\alpha = \left(\mathcal{Q}u, P_{2k+1}^{(\alpha)} \right)_\alpha = \gamma_{2k+1}^\alpha \widehat{(\mathcal{Q}u)}_{2k+1}^\alpha = 0 \quad N < k \leq K ,$$

which is a linear system of $(K - N) \times (M + 1)$, with the unknowns being the coefficients of the denominator \mathcal{Q} , i.e., \hat{q}_{2m}^α . Clearly for $K \leq N + M$ this linear system will always have at least one nontrivial solution. In what remains we shall restrict ourselves to the special case

$$K = N + M .$$

To compute the numerator, we observe that

$$\hat{p}_{2k+1}^\alpha - \widehat{(\mathcal{Q}u)}_{2k+1}^\alpha = 0 , \quad 0 \leq k \leq N , \quad (12)$$

i.e., once \mathcal{Q} , and hence, $\widehat{\mathcal{Q}u}_{2k+1}^\alpha$ is computed, the numerator follows immediately.

Let us thus focus on the computation of the denominator, or rather its coefficients, satisfying

$$\widehat{(\mathcal{Q}u)}_{2k+1}^\alpha = \sum_{m=0}^M \hat{q}_{2m}^\alpha \frac{1}{\gamma_{2k+1}^\alpha} \left(x^{2m} u, P_{2k+1}^{(\alpha)} \right)_\alpha = 0 .$$

We shall need the following results. Define for the integers k and m

$$\mathcal{I}_{m,k}^\alpha = \int_0^1 x^{2m} P_{2k+1}^{(\alpha)} \omega_\alpha dx \quad \text{and} \quad \mathcal{J}_{m,k}^\alpha = \int_0^1 x^{2m+1} P_{2k+1}^{(\alpha)} \omega_\alpha dx, \quad (13)$$

Lemma 3.1 For $m \geq 1$ and $k \geq 0$

$$\mathcal{I}_{m,k}^\alpha = \mathcal{I}_{0,k}^\alpha \frac{m!(1/2)_m}{(-k+1/2)_m(k+\alpha+2)_m}$$

with

$$\mathcal{I}_{0,k}^\alpha = \frac{k+\alpha+1}{\lambda_{2k+1}^\alpha} P_{2k}^{(\alpha+1)}(0) = \frac{(-1)^k}{2k+1} \frac{1}{2^{2k+1}} \frac{1}{k!} \frac{\Gamma(2k+\alpha+2)}{\Gamma(k+\alpha+2)}.$$

For $m < k$, $\mathcal{J}_{m,k}^\alpha = 0$ and

$$\forall j \in \mathbb{N} \quad : \quad \mathcal{J}_{k+j,k}^\alpha = \frac{(k+1)_j}{j!} \frac{(k+3/2)_j}{(2k+\alpha+5/2)_j} \mathcal{J}_{k,k}^\alpha \quad (14)$$

with

$$\mathcal{J}_{k,k}^\alpha = \frac{1}{2} \frac{\gamma_{2k+1}^{(\alpha)}}{\theta_{2k+1}^{(\alpha)}}$$

and θ_j^α the coefficient of x^j in $P_j^{(\alpha)}(x)$.

Proof: Using Eq.(1) we have

$$\lambda_{2k+1}^\alpha \mathcal{I}_{m,k}^\alpha = \int_0^1 x^{2m} \mathcal{A}_\alpha P_{2k+1}^{(\alpha)} \omega_\alpha dx = - \int_0^1 x^{2m} \left(\omega_{\alpha+1} \left(P_{2k+1}^{(\alpha)} \right)' \right)' dx.$$

Recalling the singular nature of ω_α , integration by parts twice yields

$$\begin{aligned} \lambda_{2k+1}^\alpha \mathcal{I}_{m,k}^\alpha &= 2m \int_0^1 x^{2m-1} \omega_{\alpha+1} \left(P_{2k+1}^{(\alpha)} \right)'(x) dx \\ &= -2m \int_0^1 (\omega_{\alpha+1} x^{2m-1})' P_{2k+1}^{(\alpha)} dx \\ &= -2m(2m-1) \mathcal{I}_{m-1,k}^\alpha + 2m(2m+2\alpha+1) \mathcal{I}_{m,k}^\alpha. \end{aligned}$$

From this, we recover the recurrence

$$\mathcal{I}_{m,k}^\alpha = \frac{-2m(2m-1)}{\lambda_{2k+1}^\alpha - 2m(2m+2\alpha+1)} \mathcal{I}_{m-1,k}^\alpha = \frac{-m(2m-1)}{(2k-2m+1)(k+m+\alpha+1)} \mathcal{I}_{m-1,k}^\alpha.$$

We finally note that

$$\begin{aligned} \lambda_{2k+1}^\alpha \mathcal{I}_{0,k}^\alpha &= - \int_0^1 \left(\omega_{\alpha+1} \left(P_{2k+1}^{(\alpha)} \right)' \right)' dx = - \left[\omega_{\alpha+1} \left(P_{2k+1}^{(\alpha)} \right)' \right]_0^1 \\ &= \left(P_{2k+1}^{(\alpha)} \right)'(0) = (k+\alpha+1) P_{2k}^{\alpha+1}(0). \end{aligned}$$

Combining Eq.(7) with the above result yields

$$\begin{aligned} \mathcal{I}_{m,k}^\alpha &= \frac{m(m-1/2)}{(-k+m-1/2)(k+m+\alpha+1)} \mathcal{I}_{m-1,k}^\alpha \\ &= \mathcal{I}_{0,k}^\alpha \frac{m!(1/2)_m}{(-k+1/2)_m(k+\alpha+2)_m}. \end{aligned}$$

Concerning $\mathcal{J}_{m,k}^\alpha$, we observe that $P_j^{(\alpha)}(x) = \theta_j^{(\alpha)} x^j + q_j^{(\alpha)}$ with $q_j^{(\alpha)} \in \mathbf{P}_{j-1}$:

$$\mathcal{J}_{k,k}^\alpha = \frac{1}{2} \int_{-1}^1 \frac{P_{2k+1}^{(\alpha)} - q_{2k+1}^{(\alpha)}}{\theta_{2k+1}^{(\alpha)}} P_{2k+1}^{(\alpha)} \omega_\alpha dx = \frac{1}{2} \frac{\gamma_{2k+1}^{(\alpha)}}{\theta_{2k+1}^{(\alpha)}} , \quad \theta_j^\alpha = \frac{1}{2^j} \frac{1}{j!} \frac{\Gamma(2j+2\alpha+1)}{\Gamma(j+2\alpha+1)} .$$

The proof of (14) follows the same lines. \square Thus, to find a Padé-Jacobi approximant to the Sign function, we must seek a solution to

$$\sum_{m=0}^M \hat{q}_{2m}^\alpha \frac{m!(1/2)_m}{(-k+1/2)_m (k+\alpha+2)_m} = 0. \quad (15)$$

One non-unique solution is given in the following

Proposition 3.2 *The coefficients, \hat{q}_{2m}^α , defined for $m \in [0, M]$ as*

$$\hat{q}_{2m}^\alpha = \frac{(-M)_m (A)_m (-A+M+\alpha+3/2)_m}{(m!)^2 (1/2)_m},$$

is a solution to Eq.(15) with

$$A = -(N+1/2) .$$

Proof: Inserting the above result into Eq.(15) yields

$$\sum_{m=0}^M \frac{(-M)_m (A)_m (-A+M+\alpha+3/2)_m}{m! (-k+1/2)_m (k+\alpha+2)_m} .$$

Using Eq.(3) this can be written as

$$\frac{(k+\alpha+2-A)_M (k-M-1/2+A)_M}{(k+\alpha+2)_M (k-M-1/2)_M} .$$

Recalling Eq.(2) we immediately get two solutions to Eq.(15) from each of the two terms in the numerator

$$A = M+N+\alpha+2 , \quad A = -N-1/2 .$$

In both cases, we have $(k-M-N)_M$ which vanishes for all $k \in (N+1, M+N)$. \square Thus, the denominator takes the form

$$\begin{aligned} \mathcal{Q}(x) &= \sum_{m=0}^M \frac{(-M)_m (-N-1/2)_m (N+M+\alpha+2)_m}{(m!)^2 (1/2)_m} x^{2m} \\ &= {}_3F_2(-M, -N-1/2, N+M+\alpha+2; 1, 1/2; x^2) . \end{aligned} \quad (16)$$

Before we continue with the development of the Padé-Jacobi approximation, let us consider a few properties of $Q(x)$.

Lemma 3.3 *Provided $N+3/2 > M$ we have*

$$\forall m \in [0, M] : \hat{q}_{2m} > 0,$$

and, hence,

$$Q(x) \geq Q(0) = \hat{q}_0 = 1 .$$

Proof: Consider

$$\hat{q}_{2m} = \frac{(-M)_m(-N-1/2)_m(N+M+\alpha+2)_m}{(m!)^2(1/2)_m} .$$

Clearly,

$$(-M)_m = -M(-M+1)(-M+2)\dots(-M+m-1) = (-1)^m \frac{M!}{(M-m)!} ,$$

and

$$\begin{aligned} (-N-1/2)_m &= (-N-1/2)(-N+1/2)(-N+3/2)\dots(-N-3/2+m) \\ &= (-1)^m \frac{\Gamma(N+3/2)}{\Gamma(N+3/2-m)} . \end{aligned}$$

As $m \in [0, M]$, $(-M)_m(-N-1/2)_m > 0$ provided only that $N+3/2-M > 0$, hence completing the proof. \square Thus all

roots of $Q(x)$ are complex, ensuring that the Padé-Jacobi approximation to the Sign-function always exists.

The location of the roots can be specified a bit more

Lemma 3.4 *Assume that $N \gg M$, α fixed, and $z \in \mathbb{C}$ be a root of $Q(x)$. Then*

$$\frac{1}{2M} \leq N^2|z|^2 \leq M^3 .$$

For a proof of this Lemma, see Proposition 4.7 of [16]. As we shall see shortly, this results also gives some indications of how well one can expect the to approximate the Sign function since there is a direct relation between the position of the poles and the ability of the approximation to reproduce the discontinuity.

The sharpness in N can be realized by considering the limit of large N in which case

$$(-N-1/2)_m(N+M+\alpha+2)_m \simeq (-1)^m N^{2m} ,$$

such that

$$\mathcal{Q}(x) = {}_3F_2(-M, -N-1/2, N+M+\alpha+2; 1, 1/2; x^2) \simeq {}_1F_2(-M; 1, 1/2; -N^2x^2) .$$

However, since the ${}_1F_2(a1; b1, b2; z)$ is independent of N , the roots of $\mathcal{Q}(x)$ can not decay faster than N^{-1} . It is worth emphasizing that this result assumes that M is fixed, i.e., making $M \propto N$ and/or $\alpha \propto N$ may yield qualitative differences in the approximation as we shall indeed see shortly.

Let us now return to the determination of the numerator,

$$P(x) = \sum_{n=0}^N \hat{p}_{2n+1} P_{2n+1}^{(\alpha)}(x). \quad (17)$$

The coefficients of this polynomial are given in the following

Lemma 3.5 *The coefficients \hat{p}_{2n+1}^α in (17) are defined for $n \in [0, N]$ as*

$$\hat{p}_{2n+1}^\alpha = 2 \frac{\mathcal{I}_{0,n}^\alpha}{\gamma_{2n+1}^\alpha} \frac{(N+n+\alpha+5/2)_M (n-N-M)_M}{(n+\alpha+2)_M (n+1/2-M)_M},$$

where $\mathcal{I}_{0,n}^\alpha$ is given in Lemma 3.1.

Proof: From the definition of the Padé-Jacobi approximation and the orthogonality of the Jacobi polynomials, we immediately recover from (12)

$$\begin{aligned} \gamma_{2n+1}^\alpha \hat{p}_{2n+1}^\alpha &= 2 \sum_{m=0}^M \hat{q}_{2m} \mathcal{I}_m^\alpha \\ &= 2 \mathcal{I}_{0,n}^\alpha \sum_{m=0}^M \frac{(-M)_m (-N-1/2)_m (N+M+\alpha+2)_m}{m! (1/2-n)_m (n+\alpha+2)_m} \\ &= 2 \mathcal{I}_{0,n}^\alpha \frac{(N+n+\alpha+5/2)_M (n-N-M)_M}{(n+\alpha+2)_M (n+1/2-M)_M}, \end{aligned}$$

where the last reduction follows from the Saalchütz's formula (3). \square

Using the identity $(n-N-M)_M = (-1)^M (N-n+M)!/(N-n)!$, one can express the coefficients \hat{p}_{2n+1}^α in the form

$$\hat{p}_{2n+1}^\alpha = 2(-1)^M \frac{\mathcal{I}_{0,n}^\alpha}{\gamma_{2n+1}^\alpha} \frac{(N+n+\alpha+5/2)_M (N-n+M)!}{(n+\alpha+2)_M (n+1/2-M)_M (N-n)!}. \quad (18)$$

In (16), the denominator was written in a geometric form. Now we seek the numerator in the form : $\mathcal{P}(x) = x S_{N,M}^\alpha {}_3F_2(\cdot, \cdot, \cdot; \cdot, \cdot; x^2)$, with $S_{N,M}^\alpha$ being the steepness, i.e., the value at $x = 0$ of the derivative of the rational approximation $\mathcal{R}(x) = \mathcal{P}(x)/\mathcal{Q}(x)$. The following Proposition allows us to recover the numerator on a hypergeometric form.

Proposition 3.6 *The numerator, $\mathcal{P}(x)$, for the Padé-Jacobi approximation to the Sign function, Eq.(11), takes the form*

$$\mathcal{P}(x) = x S_{N,M}^\alpha {}_3F_2(-N, -M+1/2, N+M+\alpha+5/2; 3/2, 3/2; x^2).$$

where the steepness, $S_{N,M}^\alpha$ is given as

$$S_{N,M}^\alpha = \frac{4}{\sqrt{\pi}} \frac{M!}{N!} \frac{\Gamma(N+3/2)\Gamma(N+M+\alpha+5/2)}{\Gamma(M+1/2)\Gamma(N+M+\alpha+2)}.$$

Proof: Let us compute the Jacobi coefficients of $\mathcal{T} \in \mathcal{P}_{2N+1}$ defined by

$$\mathcal{T}(x) = x S_{N,M}^\alpha {}_3F_2(-N, -M+1/2, N+M+\alpha+5/2; 3/2, 3/2; x^2).$$

For $n = 0 \dots, N$

$$\begin{aligned} \frac{\gamma_{2n+1}^\alpha}{S_{N,M}^\alpha} \hat{t}_{2n+1}^\alpha &= \int_{-1}^1 x {}_3F_2(-N, -M+1/2, N+M+\alpha+5/2; 3/2, 3/2; x^2) P_{2n+1}^{(\alpha)}(x) \omega_\alpha(x) dx \\ &= 2 \sum_{k=0}^N \frac{(-N)_k (-M+1/2)_k (N+M+\alpha+5/2)_k}{(3/2)_k (3/2)_k} \frac{1}{k!} \mathcal{J}_{k,n} \end{aligned}$$

with $\mathcal{J}_{k,n}^\alpha$ defined in (13). Using Lemma 3.1, we get

$$\frac{\gamma_{2n+1}^\alpha}{S_{N,M}^\alpha} \hat{t}_{2n+1}^\alpha = 2 \sum_{p=0}^{N-n} \frac{(-N)_{n+p} (-M+1/2)_{n+p} (N+M+\alpha+5/2)_{n+p}}{(3/2)_{n+p} (3/2)_{n+p} (1)_{n+p}} \mathcal{J}_{n+p,n}.$$

By the identity $(z)_{n+p} = (z)_n (z+n)_p$, we get

$$\hat{t}_{2n+1}^\alpha = 2X_{N,M,n}^\alpha \sum_{p=0}^{N-n} \frac{(-N+n)_p (-M+1/2+n)_p (N+M+\alpha+5/2+n)_p}{(3/2+n)_p (3/2+n)_p (1+n)_p} \mathcal{J}_{n+p,n}$$

with

$$X_{N,M,n}^\alpha = \frac{S_{N,M}^\alpha}{\gamma_{2n+1}^\alpha} \frac{(-N)_n (-M+1/2)_n (N+M+\alpha+5/2)_n}{(3/2)_n (3/2)_n (1)_n}.$$

Using (14), we get

$$\hat{t}_{2n+1}^\alpha = Y_{N,M,n}^\alpha \sum_{p=0}^{N-n} \frac{(-N+n)_p (-M+1/2+n)_p (N+M+\alpha+5/2+n)_p}{(3/2+n)_p (2n+\alpha+5/2)_p} \frac{1}{p!}$$

with $Y_{N,M,n}^\alpha = 2X_{N,M,n}^\alpha \mathcal{J}_{n,n}^\alpha$. By use of the Saalchütz's formula (3), we obtain

$$\begin{aligned} \hat{t}_{2n+1}^\alpha &= Y_{N,M,n}^\alpha \frac{(M+1)_{N-n} (-N-M-\alpha-1)_{N-n}}{(3/2+n)_{N-n} (-N-n-\alpha-3/2)_{N-n}} \\ &= Y_{N,M,n}^\alpha \frac{\Gamma(n+3/2)\Gamma(2n+\alpha+5/2)}{\Gamma(N+3/2)\Gamma(N+n+\alpha+5/2)} \frac{\Gamma(N+M-n+1)\Gamma(N+M+\alpha+2)}{\Gamma(M+1)\Gamma(n+M+\alpha+2)}. \end{aligned}$$

with

$$\begin{aligned} Y_{N,M,n}^\alpha &= \frac{\pi}{4\theta_{2n+1}^\alpha} \frac{(-N)_n (-M+1/2)_n (N+M+\alpha+5/2)_n}{n!\Gamma(3/2+n)^2} S_{N,M}^\alpha \\ &= \sqrt{\pi} \frac{(-1)^n}{\theta_{2n+1}^\alpha} \frac{\Gamma(N+3/2)}{n!\Gamma(n+3/2)^2 (N-n)!} \\ &= \times \frac{\Gamma(M+1)\Gamma(-M+1/2+n)}{\Gamma(M+1/2)\Gamma(-M+1/2)} \frac{\Gamma(N+M+n+\alpha+5/2)}{\Gamma(N+M+\alpha+2)} \\ &= \frac{1}{\sqrt{\pi}} \frac{(-1)^n}{\theta_{2n+1}^\alpha} \frac{\Gamma(N+3/2)}{n!\Gamma(n+3/2)^2 (N-n)!} \\ &\quad \times (-1)^M \frac{\Gamma(M+1)\Gamma(-M+1/2+n)}{\Gamma(N+M+\alpha+2)} \Gamma(N+M+n+\alpha+5/2). \end{aligned}$$

The last simplification follows from the reflection formula the reflection formula

$$\Gamma(z)\Gamma(1-z) = \frac{\pi}{\sin(\pi z)}.$$

Finally

$$\begin{aligned} \hat{t}_{2n+1}^\alpha &= \frac{1}{\sqrt{\pi}} \frac{(-1)^n}{\theta_{2n+1}^\alpha} \frac{\Gamma(2n+\alpha+5/2)}{n!\Gamma(n+3/2)(N-n)!\Gamma(N+n+\alpha+5/2)} \\ &\quad \times \frac{(-1)^M \Gamma(-M+1/2+n) \Gamma(N+M-n+1) \Gamma(N+M+n+\alpha+5/2)}{\Gamma(n+M+\alpha+2)} \end{aligned}$$

This is to be compared with the coefficients \hat{p}_{2n+1}^α given by Lemma 3.5 (see also (18)):

$$\begin{aligned}\hat{p}_{2n+1}^\alpha &= 2 \frac{\mathcal{I}_{0,n}^\alpha}{\gamma_{2n+1}^\alpha} \frac{\Gamma(n+\alpha+2)}{(N-n)!\Gamma(n+1/2)\Gamma(N+n+\alpha+5/2)} \\ &\quad \times (-1)^M \frac{\Gamma(N+M+n+\alpha+5/2)(N+M-n)!\Gamma(n+1/2-M)}{\Gamma(n+M+\alpha+2)}\end{aligned}$$

The ratio of the two coefficients is

$$\begin{aligned}\frac{\hat{i}_{2n+1}^\alpha}{\hat{p}_{2n+1}^\alpha} &= \frac{1}{2\sqrt{\pi}} \frac{(-1)^n \gamma_{2n+1}^\alpha}{\mathcal{I}_{0,n}^\alpha \theta_{2n+1}^\alpha} \frac{\Gamma(2n+\alpha+5/2)}{n!\Gamma(n+3/2)} \frac{\Gamma(n+1/2)}{\Gamma(n+\alpha+2)} \\ &= \frac{1}{2\sqrt{\pi}} \frac{(-1)^n \gamma_{2n+1}^\alpha}{\mathcal{I}_{0,n}^\alpha \theta_{2n+1}^\alpha} \frac{\Gamma(2n+\alpha+5/2)}{n!(n+1/2)\Gamma(n+\alpha+2)}.\end{aligned}$$

Straightforward computations give

$$\frac{\gamma_{2n+1}^\alpha}{\mathcal{I}_{0,n}^\alpha \theta_{2n+1}^\alpha} = (-1)^n 2^{4n+2\alpha+3} n!(2n+1)\Gamma(n+\alpha+2) \frac{\Gamma(2n+\alpha+2)}{\Gamma(4n+2\alpha+4)}$$

By the Legendre duplication formula¹, we have

$$\frac{\Gamma(2n+\alpha+2)}{\Gamma(4n+2\alpha+4)} = \frac{\sqrt{\pi}}{2^{4n+2\alpha+3}} \frac{1}{\Gamma(2n+\alpha+5/2)}$$

and

$$\frac{\gamma_{2n+1}^\alpha}{\mathcal{I}_{0,n}^\alpha \theta_{2n+1}^\alpha} = (-1)^n \sqrt{\pi} n!(2n+1) \frac{\Gamma(n+\alpha+2)}{\Gamma(2n+\alpha+5/2)}.$$

Hence $\hat{i}_{2n+1}^\alpha = \hat{p}_{2n+1}^\alpha$, which means the equality of the two odd polynomials, \mathcal{T} and \mathcal{P} . \square Lemma 3.6 and Eq. (16) gives the main

result in the form of a hypergeometric representation of the Padé-Jacobi approximation of Eq.(11)

Theorem 3.7 *For all integers N and M ,*

$$\mathcal{R}_{N,M}^\alpha(x) = S_{N,M}^\alpha x \frac{{}_3F_2(-N, -M+1/2, N+M+\alpha+5/2; 3/2, 3/2; x^2)}{{}_3F_2(-M, -N-1/2, N+M+\alpha+2; 1, 1/2; x^2)},$$

with the steepness $S_{N,M}^\alpha$ defined in Proposition 3.6, is a Padé-Jacobi approximation of order (N, M) to the Sign function, Eq.(11).

Note in particular that for the special case of $M = 0$, this includes the Jacobi polynomial approximation of a step function and, thus, enables the general analysis of the Gibbs phenomenon for this case also.

Let us first consider two extremal cases: $M = 0$ (polynomial approximation) and $N = 0$ (reciprocal polynomial approximation).

1. Polynomial (Jacobi) approximation.

$$\mathcal{R}_{N,0}^\alpha = S_{N,0}^\alpha x {}_3F_2(-N, 1/2, N+\alpha+5/2; 3/2, 3/2; x^2)$$

$\mathcal{R}_{N,0}^\alpha$ is nothing but the orthogonal projection (with respect to the inner product (4)) of the Sign function onto \mathcal{P}_{2N+1} defined in (9).

¹Legendre duplication formula: $\sqrt{\pi}\Gamma(2z) = 2^{2z-1}\Gamma(z)\Gamma(z+1/2)$

It has been shown in [15] that for all Jacobi approximants $S_{N,0}^\alpha \simeq 4N/\pi$. We give here the precise value of the steepness

$$S_{N,0}^\alpha = \frac{4}{\pi} \frac{\Gamma(N+3/2)}{N!} \frac{\Gamma(N+\alpha+5/2)}{\Gamma(N+\alpha+2)} \simeq \frac{4}{\pi} N.$$

2. Reciprocal polynomial approximation.

$$\mathcal{R}_{0,M}^\alpha(x) = S_{0,M}^\alpha x \frac{1}{{}_3F_2(-M, -1/2, M+\alpha+2; 1, 1/2; x^2)}$$

with the steepness

$$S_{0,M}^\alpha = 2 \frac{M!}{\Gamma(M+1/2)} \frac{\Gamma(M+\alpha+5/2)}{\Gamma(M+\alpha+2)} \simeq 2M.$$

4 Optimized Approximations

Let us consider in a bit more detail the question of the Gibbs phenomena and attempt to understand whether it is better behaved in the Padé-Jacobi approximation as compared to the classical polynomial approximations above.

We shall in particular consider the questions of steepness of the Padé-Jacobi approximation and the size of the overshoot as measured by the Gibbs constant, $G_{N,M}^\alpha$.

We consider five cases

1. **Case 1.** The classic polynomial case, $M = 0$, with the Gibbs constant denoted $G_{\cdot,0}^\alpha$.
2. **Case 2.** The reciprocal polynomial case, $N = 0$, with the Gibbs constant denoted $G_{0,\cdot}^\alpha$.
3. **Case 3.** The case of M going to infinity with N , however constrained such that $M = cN^s$ and fixed $c > 0$ and $s > 0$. In this case we shall denote the Gibbs constant as $G_{\cdot,c,s}^\alpha$.
4. **Case 4.** The case of M fixed ($\neq 0$) with the Gibbs constant denoted $G_{\cdot,M}^\alpha$.
5. **Case 5.** The case of N fixed ($\neq 0$) with the Gibbs constant denoted $G_{N,\cdot}^\alpha$.

4.1 Optimize the steepness

If we first consider the steepness of the Padé-Jacobi approximation, then this is defined as

$$\frac{d}{dx} \mathcal{R}_{N,M}^\alpha(0) = \mathcal{P}'_N(0) = S_{N,M}^\alpha,$$

due to the symmetry of the problem, i.e., the steepness measures the ability to reproduce the discontinuity.

Cases 1, 4. From Lemma 3.6 we immediately get for a fixed M

$$\mathcal{P}'_N(0) = S_{N,M}^\alpha \simeq \frac{4}{\sqrt{\pi}} \frac{\Gamma(M+1)}{\Gamma(M+1/2)} \sqrt{N(N+M+\alpha)} \simeq \frac{4}{\sqrt{\pi}} \frac{\Gamma(M+1)}{\Gamma(M+1/2)} N,$$

for large N and fixed M and α . Hence, in this case there is no qualitative difference between the pure polynomial case ($M = 0$) and a fixed value of M . All polynomials behave, asymptotically, as the Legendre case of $\alpha = 0$.

Cases 2, 5. For a fixed N (and $M \rightarrow +\infty$), the steepness grows like

$$\frac{4}{\sqrt{\pi}} \frac{\Gamma(N + 3/2)}{N!} M.$$

Case 3. If $M \propto N$, we get a more interesting result. For $M = cN$, with positive constant c

$$S_{N,cN}^\alpha \simeq \frac{4}{\sqrt{\pi}} \sqrt{c(c+1)} N^{3/2}$$

which improve drastically the steepness in comparison with the polynomial case. It is important to note that $M = N$ is not needed to recover this improved steepness, simply that $M \propto N$. In the case $M = cN^s$, the steepness is

$$S_{N,cN^s}^\alpha \simeq \frac{4}{\sqrt{\pi}} c_s N^{s+1/2},$$

with $c_s = c$ except for $s = 1$: $c_1 = \sqrt{c(c+1)}$.

We remind the reader of Lemma 3.4 which reflects these steepness results in a slightly different way, i.e., for $s = 1$ we can not hope for better than $S_{N,cN}^\alpha \simeq N^{-3/2}$ as reflected in how quickly the poles approach zero.

4.2 Optimize the Gibbs constant

Let us now also consider the Gibbs constant, defined as the maximum overshoot of the Padé-Jacobi approximation. We proceed as in [17] for the Chebyshev approximation and seek an $\eta > 0$ such that the error function $x \in [0, 1] \mapsto \mathcal{R}_{N,M}(x) - 1$ takes its maximum at the point $x = \eta/N^\beta$ as N goes to infinity (β is a fixed real number to be made precise shortly). We shall call this limit the Gibbs constant, $G_{N,M}^\alpha$.

Case 1. ($M = 0$). We seek $\eta > 0$ such that the error function takes its maximum at the point $x = \eta/N$ as N goes to infinity. In this case, we have

$$\begin{aligned} 1 + G_{\cdot,0}^\alpha &= \lim_{N \rightarrow +\infty} \mathcal{R}_{N,0}\left(\frac{\eta}{N}\right) \\ &= \frac{4}{\pi} \eta \lim_{N \rightarrow +\infty} {}_3F_2(-N, 1/2, N + \alpha + 5/2; 3/2, 3/2; (\frac{\eta}{N})^2) \\ &= \frac{4}{\pi} \eta {}_1F_2(1/2; 3/2, 3/2; -\eta^2). \end{aligned}$$

η is determined as the smallest positive solution of

$$\frac{d}{d\eta} \left[\eta \sum_{k \geq 0} \frac{(1/2)_k}{(3/2)_k (3/2)_k} \frac{1}{k!} (-\eta^2)^k \right] = 0.$$

Using the identities $(1/2)_k = \frac{1}{2^{2k}} \frac{(2k)!}{k!}$ and $(3/2)_k = (2k+1)(1/2)_k$, we obtain

$$\begin{aligned} 0 &= \frac{d}{d\eta} \left[\eta \sum_{k \geq 0} \frac{2^{2k}}{(2k+1)^2} \frac{(-\eta^2)^k}{(2k)!} \right] = \frac{1}{2} \frac{d}{d\eta} \left[\sum_{k \geq 0} \frac{(-1)^k}{2k+1} \frac{(2\eta)^{2k+1}}{(2k+1)!} \right] \\ &= \frac{1}{2\eta} \sum_{k \geq 0} (-1)^k \frac{(2\eta)^{2k+1}}{(2k+1)!} = \frac{\sin(2\eta)}{2\eta}. \end{aligned}$$

Hence $\eta = \pi/2$ and the Gibbs constant is

$$G_{\cdot,0}^\alpha = -1 + \frac{4}{\pi} \int_0^{\pi/2} \frac{\sin(2t)}{2t} dt = -1 + \frac{2}{\pi} \text{Si}(\pi) \simeq 0.178\,979\,744$$

with $\text{Si}(z) = \int_0^z \frac{\sin s}{s} ds$ being the Sine integral. This is a classic result in Fourier and Chebyshev approximations (see e.g. [9]). It has also been shown for general Jacobi approximations in [15] using properties of orthogonal polynomials.

Case 2. (Case $N = 0$). The right scaling is $x = \eta/M$ as M goes to infinity. Using the same arguments as before, we obtain

$$1 + G_{0,\cdot}^\alpha = \lim_{M \rightarrow +\infty} \mathcal{R}_{M,0}\left(\frac{\eta}{M}\right) = 2\eta \frac{1}{{}_1F_2(-1/2; 1, 1/2; -\eta^2)}.$$

η is determined as the smallest positive solution of the equation

$$\frac{d}{d\eta} \left[\frac{1}{f(\eta)} \right] = 0, \quad f(\eta) = \frac{{}_1F_2(-1/2; 1, 1/2; -\eta^2)}{\eta}.$$

We observe that $f'(\eta) = J_0(2\eta)/\eta^2$ with $J_0(t)$ being the Bessel function of the first kind of order 0. Thus, 2η equals $j_{0,1}$, the first zero of J_0 , as

$$2\eta = j_{0,1} \simeq 2.404\,825\,557\,8,$$

and the Gibbs constant

$$G_{0,\cdot}^\alpha = -1 + j_{0,1} \frac{1}{{}_1F_2(-1/2; 1, 1/2; -(j_{0,1}/2)^2)} \simeq 0.051\,356\,067.$$

Here again all the Jacobi approximants give the same Gibbs constant, the one given in [17] for the Chebyshev case, $\alpha = -1/2$. An example of reciprocal polynomial approximation in Fig.2.

Case 3. ($M = cN^s$). Consider now the case $M = cN^s$ with fixed $c > 0$ and $s > 0$. In this case the steepness grows like $\frac{4}{\sqrt{\pi}} c_s N^{s+1/2}$ and we seek $\eta > 0$ that maximizes the error function $x = \eta/N^{s+1/2}$ as N goes to infinity. Defining $G_{N,c,s}^\alpha = G_{N,cN^s}^\alpha$ we have

$$\begin{aligned} 1 + G_{\cdot,c,s}^\alpha &= \frac{4}{\sqrt{\pi}} c_s \eta \\ &\lim_{N \rightarrow +\infty} \frac{{}_3F_2(-N, -cN + 1/2, (1+c)N + \alpha + 5/2; 3/2, 3/2; (\frac{\eta}{N^{s+1/2}})^2)}{{}_3F_2(-cM, -N - 1/2, (1+c)N + \alpha + 2; 1, 1/2; (\frac{\eta}{N^{s+1/2}})^2)} \\ &= f^*(c_s \eta), \end{aligned}$$

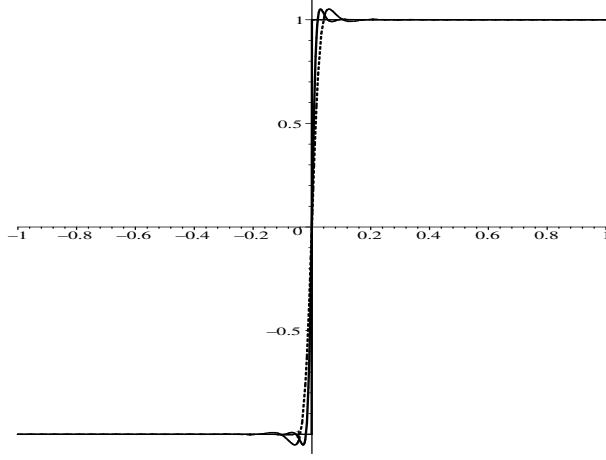


Figure 2: Padé-Jacobi approximations: $N = 0$, $\alpha = -1/2$, $M = 20$ (dashed), $M = 40$ (solid).

$$\text{with } f^*(z) = \frac{4}{\sqrt{\pi}} z \frac{{}_0F_2(; 3/2, 3/2; z^2)}{{}_0F_2(; 1, 1/2; z^2)}.$$

The unknown $\eta_{c,s}$ is determined as the smallest positive solution of the equation $\frac{d}{d\eta} [f^*(c_s \eta)] = 0$. This equation shows that $\eta_{c,s}$ does not depend on the Jacobi parameter α and

$$\eta_{c,s} = \frac{\eta^*}{c_s},$$

with η^* the first positive zero of f^* . The Gibbs constant $G_{\cdot, c, s}^\alpha = -1 + f^*(\eta^*)$ is independent of α , $c > 0$ and s .

Numerical experiments, finding the location of the first maximum of the analytic expression, yields the approximations

$$\eta_{\tilde{c}} \simeq 1.344\,947, \quad G_{\cdot, c, s}^\alpha \simeq 0.008\,149 \quad .$$

This is the value given in [17] for the special case $\alpha = 0$, $c = 1$ and $s = 1$. Examples of this type of approximation are shown in Fig.3 and Fig.4

Case 4. ($M \neq 0$). We seek $\eta > 0$ such that the error function takes its maximum at the point $x = \eta/N$ as N goes to infinity. In this case, we have

$$\begin{aligned} 1 + G_{\cdot, M}^\alpha &= \frac{4}{\sqrt{\pi}} \frac{\Gamma(M+1)}{\Gamma(M+1/2)} \eta \\ &\times \lim_{N \rightarrow +\infty} \frac{{}_3F_2(-N, -M+1/2, N+M+\alpha+5/2; 3/2, 3/2; (\frac{\eta}{N})^2)}{{}_3F_2(-M, -N-1/2, N+M+\alpha+2; 1, 1/2; (\frac{\eta}{N})^2)} \\ &= \frac{4}{\sqrt{\pi}} \frac{\Gamma(M+1)}{\Gamma(M+1/2)} \eta \frac{{}_1F_2(-M+1/2; 3/2, 3/2; -\eta^2)}{{}_1F_2(-M; 1, 1/2; -\eta^2)}. \end{aligned}$$

η_M (independent of α) is determined as the smallest positive solution of the equation $\frac{d}{dz} [f_M(z)] = 0$ with

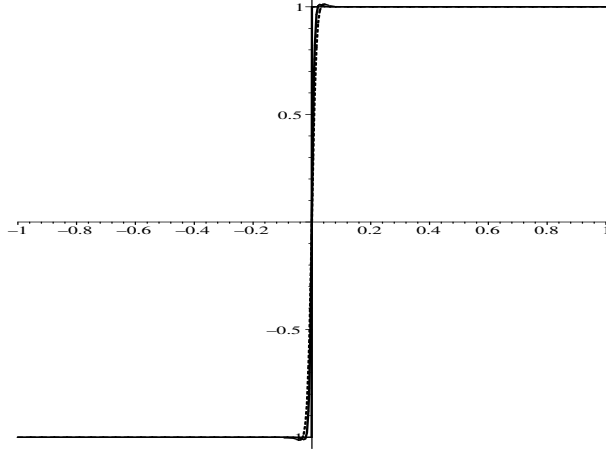


Figure 3: Padé-Jacobi approximations: $N = 10$, $\alpha = 0$, $M = 5$ (dashed), $M = 10$ (solid).

M	0	1	5	10	20	50	100	200
η_M	1.5708	1.0144	0.56	0.4098	0.2951	0.1888	0.134	0.095
$G_{.,M}^\alpha$	0.1789	0.0302	0.012	0.01	0.009	0.00851	0.0083	0.00823

Table 1: Computational evidence for the scaling of the Gibbs constant for values of fixed M

$$f_M(z) = \frac{4}{\sqrt{\pi}} \frac{\Gamma(M+1)}{\Gamma(M+1/2)} z \frac{{}_1F_2(-M+1/2; 3/2, 3/2; -z^2)}{{}_1F_2(-M; 1, 1/2; -z^2)}.$$

The Gibbs constant $G_{.,M}^\alpha = -1 + f_M(\eta_M)$ is likewise independent of α .

We have not been able to complete the analysis of this function. In Table 1 we show results for numerically finding, by seeking the position of the first maximum of the analytic expression, η_M and $G_{.,M}^\alpha$ for fixed values of M .

Based on these computations we conjecture that

$$\lim_{M \rightarrow +\infty} \sqrt{M} \eta_M = \eta^* = 1.344947$$

and

$$\lim_{M \rightarrow +\infty} G_{.,M}^\alpha = -1 + \lim_{M \rightarrow +\infty} f_M(\eta_M) = -1 + f^*(\eta^*) = G_{.,M}^\alpha \simeq 0.008149$$

The effect of changing N for fixed M can be seen by comparing Fig.3 and Fig.4.

Case 5. ($N \neq 0$). In this case the steepness grows like $\frac{4}{\sqrt{\pi}} \frac{\Gamma(N+3/2)}{N!} M$ and

$$1 + G_{N,.}^\alpha = \frac{4}{\sqrt{\pi}} \frac{\Gamma(N+3/2)}{N!} \eta \frac{{}_1F_2(-N; 3/2, 3/2; -\eta^2)}{{}_1F_2(-N-1/2; 1, 1/2; -\eta^2)}.$$

η_N (independent of α) is determined as the smallest positive solution of the equation $\frac{d}{dz} [g_N(z)] = 0$ with

$$g_N(z) = \frac{4}{\sqrt{\pi}} \frac{\Gamma(N+3/2)}{N!} z \frac{{}_1F_2(-N; 3/2, 3/2; -z^2)}{{}_1F_2(-N-1/2; 1, 1/2; -z^2)}.$$

The Gibbs constant $G_{N,.}^\alpha = -1 + g_N(\eta_N)$ is also independent of α and

$$\lim_{N \rightarrow +\infty} g_N(\eta/\sqrt{N}) = f^*(\eta^*) ,$$

based on computational experimentation as for Case 4.

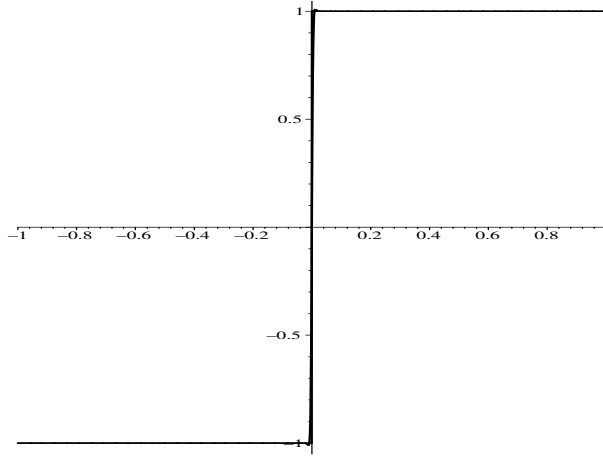


Figure 4: Padé-Jacobi approximations: $\alpha = 0$, $N = 30$, $M = 10$.

5 Concluding Remarks

We have derived an exact rational Galerkin approximation $\mathcal{R}_{N,M}^\alpha$ of the Sign function based on Jacobi expansions and investigated its ability to reduce the Gibbs phenomenon. The analysis contains the cases for M or N fixed, including the classic polynomial results, the case where both go to infinity with $M = cN^s$. The latter case is superior in terms of the Gibbs constant and steepness.

- The steepness of the approximation grows, as $N \rightarrow +\infty$, like $\frac{4}{\sqrt{\pi}} c_s N^{s+1/2}$ (in the case $c > 0$ and $s > 0$). Recall that in the polynomial case, the steepness is $\simeq \frac{4}{\pi} N$.
- The Gibbs constant is about 22 times less than that of a polynomial approximation.

One case we have not considered in detail is the one where α is a function of N and/or M . Although the analysis indicates that this could be interesting, it is less practical and unlikely to behave well numerically for high

values of α as also observed for Gibbs reconstruction methods based on Gegenbauer expansions [10, 2]. Nevertheless, we intend to consider this in more detail later. Furthermore, the rate of convergence of $\mathcal{R}_{N,M}^\alpha$ to the Sign function (the acceleration of convergence problem) remains open, yet is important to understand the value of Padé-Jacobi approximations for postprocessing. Some partial results on this can be found in [16] for the Padé-Chebyshev case and we hope to generalize these in the near future.

Acknowledgment

The work of JSH was partly supported by NSF Career Award DMS0132967, by an NSF International Award NSF-INT 0307475, and by the Alfred P. Sloan Foundation through a Sloan Research Fellowship.

References

- [1] C. BERNARDI, Y. MADAY, *Spectral methods, Handbook of numerical analysis*, Vol. V, North-Holland, Amsterdam, 1997.
- [2] J.P. BOYD, *Trouble with Gegenbauer reconstruction for defeating Gibbs phenomenon: Runge phenomenon in the diagonal limit of Gegenbauer polynomial approximations*, J. Comput. Phys. **204**(2005), pp. 253-264
- [3] C.W. CLENSHAW AND K. LORD, *Rational approximations from Chebyshev series*, in *Studies in Numerical Analysis*, B.K.P. Scaife ed., Academic Press, London, 1974, pp. 95-113.
- [4] , W.S. DON, S.M. KABER, M.S. SUN, *Fourier-Padé Approximations and Filtering for the Spectral Simulations of Incompressible Boussinesq Convection Problem*, accepted for publication in Math of comp.
- [5] L. EMMEL, S.M. KABER, Y. MADAY, *Padé-Jacobi filtering for spectral approximations of discontinuous solutions*. Numer. Algo. **33**(2003),pp. 251-264.
- [6] T. DRISCOLL AND B. FORNBERG, *A Padé-based algorithm for overcoming the Gibbs phenomenon*, Numer. Algo. **26**(2001), pp. 77-92.
- [7] A. ERDÉLYI, W. MAGNUS, F. OBERHETTINGER, AND F.G. TRICOMI, *Higher Transcendental Functions I+II+III*, Robert E. Krieger Publishing, Malabar, FL, 1981.
- [8] J.F. GEER, *Rational Trigonometric Approximations using Fourier Series Partial Sums*, J. Sci. Comput. **10**(1995), pp. 325-356.
- [9] D. GOTTLIEB AND S. A. ORSZAG, *Numerical Analysis of Spectral Methods: Theory and Applications*. CBMS-NSF **26**. SIAM, Philadelphia, 1978.
- [10] D. GOTTLIEB AND C.W. SHU, *On the Gibbs Phenomenon and its Resolution*, SIAM Review **39**(1997), pp. 644-668.

- [11] D. GOTTLIEB AND E. TADMOR, *Recovering Pointwise Values of Discontinuous Data with Spectral Accuracy*. In PROGRESS AND SUPERCOMPUTING IN COMPUTATIONAL FLUID DYNAMICS. Birkhäuser, Boston, 1984. pp. 357-375.
- [12] , J. S. HESTHAVEN, S.M. KABER, AND L. LURATI, *The Padé-Legendre Interpolation Method*, In preparation.
- [13] J.S. HESTHAVEN AND M. KIRBY, *Filtering in Legendre Spectral Methods*, Math. Comp. 2005 – submitted.
- [14] A. J. JERRI, *The Gibbs phenomenon in Fourier analysis, splines and wavelet approximations*, Kluwer Academic Publishers, Dordrecht, 1998.
- [15] S.M. KABER, *The Gibbs phenomenon for Jacobi expansions*, Technical rep. Université Pierre & Marie Curie, 2004. Available at <http://www.ann.jussieu.fr/publications/2005/R05003.html>.
- [16] S.M. KABER AND Y. MADAY, *Analysis of some Padé-Chebyshev approximants*. Accepted for publication in SIAM Journal on Numerical Analysis, 2004.
- [17] G. NÉMETH, AND G. PARIS, *The Gibbs phenomenon in generalized Padé approximation*. J. Math. Phys. 26, 1985.
- [18] G. SZEGÖ, *Orthogonal Polynomials*, American Mathematical Society, 1930.
- [19] J. TANNER AND E. TADMOR, *Adaptive Mollifiers - High Resolution Recover of Piecewise Smooth Data from its Spectral Information*, Found. Comput. Math. **2**(2002), pp. 155-189.
- [20] H. VANDEVEN, *Family of Spectral Filters for Discontinuous Problems*, J. Sci. Comput. **8**(1991), pp. 159-192.