

Exchange algorithm for evaluation and approximation error-optimized polynomials

Denis Arzelier, Florent Bréhard, Mioara Joldes

► **To cite this version:**

Denis Arzelier, Florent Bréhard, Mioara Joldes. Exchange algorithm for evaluation and approximation error-optimized polynomials. Rapport LAAS n° 19038. 2019. <hal-02006606>

HAL Id: hal-02006606

<https://hal.archives-ouvertes.fr/hal-02006606>

Submitted on 4 Feb 2019

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.

Exchange algorithm for evaluation and approximation error-optimized polynomials

Denis Arzelier¹, Florent Bréhard², and Mioara Joldes³

¹LAAS-CNRS, Toulouse, France, arzelier@laas.fr

²ENS de Lyon & LAAS-CNRS, Lyon & Toulouse, France,
florent.brehard@ens-lyon.fr

³LAAS-CNRS, Toulouse, France, joldes@laas.fr

February 4, 2019

Abstract

Machine implementation of mathematical functions often relies on polynomial approximations. The particularity is that rounding errors occur both when representing the polynomial coefficients on a finite number of bits, and when evaluating it in finite precision. Hence, for finding the best polynomial (for a given fixed degree, norm and interval), one has to consider both types of errors: approximation and evaluation. While efficient algorithms were already developed for taking into account the approximation error, the evaluation part is usually a posteriori handled, in an ad-hoc manner. Here, we formulate a semi-infinite linear optimization problem whose solution is the best polynomial with respect to the supremum norm of the sum of both errors. This problem is then solved with an iterative exchange algorithm, which can be seen as an extension of the well-known Remez algorithm. A discussion and comparison of the obtained results on different examples are finally presented.

1 Introduction

Polynomials are often used for approximating functions on computers [1,21]. Their evaluation only requires additions and multiplications, which are efficiently implemented in hardware floating-point (FP) arithmetic units. FP

operations are specified by the IEEE 754-2008 [17] standard, which requires, among others, correctly rounded basic arithmetic operations $(+, -, *, /, \sqrt{})$ for several precision formats, and recommends correctly rounded elementary functions like \exp , \sin , \cos . Very efficient fixed FP precision implementations exist [15, 20] for such functions and are collected in mathematical libraries (*libms*), which can be nowadays almost automatically generated and tuned [12, 18]. Recently, in [19], such code generating techniques were extended to larger classes of special functions, which are widely used in scientific and technical applications (like Bessel, Airy, Erf, etc.).

The problem of evaluating a function for the whole FP input range is firstly reduced to the evaluation of an approximation valid in a rather small compact domain I . This can be done for instance, by argument reduction techniques, which are available only for specific elementary functions, and/or by piecewise polynomial approximations [22]. Then, the implementation task becomes: given a description of a function f , an input interval I , and a target accuracy $\varepsilon > 0$, one is requested a source code which provides a function \tilde{f} , such that: $\left\| (f - \tilde{f})/f \right\|_I \leq \varepsilon$, where we denote by $\|g\|_I := \sup_{t \in I} |g(t)|$ the supremum norm of g on I .

Typically, this is handled in two main steps:

Approximation an approximation polynomial p is searched for, such that two main requirements are met: its coefficients are representable with a specified fixed precision format (usually, binary32, binary64, or an unevaluated sum of such formats) and the approximation error is less than a target $\varepsilon_{\text{approx}}$, whether absolute $\|f - p\|_I \leq \varepsilon_{\text{approx}}$ or relative $\|(f - p)/f\|_I \leq \varepsilon_{\text{approx}}$.

For that, efficient algorithms were developed in [5, 6]. In the simpler case of polynomials p with real coefficients and given degree n , $p = \sum_{i=0}^n a_i t^i$, this boils down to the so-called *minimax* problem:

$$\min_{\substack{a_i \in \mathbb{R}, \\ i \in [0..n]}} \max_{t \in I} |f(t) - p(t)|, \quad (P_{\text{minimax}})$$

which can be solved by the Remez algorithm (see [5, 10] and references therein). This iterative algorithm has quadratic convergence and rather low complexity, since it involves solving a linear system of size $n + 2$ at each step, together with numerically computing the extrema of $f - p$ over I .

Evaluation an efficient evaluation scheme for p is searched for; since after each addition or multiplication, rounding errors occur, one must ensure that

the computed value \tilde{p} satisfies $\|p - \tilde{p}\|_I \leq \varepsilon_{\text{eval}}$ (or $\|(p - \tilde{p})/p\|_I \leq \varepsilon_{\text{eval}}$) for a given threshold $\varepsilon_{\text{eval}}$.

Heuristics presented in [20] extend the precision of the *important* coefficients, such that the evaluation error remains below $\varepsilon_{\text{eval}}$. For instance, Sollya command *implementpoly* uses a Horner-based evaluation scheme, which behaves rather well when the evaluation interval is sufficiently small and contains zero. Otherwise, consider step i of Horner evaluation $a_i + t\tilde{p}_i(t)$, where \tilde{p}_i is the already computed partial polynomial evaluation: when the argument $|t| \gg 1$, the accumulated evaluation error is much amplified when multiplying by t . Another heuristic is a ratio test between a_i and $t\tilde{p}_i(t)$, to check for cancellation issues which appear when both terms have the same order of magnitude and opposite signs.

Once the coefficients have been chosen, the approximation and the evaluation error can *a posteriori* be certified by several existing algorithms and tools, like Sollya [11], Gappa [14], Rosa [13] or Real2Float [25].

It is important to note that steps *a*) and *b*) are usually independently considered. An exception occurs for the case of very small precisions or polynomial degrees, where an exhaustive search on the rounded coefficients is possible [28].

However, as explicitly mentioned in [6], *one would like to take into account the roundoff error that occurs during polynomial evaluation: getting the polynomial, with constraints on the size of the coefficients, that minimizes the total (approximation plus roundoff) error would be extremely useful.*

The purpose of this article is to make progress on this open question: we search for the coefficients of a polynomial $p(t) = \sum_{i=0}^n a_i t^i$, of given degree n , which minimizes the maximum of the sum of both approximation and evaluation errors over an input interval I , with respect to f . We consider a *black-box* description of f i.e., one disposes of values $f(t)$, up to any required accuracy [20]. This allows for handling very general functions (elementary, special, etc.), but also implies that no argument reduction step is usually possible. For simplicity we state the problem for the absolute error case:¹

$$\min_{\substack{a_i \in \mathbb{R}, \\ i \in [0..n]}} \max_{t \in I} (|f(t) - p(t)| + |\tilde{p}(t) - p(t)|) \quad (P_{\text{general}})$$

In Section 2, we give a linearized bound for the evaluation error $|\tilde{p}(t) - p(t)|$. Based on [23], the performance of a given arbitrary evaluation scheme is recursively assessed by bounding the rounding error of each elementary

¹The relative error can be similarly handled from a theoretical standpoint (cf. Appendix .2).

operation. This leads to the formulation in Section 3 of Problem P_{general} as a linear semi-infinite programming (SIP) problem [24, 26].

In this context, we show two results: on the theoretical side, based on the duality theory, we revisit, explain and extend an *exchange algorithm* [7–9, 29], which solves this problem in Section 4. On the practical side, the solution of this problem provides a first attempt on simultaneously optimizing over both errors: we show that in some cases the evaluation error can be improved. We also show that in some other cases, the minimax polynomial solution of problem P_{minimax} is very close to the solution of P_{general} . Numerical examples and a discussion are provided in Section 5.

2 Evaluation error

Consider first some basic notation used for error analysis [22]. Firstly, assume radix-2, precision- p , floating-point arithmetic with unbounded exponent range i.e, provided that overflows and underflows do not occur. If $t \in \mathbb{R}$, define $\text{RN}(t)$ as t rounded to nearest. This is the default rounding mode in IEEE-754 arithmetic [17], so that, given two FP numbers a and b , when the instruction $c = a \top b$ appears in a program, what is effectively computed is $c = \text{RN}(a \top b)$, for any arithmetic operation $\top \in \{+, -, \times, \div\}$. We have

$$\frac{|t - \text{RN}(t)|}{|t|} \leq \frac{u}{1 + u} < u, \quad (1)$$

where $u = 2^{-p}$ is called the *rounding unit*.

Moreover, there exists a real number ϵ such that

$$\text{RN}(a \top b) = (a \top b)(1 + \epsilon), \quad |\epsilon| \leq u. \quad (2)$$

Based on the previous property, the error of any arithmetic expression can be recursively bounded. Firstly, for specific evaluation schemes, like Horner, bounds date back to the work of Oliver [23], which is detailed below as an example. More recently, several works insisted on the automatic algorithmic approach via operator overloading similar to automatic differentiation [4, 27]. Based on this, we propose, for completeness², Algorithm 2, which automatically computes linearized expressions for the evaluation error (like in (6)), for any given symbolic expression tree e , provided with symbolic rounding errors for each tree node. Let us exemplify on the evaluation of the polyno-

²While the general ideas are the same as in [4, 27] and references therein, we could not find the exact pseudo-code in literature, so we state it in order to provide a complete algorithmic solution for problem P_{general} .

Algorithm 1 HORNER(p,t).

```

1:  $r_n \leftarrow a_n$ 
2: for  $i = n - 1$  downto 0 do
3:    $r_i \leftarrow \text{RN}(\text{RN}(r_{i+1} \times t) + a_i)$ 
4: end for
5: return  $r_0$ 

```

mial $p(t) = a_n t^n + a_{n-1} t^{n-1} + \dots + a_0$ using Horner's rule, assuming that a Fused Multiply Add (FMA) instruction is not employed. The actual machine operations are recalled in Algorithm 1. We have:

$$r_n = a_n, \quad (3a)$$

$$r_{n-1} = (tr_n(1 + \epsilon_{n-1}^\times) + a_{n-1})(1 + \epsilon_{n-1}^+), \quad (3b)$$

where ϵ_{n-1}^\times and ϵ_{n-1}^+ model the rounding errors for multiplication and addition at step $n - 1$. By induction, one obtains:

$$r_k = \sum_{i=k}^n \left((1 + \epsilon_i^+) \prod_{j=k}^{i-1} (1 + \epsilon_j^+) (1 + \epsilon_j^\times) \right) a_i t^{i-k}, \quad (4)$$

where we define $\epsilon_n^+ := 0$ and $\prod_{j=k}^{k-1} (1 + \epsilon_j^+) (1 + \epsilon_j^\times) := 1$. This implies that the total evaluation error is:

$$r_0 - \sum_{i=0}^n a_i t^i = \sum_{i=0}^n \left((1 + \epsilon_i^+) \prod_{j=0}^{i-1} (1 + \epsilon_j^+) (1 + \epsilon_j^\times) - 1 \right) a_i t^i. \quad (5)$$

Here, we consider only a linear approximation θ_{lin} of the evaluation error, function of ϵ_i^+ and ϵ_i^\times , in what follows. This gives, for our Horner example:

$$\theta_{\text{lin}}^{(\text{Horner})} := \sum_{j=0}^{n-1} \left(\sum_{i=j+1}^n a_i t^i \right) \epsilon_j^\times + \sum_{j=0}^{n-1} \left(\sum_{i=j}^n a_i t^i \right) \epsilon_j^+. \quad (6)$$

Moreover, provided bounds are specified for each rounding error, depending on the precision employed, one obtains upper bounds for the linearized absolute evaluation error. For instance, if binary64 is used for all the computations in Algorithm 1, with $u = 2^{-53}$, one has:

$$\left| \theta_{\text{lin}}^{(\text{Horner})} \right| \leq 2u \sum_{j=0}^{n-1} \left\| \sum_{i=j}^n a_i t^i \right\|, \quad (7)$$

where the double superscript indicates that the first and last terms in the summation are to be halved.

In order to automate this evaluation error analysis for the general case, we firstly associate to a "pure" mathematical expression e^* , a given symbolic evaluation scheme with roundings e , composed of terms $\text{RN}(e', u)$. This means that e' is rounded with a relative error bounded by u . This formulation models both possible rounding errors on an input variable ($e' \in \mathcal{V}$, where \mathcal{V} denotes the set of input variables) and the rounding errors of arithmetic operations (if $e' = a_1 \top e_2$). Then, we build an expression \tilde{e} , as in (4), by recursively replacing terms $\text{RN}(e', u)$ in e , with $\tilde{e}'(1 + \epsilon_{e'}^{[u]})$ where $|\epsilon_{e'}^{[u]}| \leq u$. Several examples are shown in Table 1.

$e_1 = \text{RN}(a + \text{RN}(b \times c, u), u)$ $\tilde{e}_1 = (a + bc(1 + \epsilon_{b \times c}^{[u]}))(1 + \epsilon_{a + \text{RN}(b \times c, u)}^{[u]})$ $\theta_{\text{lin}}^{(e_1)} = bc\epsilon_{b \times c}^{[u]} + (a + bc)\epsilon_{a + \text{RN}(b \times c, u)}^{[u]}$ $ \theta_{\text{lin}}^{(e_1)} \leq (bc + a + bc)u$ <p>▷ <i>Arithmetic operations in double precision.</i></p>	$e_1^* = a + bc$
$e_2 = \text{RN}(a + b \times c, u)$ $\tilde{e}_2 = (a + bc)(1 + \epsilon_{a + b \times c}^{[u]})$ $\theta_{\text{lin}}^{(e_2)} = (a + bc)\epsilon_{a + b \times c}^{[u]}$ $ \theta_{\text{lin}}^{(e_2)} \leq a + bc u$ <p>▷ <i>Fused multiply-add (FMA) in double precision.</i></p>	$e_2^* = a + bc$
$e_3 = \text{RN}(\text{RN}(a, u') + b \times \text{RN}(a, u'), u)$ $\tilde{e}_3 = (a(1 + \epsilon_a^{[u']}) + ba(1 + \epsilon_a^{[u']}))(1 + \epsilon_{\text{RN}(a, u') + b \times \text{RN}(a, u')}^{[u]})$ $\theta_{\text{lin}}^{(e_3)} = (a + ba)\epsilon_a^{[u']} + (a + ba)\epsilon_{\text{RN}(a, u') + b \times \text{RN}(a, u')}^{[u]}$ $ \theta_{\text{lin}}^{(e_3)} \leq a + ba (u + u')$ <p>▷ <i>Fused multiply-add (FMA) in double precision, with input a rounded to single precision.</i></p>	$e_3^* = a + ba$

Table 1: Evaluation error examples ($u = 2^{-53}$, $u' = 2^{-24}$).

Finally, automatic linearized evaluation error expressions θ_{lin} , such as in (6), are obtained using Algorithm 2. Specifically, for an arithmetic expression with roundings e , this algorithm recursively computes an expression of the form $\theta_{\text{lin}} = \sum_{i=1}^k \theta_{\text{lin},i} \epsilon_{e_i}^{[u_i]}$, with symbolic $\epsilon_{e_i}^{[u_i]}$ ($i \in [1..k]$) for each term $\text{RN}(e_i, u_i)$ in e . The coefficients $\theta_{\text{lin},i}$ are arithmetic expressions depending only on the input variables in e . Note that $\text{RN}(e_i, u_i)$ may occur several

times in e , but the error variable $\epsilon_{e_i}^{[u_i]}$ is unique since the rounding operation is deterministic. This allows to bound the (linearized) evaluation error as in (7).

Proposition 1 (Correctness of Algorithm 2). *Let e be an arithmetic expression with roundings, and $\theta_{lin} = \sum_{i=1}^k \theta_{lin,i} \epsilon_{e_i}^{[u_i]}$ the linearized expression for the evaluation error returned by Algorithm 2. If $u_i \leq u$ for all $i \in [1 \dots k]$, then:*

$$|\tilde{e} - e^*| \leq \sum_{i=1}^k |\theta_{lin,i}| u_i + \mathcal{O}(u^2), \quad \text{as } u \rightarrow 0.$$

Usually, for polynomial evaluation schemes, the functions $\theta_{lin,i}$ are linear with respect to the coefficients \mathbf{a} of $p(t)$, that is $u_i \theta_{lin,i} = \boldsymbol{\pi}_i(t)^T \mathbf{a}$ for some $\boldsymbol{\pi}_i(t) \in \mathbb{R}^{n+1}$. Hence we obtain a linearized bound of the evaluation error of the form:

$$|\theta_{lin}(\mathbf{a}, t)| \leq \sum_{i=1}^k |\boldsymbol{\pi}_i(t)^T \mathbf{a}|, \quad \mathbf{a} \in \mathbb{R}^{n+1}, t \in \mathbb{R}. \quad (8)$$

We denote the right-hand side of equation (9) by

$$\theta(\mathbf{a}, t) = \sum_{i=1}^k |\boldsymbol{\pi}_i(t)^T \mathbf{a}|, \quad \mathbf{a} \in \mathbb{R}^{n+1}, t \in \mathbb{R}. \quad (9)$$

Example 1 (Bound evaluation error for Horner). *In particular, for Horner evaluation, we have from equation (7):*

$$\begin{aligned} \boldsymbol{\pi}_1(t)^T &= (u, ut, \dots, ut^{n-1}, ut^n) \\ \boldsymbol{\pi}_2(t)^T &= (0, 2ut, \dots, 2ut^{n-1}, 2ut^n) \\ &\dots \\ \boldsymbol{\pi}_n(t)^T &= (0, 0, \dots, 2ut^{n-1}, 2ut^n) \\ \boldsymbol{\pi}_{n+1}(t)^T &= (0, 0, \dots, 0, ut^n) \end{aligned}$$

Algorithm 2 LINEVALERROR(e)

Input: e an arithmetic expression with explicit roundings.

Output: θ_{lin} the linearized evaluation error of e .

```
if  $e \in \mathcal{V}$  then
  return 0
else if  $e = \text{RN}(f, u)$  then
   $\theta'_{\text{lin}} \leftarrow \text{LINEVALERROR}(f)$ 
  return  $\theta'_{\text{lin}} + f^* \epsilon_f^{[u]}$ 
else if  $e = -f$  then
   $\theta'_{\text{lin}} \leftarrow \text{LINEVALERROR}(f)$ 
  return  $-\theta'_{\text{lin}}$ 
else if  $e = f + g$  then
   $\theta'_{\text{lin}} \leftarrow \text{LINEVALERROR}(f)$ 
   $\theta''_{\text{lin}} \leftarrow \text{LINEVALERROR}(g)$ 
  return  $\theta'_{\text{lin}} + \theta''_{\text{lin}}$ 
else if  $e = f \times g$  then
   $\theta'_{\text{lin}} \leftarrow \text{LINEVALERROR}(f)$ 
   $\theta''_{\text{lin}} \leftarrow \text{LINEVALERROR}(g)$ 
  return  $g^* \theta'_{\text{lin}} + f^* \theta''_{\text{lin}}$ 
end if
```

3 Semi-Infinite Programming formulation

3.1 Problem (P_{general}) as a linear SIP

Noting that Problem (P_{general}) is a piecewise-linear optimization problem and using the convex evaluation error formula $\theta(\mathbf{a}, t)$ at point $t \in [t_l, t_r]$ obtained in Section 2, Problem (P_{general}) becomes Problem (P'_{general}) (see [3, Section 4.3.1] for instance), with the compact index set $I = [t_l, t_r]$ and the monomial basis $\boldsymbol{\pi}_0(t) = (1, \dots, t^n)^T$.

$$\begin{aligned} \min_{(\bar{a}, \mathbf{a}) \in \mathbb{R}^{n+2}} \quad & \bar{a} \\ \text{s.t.} \quad & |f(t) - \boldsymbol{\pi}_0(t)^T \mathbf{a}| + \theta(\mathbf{a}, t) - \bar{a} \leq 0, \quad t \in I. \end{aligned} \quad (P'_{\text{general}})$$

Problem (P'_{general}) is a convex Semi-Infinite Programming (SIP) problem (see [24] which provides a comprehensive overview of SIP) that can be reformulated as a linear SIP problem, at the expense of a different index set Ω replacing the previous index set I . Here, the set of constraints of (P'_{general}) involving absolute values is replaced by as many linear constraints as required

to represent all possible sign combinations. More precisely, the evaluation

error is as in equation (9), $\theta(\mathbf{a}, t) = \sum_{i=1}^k |\mathbf{a}^T \boldsymbol{\pi}_i(t)|$, and define:

$$\mathbf{x} = (\bar{a}, \mathbf{a}) \in \mathbb{R}^{n+2}, \quad \mathbf{z} = (1, 0, \dots, 0) \in \mathbb{R}^{n+2},$$

$$\boldsymbol{\alpha}(t, \sigma_0, \dots, \sigma_k) = (1, \sigma_0 \boldsymbol{\pi}_0^T(t) + \sum_{i=1}^k \sigma_i \boldsymbol{\pi}_i^T(t))^T \in \mathbb{R}^{n+2}, \quad (10)$$

$$\mathfrak{S} = \{-1, 0, 1\}^{k+1}, \quad \omega = (t, \sigma_0, \dots, \sigma_k) \in \Omega := I \times \mathfrak{S}.$$

Then, Problem (P'_{general}) is exactly given by the following linear SIP:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & \mathbf{z}^T \mathbf{x} \\ \text{s.t.} \quad & \boldsymbol{\alpha}(\omega)^T \mathbf{x} \geq c(\omega), \quad \omega \in \Omega, \end{aligned} \quad (P)$$

where $d = n+2$, $c(\omega) = \sigma_0 f(t)$, Ω is a compact metric space and the function $g(\mathbf{x}, \omega) = c(\omega) - \boldsymbol{\alpha}(\omega)^T \mathbf{x} \leq 0$ defining the feasible set is a continuous function from $\mathbb{R}^{n+2} \times \Omega$ into \mathbb{R} . Note that for $\mathfrak{S}' = \{-1, 1\} \times \{0\}^k$ and $\Omega' = I \times \mathfrak{S}' \subseteq \Omega$, (P_{minimax}) is exactly retrieved as shown in the next example.

Example 2. For $n = 5$, Problem (P_{minimax}) is:

$$\begin{aligned} \min_{(\bar{a}, \mathbf{a}) \in \mathbb{R}^7} \quad & \bar{a} \\ \text{s.t.} \quad & (1, \sigma_0 1, \sigma_0 t, \dots, \sigma_0 t^5)(\bar{a}, a_0, a_1, \dots, a_5)^T \geq \sigma_0 f(t), \\ & \sigma_0 = \mp 1, \quad t \in I. \end{aligned} \quad (\text{Example 2 (a)})$$

while Problem (P'_{general}), assuming Horner evaluation is:

$$\begin{aligned} \min_{(\bar{a}, \mathbf{a}) \in \mathbb{R}^7} \quad & \bar{a} \\ \text{s.t.} \quad & (1, \sigma_0 + \sigma_1 u, (\sigma_0 + \sigma_1 u + \sigma_2 2u)t, \dots, \\ & (\sigma_0 + \sigma_1 u + \dots + \sigma_5 u)t^5)(\bar{a}, a_0, a_1, \dots, a_5)^T \\ & \geq \sigma_0 f(t), \\ & \sigma_0 = \mp 1, \sigma_1 = \mp 1, \dots, \sigma_5 = \mp 1, \quad t \in I. \end{aligned} \quad (\text{Example 2 (b)})$$

We propose in Section 4 an exchange algorithm which solves Problem (P'_{general}) and can be seen as a generalization (in the above framework) of the classical Remez algorithm, which solves Problem (P_{minimax}). In order to prove its correctness, we need some important discretization properties of linear SIP problems, which are based on conjugate duality theory. Our presentation closely follows the survey [26].

3.2 Duality and Discretization for SIP

For a Problem (P) , we denote respectively by $\text{val}(P)$ and $\text{Sol}(P)$, its optimal value and the set of its optimal solutions.

A discretization (P_m) of (P) for a set $\omega = \{\omega_1, \dots, \omega_m\} \subseteq \Omega$ is the following linear program:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & \mathbf{z}^T \mathbf{x} \\ \text{s.t.} \quad & \boldsymbol{\alpha}(\omega_j)^T \mathbf{x} \geq c(\omega_j), \quad j = 1, \dots, m. \end{aligned} \tag{P_m}$$

Since the feasible set of (P) is included in the feasible set of (P_m) , we have that $\text{val}(P_m) \leq \text{val}(P)$. The existence of a discretization (P_m) such that the equality holds is a particularly appealing feature of some linear SIPs since the solution of (P) may be obtained by the solution of (P_m) if we are able to find the corresponding set ω .

Definition 1. [26] (P) is said to be reducible if there exists a discretization (P_m) defined by the subset $\{\omega_1, \dots, \omega_m\} \subseteq \Omega$ such that $\text{val}(P_m) = \text{val}(P)$.

The characterization of reducible SIP problems relies on the central notion of *duality* that rules the interplay between two optimization problems. This notion has its roots in the interrelations between a normed linear space and its topological dual. Let us define the continuous mapping $h : \mathbf{x} \mapsto g(\mathbf{x}, \cdot)$ from \mathbb{R}^d to the Banach space of continuous functions $\mathcal{C}(\Omega)$, equipped with the uniform norm $\|h\|_\Omega = \sup_{\omega \in \Omega} |h(\omega)|$, then the topological dual of $\mathcal{C}(\Omega)$ is the space $\mathcal{C}(\Omega)^*$ of signed Borel measures μ over $(\Omega, \mathcal{B}(\mathbb{R}^{k+2}))$ [16, Section 21.5]. For a measure $\mu \in \mathcal{C}(\Omega)^*$, its support is the smallest closed subset Γ of Ω such that $|\mu|(\Omega \setminus \Gamma) = 0$. A positive measure μ is denoted by $\mu \succeq 0$. A classical example of a positive measure with discrete support is the Dirac measure of support $\{\omega_j\}$:

$$\delta_{\omega_j}(A) = \begin{cases} 0 & \text{if } \omega_j \notin A, \\ 1 & \text{if } \omega_j \in A. \end{cases} \quad A \subseteq \Omega. \tag{11}$$

Defining the bilinear form pairing $\mathcal{C}(\Omega)$ and $\mathcal{C}(\Omega)^*$ by the duality bracket:

$$\langle h, \mu \rangle = \int_{\Omega} h(\omega) d\mu(\omega), \tag{12}$$

the *dual* problem (D) associated to the *primal* problem (P) is obtained as:

$$\begin{aligned} \max_{\mu \succeq 0} \quad & \int_{\Omega} c(\omega) d\mu(\omega) \\ \text{s.t.} \quad & \int_{\Omega} \boldsymbol{\alpha}(\omega) d\mu(\omega) = \mathbf{z}. \end{aligned} \tag{D}$$

The *weak duality* property, that is $\text{val}(D) \leq \text{val}(P)$ always holds. Problem (D) is an LP problem defined in the space of positive measures which is hard to solve. By restricting the support of $\mu \succeq 0$ to $\{\omega_1, \dots, \omega_m\}$, that is $\mu = \sum_{j=1}^m y_j \delta_{\omega_j}$ with $y_j \geq 0$, a discretized counterpart (D_m) of (D) is obtained:

$$\begin{aligned} \max_{\substack{y_j \geq 0 \\ j \in [1..m]}} & \sum_{j=1}^m c(\omega_j) y_j \\ \text{s.t.} & \sum_{j=1}^m y_j \boldsymbol{\alpha}(\omega_j) = \mathbf{z}, \end{aligned} \tag{D_m}$$

with $\text{val}(D_m) \leq \text{val}(D)$. It is important to note that the LP dual of the discretized problem (P_m) is exactly (D_m) which implies that $\text{val}(D_m) = \text{val}(P_m)$ (strong duality holds) provided that none of (P_m) or (D_m) is infeasible.

So far, under these mild assumptions, we have that $\text{val}(D_m) = \text{val}(P_m) \leq \text{val}(D) \leq \text{val}(P)$ and conditions for having only equalities (respectively reducibility and strong duality properties) may be obtained by using conjugate duality theory as developed in [26, Theorems 2.2, 2.3 and 3.2].

Theorem 1. [26, Thm. 2.2, 2.3, 3.2] *Under the assumptions:*

A1 Ω is a compact metric space, $\boldsymbol{\alpha} : \Omega \rightarrow \mathbb{R}^d$ and $c : \Omega \rightarrow \mathbb{R}$ are continuous functions;

A2 $\text{val}(P)$ is finite;

A3 (Slater's condition): there exists \mathbf{x}° such that:

$$\boldsymbol{\alpha}(\omega)^T \mathbf{x}^\circ > c(\omega), \quad \text{for all } \omega \in \Omega; \tag{13}$$

A4 There exist $\omega_1, \dots, \omega_d \in \Omega$ with $(\boldsymbol{\alpha}(\omega_1), \dots, \boldsymbol{\alpha}(\omega_d))$ linearly independent such that:

$$\exists y_1, \dots, y_d > 0, \quad \mathbf{z} = \sum_{j=1}^d y_j \boldsymbol{\alpha}(\omega_j), \tag{14}$$

the following statements are true:

(i) $\text{Sol}(P) \neq \emptyset$ and bounded;

(ii) $\text{Sol}(D) \neq \emptyset$ and bounded;

(iii) Problem (P) is reducible to a Problem (P_m) with $m \leq d$;

(iv) $\text{val}(P) = \text{val}(D) = \text{val}(P_m) = \text{val}(D_m)$.

Proposition 2. *Assumptions A1-A4 are satisfied for our Problem (P'_{general}) and therefore results (i)-(iv) of Theorem 1 apply.*

Proof.

- A1 By construction, our set Ω is a compact metric space and α and c are polynomials and therefore continuous on Ω ;
- A2 $\text{val}(P) = +\infty$ means that the primal problem (P'_{general}) is not feasible but $\mathbf{x} = (\max_{t \in I} |f(t)|, 0 \cdots, 0)$ is a feasible point for (P'_{general}) , therefore $\text{val}(P) < +\infty$. In addition, $\text{val}(P) > -\infty$ since $\bar{a} \geq 0$ by construction and for all feasible points of (P'_{general}) ;
- A3 It may be easily deduced from the proof of A2 that $\mathbf{x}^\circ = (\max_{t \in I} |f(t)| + \varsigma, 0 \cdots, 0)$ is a strictly feasible point for (P'_{general}) for any $\varsigma > 0$;
- A4 An instance for $\{\omega_1, \dots, \omega_{n+2}\}$ is provided by Algorithm INIT in Section 4 (see Lemma 1 in Appendix .1).

□

The conclusion of this theorem, that is both (P) and (D) are reducible to a discretization of size of at most d , allows us to recast the problem of solving (P'_{general}) as the problem of finding the right discretization $\{\omega_1, \dots, \omega_d\}$ such that item (iv) of Theorem 1 applies and to solve the associated (P_m) and/or (D_m) . This goal may be reached by tailoring the general exchange algorithm for semi-infinite linear programs presented in [7] to our specific case.

This algorithm can be seen as a generalization of the dual simplex algorithm for Problem (D) . The main idea consists first in finding at each iteration ℓ , the solution $\mathbf{y}^{(\ell)}$ of $(D_{n+2}^{(\ell)})$, with $\boldsymbol{\omega}^{(\ell)} = \{\omega_j^{(\ell)}\}_{j=1}^{n+2}$. Such a solution is a feasible (but not necessarily optimal) point of the dual Problem (D) . Moreover, the objective value $\mathbf{z}^T \mathbf{x}^{(\ell)}$ of $(P_m^{(\ell)})$ and (P) for the instance $\mathbf{x}^{(\ell)} := (\bar{a}^{(\ell)}, \mathbf{a}^{(\ell)})$ is equal to the objective value of (D) for the instance $\mathbf{y}^{(\ell)}$ ³. Hence, either $\mathbf{x}^{(\ell)}$ is a feasible solution of Problem (P) by Theorem 1, or it is an infeasible point of Problem (P) . In the latter case, one of these constraints is replaced by a new one, indexed by $\omega_*^{(\ell)}$, in an exchange step in order to increase the objective value of the dual and works towards primal feasibility.

³The feasible set of (P) is included in the feasible set of $(P_m^{(\ell)})$, for all ℓ .

4 Iterative exchange algorithm

Algorithm EVAL&APPROXOPTIMIZE computes the degree- n best polynomial approximation with respect to both evaluation and approximation errors, i.e. it solves (P'_{general}) based on the theoretical developments from Section 3.2. Before entering the details, we provide an analogy with Remez algorithm. Roughly speaking, the new algorithm consists of the same main steps:

- INIT provides *a good set of initial points*.
- At each step, SOLVEPRIMAL solves a linear system of equations (built w.r.t. the current set of points), where the variables are the polynomial coefficients.
- Then, FINDNEWINDEX finds a new point where the *total error* is maximal.
- Finally, EXCHANGE replaces one point from the current set with this new point.
- This process repeats until a certain tolerance threshold is reached concerning the *total error*.

However, when considering both errors, one can not only rely on the primal problem (coefficients reconstruction), but also needs the dual problem. This implies:

- Besides classical points, a combination of signs (signatures) is required at each step.
- INIT and EXCHANGE need the solution of the dual problem.

A running example for this algorithm is given in Section 5. We focus now on its correctness, which is stated in Theorem 2. For this, one needs an assumption on the dual solution, which always holds in the Remez algorithm. It is not proven in our setting, but it never failed in practice.

Assumption 1. *At each iteration ℓ , the solution $\mathbf{y}^{(\ell)}$ of the dual discretized Problem $(D_{n+2}^{(\ell)})$ is an interior point, that is $y_j^{(\ell)} > 0$ for all $j \in [1 \dots n + 2]$.*

Theorem 2. *Let f be a continuous function over an interval $I = [t_l, t_r]$, a degree $n \geq 0$, a linearized evaluation error bound θ and a tolerance parameter*

$\tau > 0$. Under Assumption 1, `EVAL&APPROXOPTIMIZE`(f, n, I, θ, τ) terminates and returns a degree- n polynomial approximation for f with a total error ε (approximation and evaluation) satisfying:

$$\varepsilon^* \leq \varepsilon \leq (1 + \tau)\varepsilon^*, \quad (15)$$

where ε^* is the total error of the best degree- n polynomial approximation of f .

Proof. • First, we prove by induction that the following properties hold at each iteration $\ell \geq 0$:

- (i) $\{\boldsymbol{\alpha}(\omega_j^{(\ell)})\}_{j=1}^{n+2}$ is a basis of \mathbb{R}^{n+2} ;
- (ii) $\mathbf{y}^{(\ell)}$ is the optimal solution of Problem (D_m) for $\boldsymbol{\omega}^{(\ell)}$;
- (iii) $\mathbf{x}^{(\ell)}$ is the optimal solution of Problem (P_m) for $\boldsymbol{\omega}^{(\ell)}$;
- (iv) $\omega_*^{(\ell)} = \arg \max_{\omega \in \Omega} (c(\omega) - \boldsymbol{\alpha}(\omega)^T \mathbf{x}^{(\ell)})$;

using the correctness Lemmas given in Appendix .1. For $\ell = 0$, `INIT`(n, I) returns $\boldsymbol{\omega}^{(0)}, \mathbf{y}^{(0)}$ satisfying (i) and (ii). Then `SOLVEPRIMAL`($f, n, \theta, \boldsymbol{\omega}^{(0)}$) computes $\mathbf{x}^{(0)} = (\bar{a}^{(0)}, \mathbf{a}^{(0)})$ satisfying (iii). Finally, `FINDNEWINDEX`($f, n, I, \theta, \mathbf{a}^{(0)}$) gives $\omega_*^{(0)}, \bar{a}_*^{(0)}$ satisfying (iv).

For the inductive step, `EXCHANGE`($n, \theta, \boldsymbol{\omega}^{(\ell)}, \mathbf{y}^{(\ell)}, \omega_*^{(\ell)}$) computes $\boldsymbol{\omega}^{(\ell+1)}, \mathbf{y}^{(\ell+1)}$ satisfying (i) and (ii), by induction hypothesis on $\boldsymbol{\omega}^{(\ell)}, \mathbf{y}^{(\ell)}, \omega_*^{(\ell)}$. Then, `SOLVEPRIMAL`($f, n, \theta, \boldsymbol{\omega}^{(\ell+1)}$) and `FINDNEWINDEX`($f, n, I, \theta, \mathbf{a}^{(\ell+1)}$) compute $\mathbf{x}^{(\ell+1)}, \omega_*^{(\ell+1)}, \bar{a}_*^{(\ell+1)}$ satisfying (iii) and (iv).

• Moreover, at each iteration ℓ , we have $\bar{a}^{(\ell)} \leq \varepsilon^* \leq \bar{a}_*^{(\ell)}$. Indeed, $\mathbf{x}^{(\ell)}$ is the optimal solution of the discretized Problem (P_m) for $\boldsymbol{\omega}^{(\ell)}$, whose objective value $\bar{a}^{(\ell)}$ is less or equal to the optimal value ε^* of Problem (P). On the other side, $\bar{a}_*^{(\ell)}$ is the total error of degree- n polynomial $\mathbf{a}^{(\ell)T} \boldsymbol{\pi}_0(t)$ and therefore, it is greater or equal to the optimal error ε^* . In addition, Lemma 5 in the Appendix proves $\bar{a}^{(\ell)} \leq \bar{a}^{(\ell+1)}$.

• Finally, the convergence of this iterative process is proved by [7, Theorem 2.1], relying on Assumption 1. Hence, Algorithm `EVAL&APPROXOPTIMIZE` terminates at some iteration ℓ , with $\bar{a}_*^{(\ell)} \leq (1 + \tau)\bar{a}^{(\ell)}$, yielding the enclosure (15). \square

Algorithm 3 EVAL&APPROXOPTIMIZE(f, n, I, θ, τ)

Input: function f , $n \geq 0$, I , $\theta(\mathbf{a}, t)$ as in (9), $\tau > 0$.

Output: (\bar{a}, \mathbf{a}) solution of Problem (P) within accuracy τ .

▷ *Initialization*

- 1: $(\boldsymbol{\omega}^{(0)}, \mathbf{y}^{(0)}) \leftarrow \text{INIT}(n, I)$
 - 2: $(\bar{a}^{(0)}, \mathbf{a}^{(0)}) \leftarrow \text{SOLVEPRIMAL}(f, n, \theta, \boldsymbol{\omega}^{(0)})$
 - 3: $(\omega_*^{(0)}, \bar{a}_*^{(0)}) \leftarrow \text{FINDNEWINDEX}(f, n, I, \theta, \mathbf{a}^{(0)})$
 - 4: $\ell \leftarrow 0$
 - ▷ *Iterate while accuracy τ not reached*
 - 5: **while** $\bar{a}_*^{(\ell)} / \bar{a}^{(\ell)} > 1 + \tau$ **do**
 - 6: $(\boldsymbol{\omega}^{(\ell+1)}, \mathbf{y}^{(\ell+1)}) \leftarrow \text{EXCHANGE}(n, \theta, \boldsymbol{\omega}^{(\ell)}, \mathbf{y}^{(\ell)}, \omega_*^{(\ell)})$
 - 7: $(\bar{a}^{(\ell+1)}, \mathbf{a}^{(\ell+1)}) \leftarrow \text{SOLVEPRIMAL}(f, n, \theta, \boldsymbol{\omega}^{(\ell+1)})$
 - 8: $(\omega_*^{(\ell+1)}, \bar{a}_*^{(\ell+1)}) \leftarrow \text{FINDNEWINDEX}(f, n, I, \theta, \mathbf{a}^{(\ell+1)})$
 - 9: $\ell \leftarrow \ell + 1$
 - 10: **end while**
 - 11: **return** $(\bar{a}^{(\ell)}, \mathbf{a}^{(\ell)})$
-

Algorithm 4 INITPOINTS(n, I)

Input: $n \geq 0$, $I = [t_l, t_r]$.

Output: $\boldsymbol{\omega} \in \Omega^{n+2}$ and solution \mathbf{y} of Problem (D_m).

▷ *Initialize with Chebyshev nodes and Remez constraints*

- 1: **for** j **in** $[1..n+2]$ **do**
 - 2: $t_j \leftarrow \frac{t_l+t_r}{2} + \cos\left(\frac{(j-1)\pi}{n+1}\right) \frac{t_l-t_r}{2}$
 - 3: $\boldsymbol{\sigma}_j \leftarrow ((-1)^j, 0, \dots, 0)$
 - 4: $\omega_j \leftarrow (t_j, \boldsymbol{\sigma}_j)$
 - 5: **end for**
 - ▷ *Compute dual solution*
 - 6: Solve for \mathbf{y} the linear system $\sum_{j=1}^{n+2} y_j \boldsymbol{\alpha}(\omega_j) = \mathbf{z}$
 - 7: **return** $(\boldsymbol{\omega}, \mathbf{y})$
-

Algorithm 5 SOLVEPRIMAL($f, n, \theta, \boldsymbol{\omega}$)

Input: function f , $n \geq 0$, θ the evaluation error, $\boldsymbol{\omega} \in \Omega^{n+2}$.

Output: (\bar{a}, \mathbf{a}) solution of Problem (P_m) for $\boldsymbol{\omega}$.

1: Solve for $(\bar{a}, \mathbf{a}) \in \mathbb{R}^{n+2}$ the linear system:

$$\boldsymbol{\alpha}(\omega_j)^T(\bar{a}, \mathbf{a}) = c(\omega_j), \quad j \in [1..n+2]$$

2: **return** (\bar{a}, \mathbf{a})

Algorithm 6 FINDNEWINDEX($f, n, I, \theta, \mathbf{a}$)

Input: function f , $n \geq 0$, $I = [t_l, t_r]$, θ the evaluation error as in (9), coefficients $\mathbf{a} \in \mathbb{R}^{n+1}$.

Output: (ω_*, \bar{a}_*) with $\omega_* = (t_*, \boldsymbol{\sigma}_*) \in \Omega$

▷ *Compute maximal error in absolute value*

1: $t_* \leftarrow \arg \max_{t_l \leq t \leq t_r} |\mathbf{a}^T \boldsymbol{\pi}_0(t) - f(t)| + \sum_{i=1}^k |\mathbf{a}^T \boldsymbol{\pi}_i(t)|$

2: $\bar{a}_* \leftarrow \max_{t_l \leq t \leq t_r} |\mathbf{a}^T \boldsymbol{\pi}_0(t) - f(t)| + \sum_{i=1}^k |\mathbf{a}^T \boldsymbol{\pi}_i(t)|$

▷ *Reconstruct signature*

3: $\sigma_{*0} \leftarrow -\text{sign}(\mathbf{a}^T \boldsymbol{\pi}_0(t_*) - f(t_*))$

4: $\sigma_{*i} \leftarrow -\text{sign}(\mathbf{a}^T \boldsymbol{\pi}_i(t_*))$, $i \in [1..k]$

5: $\omega_* \leftarrow (t_*, \boldsymbol{\sigma}_*)$

6: **return** (ω_*, \bar{a}_*)

Algorithm 7 EXCHANGE($n, \theta, \boldsymbol{\omega}, \mathbf{y}, \omega_*$)

Input: $n \geq 0$, θ the evaluation error, $\boldsymbol{\omega} \in \Omega^{n+2}$, dual solution $\mathbf{y} \in \mathbb{R}^{n+2}$, new index $\omega_* \in \Omega$.

Output: new set $\boldsymbol{\omega}' \in \Omega^{n+2}$ and dual solution $\mathbf{y}' \in \mathbb{R}^{n+2}$.

1: Solve for $\boldsymbol{\gamma} \in \mathbb{R}^{n+2}$ the linear system:

$$\sum_{j=1}^{n+2} \gamma_j \boldsymbol{\alpha}(\omega_j) = \boldsymbol{\alpha}(\omega_*)$$

▷ *Exiting index*

2: $j_0 \leftarrow \arg \min \left\{ \frac{y_j}{\gamma_j} \mid \gamma_j > 0 \right\}$

▷ *Update dual solution*

3: $\tilde{y}_* \leftarrow \frac{y_{j_0}}{\gamma_{j_0}}$

4: $\tilde{y}_j \leftarrow y_j - \gamma_j \tilde{y}_*$, $j \in [1..n+2]$

5: $\{(\boldsymbol{\omega}'_j, y'_j)\} \leftarrow \{(\omega_j, \tilde{y}_j), j \in [1..n+2] - \{j_0\} \cup \{*\}\}$,

6: **return** $(\boldsymbol{\omega}', \mathbf{y}')$

5 Examples and conclusion

To illustrate Algorithm EVAL&APPROXOPTIMIZE, a tutorial example for Airy special function is discussed (further comparisons in binary32 case are given in Appendix .3). Then approximations with binary64 coefficients of arcsin are presented.

Example 3 (Airy function). Let Ai over $I = [-2, 2]$, approximated by a polynomial of degree $n = 6$, evaluated using the Horner scheme with $u = 2^{-12}$. The terms $\{\boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_7\}$ defining the evaluation error θ are given in Example 1. We fix a tolerance $\tau = 0.01$.

At iteration 0 (Figure 1), the points $t_j^{(0)}$ are initialized with the Chebyshev nodes and the signatures $\boldsymbol{\sigma}_j^{(0)}$ define a Remez-like system of linear equations on the coefficients of the polynomial (Figure 1d). Its solution $\mathbf{x}^{(0)} = (\bar{\mathbf{a}}^{(0)}, \mathbf{a}^{(0)})$ defines a polynomial $p^{(0)}(t) = \mathbf{a}^{(0)T} \boldsymbol{\pi}_0(t)$, whose approximation error is depicted in Figure 1a. It exhibits quasi-equioscillations indicating that $p^{(0)}$ is rather close to the degree-6 minimax approximation of Ai over I . However, the total error is more important near -2 and 2 (Figure 1b), due to the evaluation depicted in green. In particular, the algorithm detects the maximum error at $t_*^{(0)} = -2$ (in orange). Note that $t_1^{(0)}$ was already equal to -2 , but $\omega_1^{(0)} \neq \omega_*^{(0)}$ since the signatures are different. To perform the exchange, the dual solution is needed (Figure 1c). It is a positive

combination of Dirac measures supported on the finite set $\omega^{(0)}$.

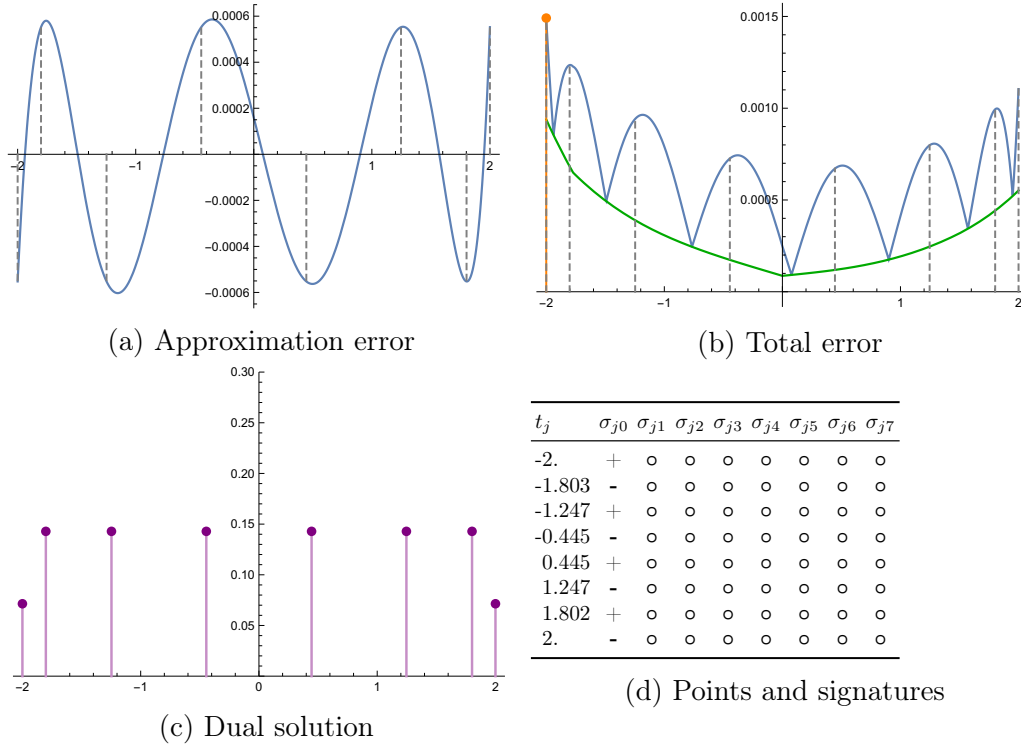
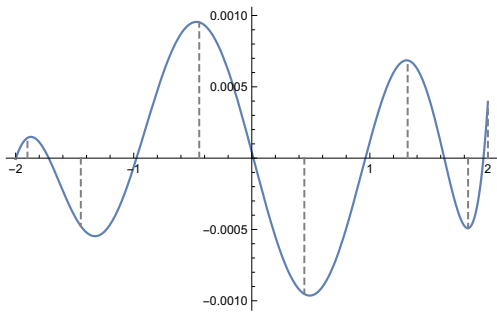
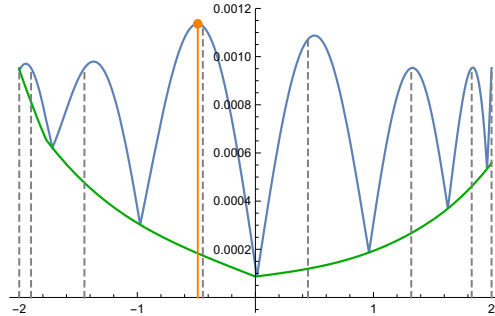


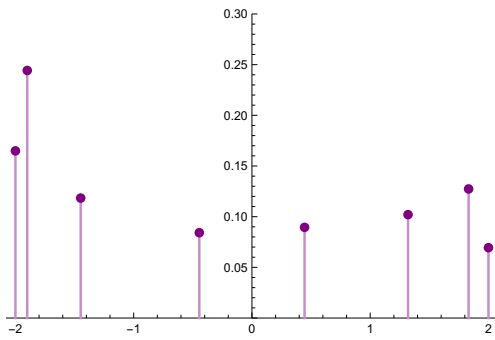
Figure 1: Approximation of A_i over $[-2, 2]$: iteration 0



(a) Approximation error



(b) Total error



(c) Dual solution

t_j	σ_{j0}	σ_{j1}	σ_{j2}	σ_{j3}	σ_{j4}	σ_{j5}	σ_{j6}	σ_{j7}
-2.	+	-	-	+	+	+	+	-
-1.9	-	-	-	+	+	+	+	-
-1.448	+	-	-	+	+	+	-	-
-0.445	-	o	o	o	o	o	o	o
0.445	+	o	o	o	o	o	o	o
1.32	-	-	+	+	-	-	+	-
1.832	+	-	-	+	-	-	+	-
2.	-	-	-	+	-	-	+	-

(d) Points and signatures

Figure 2: Approximation of A_i over $[-2, 2]$: iteration 6

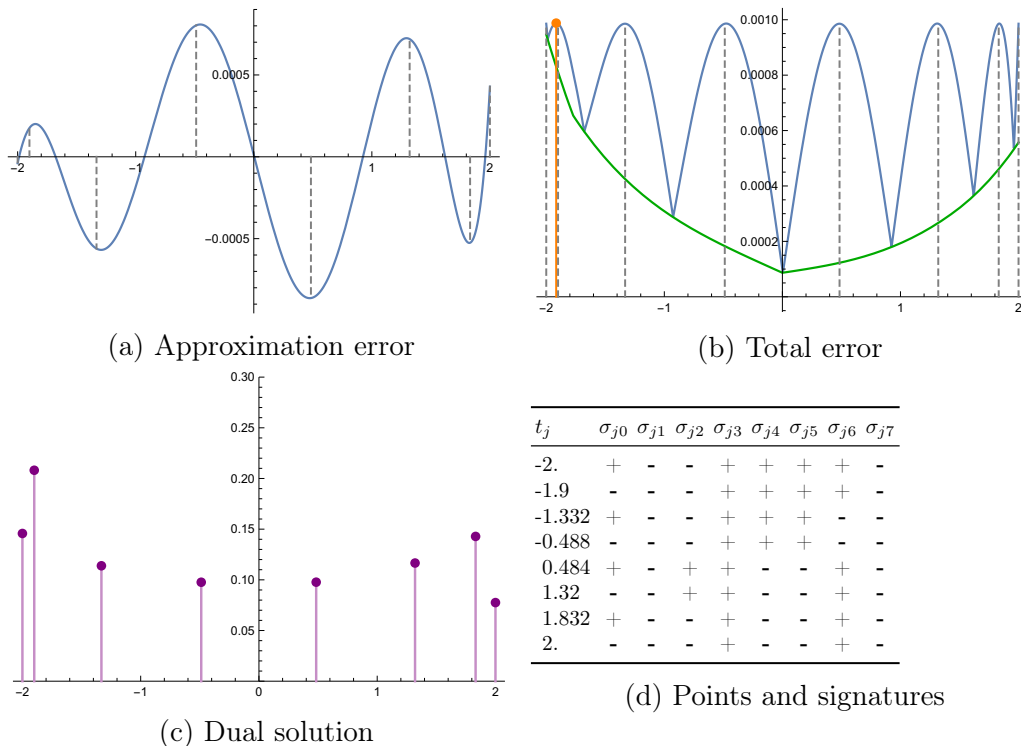


Figure 3: Approximation of A_i over $[-2, 2]$: iteration 9

Moving forward to iteration 6 (Figure 2), the total error is more balanced, though still not optimal. Both the signatures and the approximation error are now completely different from the Remez solution.

Eventually, the algorithm stops at iteration 9 (Figure 3). Indeed, the maximum total error $\bar{a}_*^{(9)}$ (in orange) is less than 1% higher than the error $\bar{a}^{(9)}$ over the discrete set $\omega^{(9)}$. Note that the total error reaches its maximum at $n + 2 = 8$ points. This became possible by unbalancing the approximation error, namely reducing the amplitude of the oscillations near -2 and 2 , at the cost of higher oscillations in the middle of I .

Example 4 (Arcsine function). Consider $f = \arcsin$, over the interval $I = [0.75; 1]$. This example is particularly insightful because f is ill-conditioned over I , and also because its values are close to 1, so absolute or relative error treatment is similar. Firstly, assume that argument reduction techniques are not available (this can be the case for any ill-conditioned function, known only through sampling via a black-box approach), operations are in binary64 and an approximation polynomial of degree 20 with binary64 coefficients is searched for.

Figure 4(a) shows the absolute approximation error between \arcsin and the best approximation polynomial p_{Remez} with *real* coefficients, for which

$\max_I |f - p_{\text{Remez}}| \leq 0.00496$. Then, the coefficients of p_{Remez} are directly rounded to binary64, which results in $\max_I |f - \text{roundCoeff}(p_{\text{Remez}})| \simeq 0.004961.15 \cdot 10^{10}$ (cf. Figure 4(b)). This is due to the high magnitude of the coefficients. To search for better coefficients, one could use the proficient FPMinimax routine of Sollya [11], which optimizes on the coefficient space of representable FP numbers. The approximation error in this case is shown in Figure 4(c) and one has $\max_I |f - p_{\text{FPMinimax}}| \simeq 1.58$. Finally, in Figure 4(d) we provide the plot of the total error between f and the obtained Eval&ApproxOpt polynomial (with binary64 coefficients), which provides the best bound in this case $\max_I |f - p_{\text{E\&A}}| \simeq 0.0081$. One can then actually prove with Gappa [14] that the Horner scheme evaluation of $p_{\text{E\&A}}$ has an absolute error less than 0.00035.

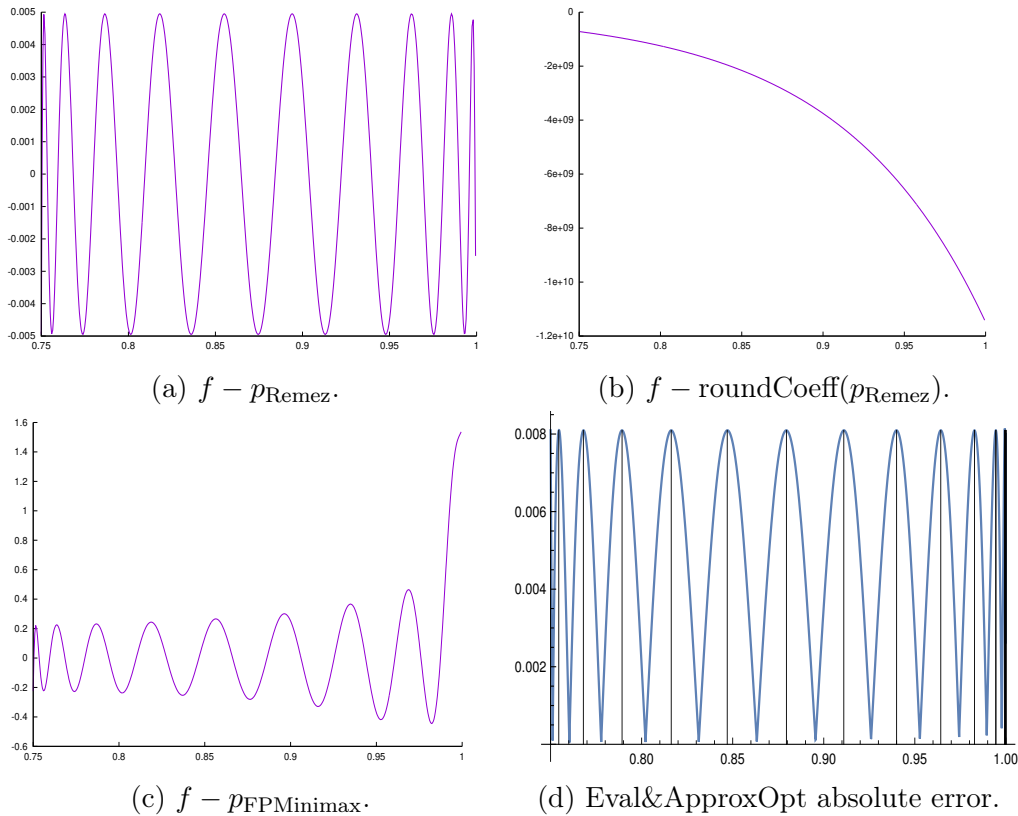


Figure 4: Error plots for different approximation polynomials of degree 20 for $f = \arcsin$ over $I = [0.75; 1]$.

Secondly, suppose that I is re-centered to $I_c = [-0.125; 0.125]$ and the other parameters remain unchanged. Thus, the function to be approximated is $f = \arcsin(x + c)$, where c is the middle point of I . In this case, it turns

out that all the methods presented above provide almost the same accuracy. In Figure 5(a), the plot of the absolute error between f and p_{Remez} is given: $\max_{I_c} |f - p_{Remez}| \leq 0.00496$; at the drawing scale, the plot is the same for the other errors and one has that $\max_{I_c} |f - \text{roundCoeff}(p_{Remez})|$, $\max_{I_c} |f - p_{\text{FPMinimax}}|$ and $\max_{I_c} |f - p_{\text{E\&A}}|$ are all very close in magnitude $\simeq 0.00496$. Similarly, one can prove with Gappa that the Horner scheme evaluation has an absolute error of $7.16227 \cdot 10^{-11}$ for every above-mentioned polynomial. Hence in this case, our algorithm does not improve the error bound. Yet, it shows (at least numerically) that if only binary64 computations and a Horner scheme are used, there is no other polynomial which performs better w.r.t. both evaluation and approximation errors. For completeness, in Figure 5(b) we show that FPMinimax and Eval&ApproxOpt polynomials are different. However, this difference is very small and this does not influence the number of correct bits provided $-\log_2(0.00496) \simeq 7.65$.

Finally, we mention that when other argument reduction techniques exist, and when the evaluation error is not an issue (very small intervals around zero), the FPMinimax method still provides better tuned FP coefficients. So this opens the question for several future extensions. A mixed-integer linear programming problem could be formulated in the provided optimization framework. However, a similar exchange procedure in this case is not obvious. Concerning precisions of the coefficients and operations, they can be variable, as mentioned in Section 2, but a more detailed study is needed to eventually take into account higher order error terms for the error estimation formula. The polynomial coefficients stay linear in such a formula, so the algorithm presented can be straightforwardly used in such a case. In addition, this formula directly allows for the estimation of evaluation errors for other numerical schemes and eventually polynomial bases [2] for which a practical study is necessary.

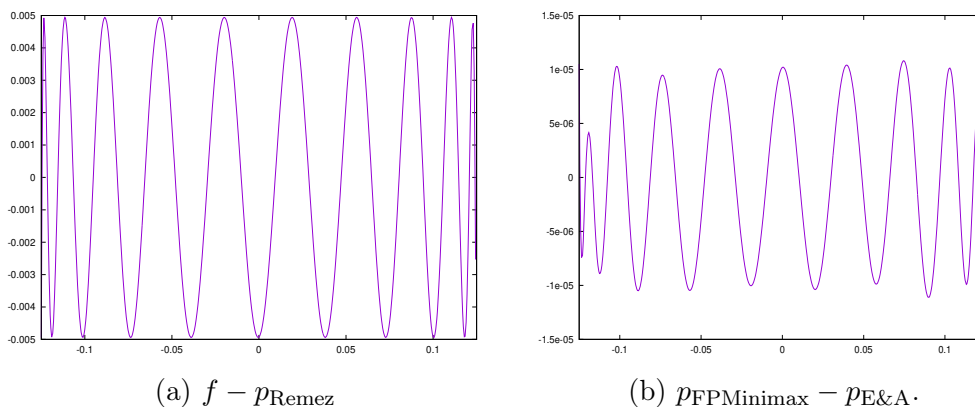


Figure 5: Error plots for different approximation polynomials of degree 20 for $f = \arcsin(x + c)$ over $I_c = [-0.125, 0.125]$.

References

- [1] M. Abramowitz and I. A. Stegun. *Handbook of mathematical functions with formulas, graphs, and mathematical tables*, volume 55 of *National Bureau of Standards Applied Mathematics Series*. Courier Corp., 1964.
- [2] R. Barrio, H. Jiang, and S. Serrano. A general condition number for polynomials. *SIAM J. Numer. Anal.*, 51(2):1280–1294, 2013.
- [3] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, New York, NY, USA, 2004.
- [4] T. Braconnier and P. Langlois. From rounding error estimation to automatic correction with automatic differentiation. In *Automatic differentiation of algorithms*, pages 351–357. Springer, 2002.
- [5] N. Brisebarre and S. Chevillard. Efficient polynomial L^∞ approximations. In *18th IEEE Symposium on Computer Arithmetic (ARITH-18)*, pages 169–176, Montpellier, France, 2007.
- [6] N. Brisebarre, J.-M. Muller, and A. Tisserand. Computing machine-efficient polynomial approximations. *ACM Trans. Math. Software*, 32(2):236–256, June 2006.
- [7] C. Carasso. *L’algorithme d’échange en optimisation convexe*. PhD thesis, Université Joseph-Fourier-Grenoble I, 1973.

- [8] C. Carasso and P. J. Laurent. Un algorithme général pour l'approximation au sens de Tchebycheff de fonctions bornées sur un ensemble quelconque. In *Approximation Theory*, number 556 in Lecture Notes in Mathematics, pages 99–121. Springer Berlin Heidelberg, 1976.
- [9] B. L. Chalmers. The Remez exchange algorithm for approximation with linear restrictions. *Trans. Amer. Math. Soc.*, 223:103–131, 1976.
- [10] S. Chevillard, J. Harrison, M. Joldeş, and C. Lauter. Efficient and accurate computation of upper bounds of approximation errors. *Theoret. Comput. Sci.*, 412(16):1523–1543, 2011.
- [11] S. Chevillard, M. Joldeş, and C. Lauter. Sollya: An environment for the development of numerical codes. In *International Congress on Mathematical Software*, pages 28–31. Springer, 2010.
- [12] C. Daramy-Loirat, D. Defour, F. de Dinechin, M. Gallet, N. Gast, C. Q. Lauter, and J.-M. Muller. CR-LIBM, a library of correctly-rounded elementary functions in double-precision. Technical report, LIP Laboratory, Arenal team, Dec. 2006.
- [13] E. Darulova and V. Kuncak. Towards a compiler for reals. *ACM Trans. Program. Lang. Syst.*, 39(2):8:1–8:28, Mar. 2017.
- [14] M. Daumas and G. Melquiond. Certification of bounds on expressions involving rounded operators. *ACM Trans. Math. Software*, 37(1):Art. 2, 20, 2010.
- [15] F. de Dinechin, C. Q. Lauter, and J.-M. Muller. Fast and correctly rounded logarithms in double-precision. *Theor. Inform. Appl.*, 41(1):85–102, 2007.
- [16] R. Halsey and F. Patrick. *Real Analysis*. Prentice Hall, 2010.
- [17] IEEE Computer Society. *IEEE Standard for Floating-Point Arithmetic*. IEEE Standard 754-2008, Aug. 2008. Available at <http://ieeexplore.ieee.org/servlet/opac?punumber=4610933>.
- [18] O. Kupriianova and C. Lauter. Metalibm: A mathematical functions code generator. In *International Congress on Mathematical Software*, pages 713–717. Springer, 2014.
- [19] C. Lauter and M. Mezzarobba. Semi-automatic floating-point implementation of special functions. In J.-M. Muller, A. Tisserand, and J. Villalba, editors, *ARITH 22*, Lyon, France, 2015. IEEE.

- [20] C. Q. Lauter. *Arrondi Correct de Fonctions Mathématiques*. PhD thesis, ÉNS de Lyon, Lyon, France, Oct. 2008.
- [21] J.-M. Muller. *Elementary Functions, Algorithms and Implementation*. Birkhäuser, Boston, 3rd edition, 2016.
- [22] J.-M. Muller, N. Brunie, F. de Dinechin, C.-P. Jeannerod, M. Joldes, V. Lefèvre, G. Melquiond, N. Revol, and S. Torres. *Handbook of Floating-Point Arithmetic*. Birkhäuser Boston, 2018.
- [23] J. Oliver. Rounding error propagation in polynomial evaluation schemes. *J. Comput. Appl. Math.*, 5(2):85–97, 1979.
- [24] R. Reemtsen and J.-J. Rückmann. *Semi-Infinite Programming*, volume 25. Springer Science & Business Media, 1998.
- [25] A. Rocca, V. Magron, and T. Dang. Certified Roundoff Error Bounds using Bernstein Expansions and Sparse Krivine-Stengle Representations. In *24th IEEE Symposium on Computer Arithmetic*. IEEE, 2017.
- [26] A. Shapiro. Semi-infinite programming, duality, discretization and optimality conditions. *Optimization*, 58(2):133–161, 2009.
- [27] A. Solovyev, M. S. Baranowski, I. Briggs, C. Jacobsen, Z. Rakamarić, and G. Gopalakrishnan. Rigorous estimation of floating-point round-off errors with symbolic taylor expansions. *ACM Trans. Program. Lang. Syst.*, 41(1):2:1–2:39, Dec. 2018.
- [28] A. Tisserand. High-performance hardware operators for polynomial evaluation. *International Journal of High Performance Systems Architecture (IJHPSA)*, 1(1):14–23, 2007.
- [29] G. A. Watson. The calculation of best restricted approximations. *SIAM J. Numer. Anal.*, 11(4):693–699, 1974.

.1 Complementary proofs for the exchange algorithm

This appendix provides auxiliary lemmas for the proof of Theorem 2.

Lemma 1 (Correctness of INIT). *INIT(n, I) computes $\boldsymbol{\omega} = \{\omega_j\}_{j=1}^{n+2} \in \Omega^{n+2}$ and $\mathbf{y} \in \mathbb{R}^{n+2}$ satisfying:*

- $\{\boldsymbol{\alpha}(\omega_j)\}_{j=1}^{n+2}$ is a basis of \mathbb{R}^{n+2} ;
- \mathbf{y} is the optimal solution of Problem (D_m) for $\boldsymbol{\omega}$;
- $y_j > 0$ for all $j \in [1..n+2]$.

Note that Algorithm INIT essentially initializes the problem with Chebyshev nodes for heuristic efficiency and signatures corresponding to the classical Remez algorithm, without the evaluation error term.

Proof. Let $\mathbf{A}(\boldsymbol{\omega})$ denote the $(n+2)$ square matrix whose columns are the $\boldsymbol{\alpha}(\omega_j)$. Since $\boldsymbol{\sigma}_j = ((-1)^j, 0, \dots, 0)$, we have

$$\mathbf{A}(\boldsymbol{\omega}) = \begin{pmatrix} 1 & \dots & \dots & 1 \\ -1 & \dots & \dots & (-1)^{n+2} \\ -t_1 & \dots & \dots & (-1)^{n+2}t_{n+2} \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ -t_1^n & \dots & \dots & (-1)^{n+2}t_{n+2}^n \end{pmatrix}.$$

First, we prove the existence of a feasible point \mathbf{y} in Problem (D_m) for $\boldsymbol{\omega}$, that is $\mathbf{A}(\boldsymbol{\omega})\mathbf{y} = \mathbf{z}$ and $\mathbf{y} \geq 0$. From Farkas' lemma [26], if such a \mathbf{y} does not exist, then there exists $\mathbf{x} = (\bar{a}, \mathbf{a}) \in \mathbb{R}^{n+2}$ s.t. $\mathbf{z}^T \mathbf{x} = \bar{a} < 0$ and $\mathbf{A}(\boldsymbol{\omega})^T \mathbf{x} \geq 0$, that is

$$\bar{a} + (-1)^j \mathbf{a}^T \boldsymbol{\pi}_0(t_j) \geq 0, \quad j \in [1..n+2].$$

Since $\bar{a} < 0$, this implies that $\text{sign}(\mathbf{a}^T \boldsymbol{\pi}_0(t_j)) = (-1)^j$. But $\mathbf{a}^T \boldsymbol{\pi}_0(t)$ is a polynomial of degree at most n , hence it cannot strictly change signs $n+2$ times. Consequently, Problem (D_m) has a feasible point \mathbf{y} .

Now, suppose that the columns of $\mathbf{A}(\boldsymbol{\omega})$ are not linearly independent, or that $y_j = 0$ for some j . Both cases imply that there exists $J \subset [1..n+2]$ of size $n+1$, and $\tilde{\mathbf{y}} \in \mathbb{R}^{n+1}$ s.t. $\sum_{j \in J} \tilde{y}_j \boldsymbol{\alpha}(\omega_j) = \mathbf{z}$. In particular, by canceling the first component, the family $\{\boldsymbol{\pi}_0(t_j)\}_{j \in J}$ is linearly dependent. But the Vandermonde determinant of this system cannot vanish since the t_j

are pairwise distinct. Therefore, $\{\boldsymbol{\alpha}(\omega_j)\}_{j=1}^{n+2}$ is a basis of \mathbb{R}^{n+2} , and $y_j > 0$ for all j .

Finally, since $\mathbf{A}(\boldsymbol{\omega})$ is invertible, \mathbf{y} is the unique optimal feasible point of (D_m) . \square

Lemma 2 (Correctness of SOLVEPRIMAL). *If $\{\boldsymbol{\alpha}(\omega_j)\}_{j=1}^{n+2}$ is a basis of \mathbb{R}^{n+2} and Problem (D_m) for $\boldsymbol{\omega}$ is feasible, then SOLVEPRIMAL($f, n, \theta, \boldsymbol{\omega}$) computes the optimal solution $\mathbf{x} = (\bar{a}, \mathbf{a})$ of Problem (P_m) .*

Proof. Algorithm SOLVEPRIMAL computes the solution $\mathbf{x} \in \mathbb{R}^{n+2}$ of $\boldsymbol{\alpha}(\omega_j)^T \mathbf{x} = c(\omega_j)$ for $j \in [1..n+2]$. We show that \mathbf{x} is the optimal solution of Problem (P_m) for $\boldsymbol{\omega}$.

Let $\tilde{\mathbf{x}}$ be any feasible point in (P_m) . Since the dual Problem (D_m) for $\boldsymbol{\omega}$ is feasible, there exists $\mathbf{y} \geq 0$ s.t. $\mathbf{z} = \sum_{j=1}^{n+2} y_j \boldsymbol{\alpha}(\omega_j)$. Then

$$\mathbf{z}^T \tilde{\mathbf{x}} = \sum_{j=1}^{n+2} y_j \boldsymbol{\alpha}(\omega_j)^T \tilde{\mathbf{x}} \geq \sum_{j=1}^{n+2} y_j c(\omega_j) = \sum_{j=1}^{n+2} y_j \boldsymbol{\alpha}(\omega_j)^T \mathbf{x} = \mathbf{z}^T \mathbf{x},$$

thereby establishing optimality of \mathbf{x} . \square

Lemma 3 (Correctness of FINDNEWINDEX). *Given $\mathbf{x} = (\bar{a}, \mathbf{a})$, FINDNEWINDEX($f, n, I, \theta, \mathbf{a}$) computes ω_* and \bar{a}_* corresponding to the most violated constraint:*

$$\begin{aligned} \omega_* &= \arg \max_{\omega \in \Omega} (c(\omega) - \boldsymbol{\alpha}(\omega)^T \mathbf{x}), \\ \bar{a}_* - \bar{a} &= \max_{\omega \in \Omega} (c(\omega) - \boldsymbol{\alpha}(\omega)^T \mathbf{x}). \end{aligned}$$

Proof. We have

$$\begin{aligned} & \max_{\omega \in \Omega} (c(\omega) - \boldsymbol{\alpha}(\omega)^T \mathbf{x}) \\ &= \max_{\omega=(t, \boldsymbol{\sigma}) \in \Omega} \left(\sigma_0 (f(t) - \mathbf{a}^T \boldsymbol{\pi}_0(t)) - \sum_{i=1}^k \sigma_i \mathbf{a}^T \boldsymbol{\pi}_i(t) \right) - \bar{a} \\ &= \max_{t_l \leq t \leq t_r} \left(|\mathbf{a}^T \boldsymbol{\pi}_0(t) - f(t)| + \sum_{i=1}^k |\mathbf{a}^T \boldsymbol{\pi}_i(t)| \right) - \bar{a}. \end{aligned}$$

Therefore, by computing t_* (line 1) and $\boldsymbol{\sigma}_*$ (lines 3-4), Algorithm FINDNEWINDEX ensures that $\omega_* := (t_*, \boldsymbol{\sigma}_*)$ is the index of the most violated constraint, with

$$c(\omega_*) - \boldsymbol{\alpha}(\omega_*)^T \mathbf{x} = \bar{a}_* - \bar{a} \geq 0.$$

\square

Lemma 4 (Correctness of EXCHANGE). *If $\{\boldsymbol{\alpha}(\omega_j)\}_{j=1}^{n+2}$ is a basis of \mathbb{R}^{n+2} and \mathbf{y} the optimal solution of Problem (D_m) for $\boldsymbol{\omega}$, then $\text{EXCHANGE}(n, \theta, \boldsymbol{\omega}, \mathbf{y}, \omega_*)$ computes new $\boldsymbol{\omega}' \in \Omega^{n+2}$ and $\mathbf{y}' \in \mathbb{R}^{n+2}$ such that:*

- $\boldsymbol{\omega}' = \{\omega_j\}_{j \in J} \cup \{\omega_*\}$ for a subset $J \subset [1 \dots n+2]$ of size $n+1$;
- $\{\boldsymbol{\alpha}(\omega_j)\}_{j=1}^{n+2}$ is a basis of \mathbb{R}^{n+2} ;
- \mathbf{y}' is the optimal solution of Problem (D_m) for $\boldsymbol{\omega}'$.

Proof. In order to increase the objective value in the dual Problem (D) , a measure supported on $\boldsymbol{\omega} \cup \{\omega_*\}$ is looked for, requiring nonnegative coefficients \tilde{y}_j for $j \in [1 \dots n+2] \cup \{*\}$ s.t.

$$\mathbf{z} = \sum_{j=1}^{n+2} \tilde{y}_j \boldsymbol{\alpha}(\omega_j) + \tilde{y}_* \boldsymbol{\alpha}(\omega_*) = \sum_{j=1}^{n+2} (\tilde{y}_j + \gamma_j \tilde{y}_*) \boldsymbol{\alpha}(\omega_j),$$

while maximizing \tilde{y}_* . But $\{\boldsymbol{\alpha}(\omega_j)\}_{j=1}^{n+2}$ being a basis of \mathbb{R}^{n+2} implies

$$\tilde{y}_j + \gamma_j \tilde{y}_* = y_j, \quad j \in [1 \dots n+2].$$

The nonnegativity constraints on the \tilde{y}_j induces the choice of the exiting index ω_{j_0} (line 2) and the values of the \tilde{y}_j (lines 3-4). Note that the first line of the linear system in line 1 says that the coefficients γ_j sum to 1. Hence, at least one of them is strictly positive, so that j_0 exists (line 2), though it is not necessarily unique.

Finally, $\{\boldsymbol{\alpha}(\omega_j)\}_{j \in [1 \dots n+2] - \{j_0\} \cup \{*\}}$ remains a basis of \mathbb{R}^{n+2} since $\gamma_{j_0} \neq 0$, that is $\boldsymbol{\alpha}(\omega_*)$ is not in the linear subspace spanned by $\{\boldsymbol{\alpha}(\omega_j)\}_{j \in [1 \dots n+2] - \{j_0\}}$. \square

Lemma 5. *The total error $\bar{a}^{(\ell)}$ computed over the discrete set $\boldsymbol{\omega}^{(\ell)}$ increases at each iteration.*

Proof. Let $\tilde{y}_*^{(\ell)}$ and $\gamma_j^{(\ell)}$ denote the variables \tilde{y}_* and γ_j for $j \in [1 \dots n+2]$ in $\text{EXCHANGE}(n, \theta, \boldsymbol{\omega}^{(\ell)}, \mathbf{y}^{(\ell)}, \omega_*^{(\ell)})$. By strong duality in linear programming, $\bar{a}^{(\ell)}$ is also the objective value of the optimal solution $\mathbf{y}^{(\ell)}$ in the discretized dual problem (D_{n+2}) for $\boldsymbol{\omega}^{(\ell)}$. Hence, by writing

$$\begin{aligned} \bar{a}^{(\ell)} &= \sum_{j=1}^{n+2} y_j^{(\ell)} c(\omega_j^{(\ell)}), \quad \text{and} \\ \bar{a}^{(\ell+1)} &= \sum_{j=1}^{n+2} y_j^{(\ell+1)} c(\omega_j^{(\ell+1)}) \\ &= \sum_{j=1}^{n+2} \left(y_j^{(\ell)} - \gamma_j^{(\ell)} \tilde{y}_*^{(\ell)} \right) c(\omega_j^{(\ell)}) + \tilde{y}_*^{(\ell)} c(\omega_*^{(\ell)}), \end{aligned}$$

we have

$$\begin{aligned}
\bar{a}^{(\ell+1)} - \bar{a}^{(\ell)} &= \tilde{y}_*^{(\ell)} \left(c(\omega_*^{(\ell)}) - \sum_{j=1}^{n+2} \gamma_j^{(\ell)} c(\omega_j^{(\ell)}) \right) \\
&= \tilde{y}_*^{(\ell)} \left(c(\omega_*^{(\ell)}) - \sum_{j=1}^{n+2} \gamma_j^{(\ell)} \boldsymbol{\alpha}(\omega_j^{(\ell)})^T \mathbf{x}^{(\ell)} \right) \\
&= \tilde{y}_*^{(\ell)} (c(\omega_*^{(\ell)}) - \boldsymbol{\alpha}(\omega_*^{(\ell)})^T \mathbf{x}^{(\ell)}) \geq 0,
\end{aligned}$$

because $\tilde{y}_*^{(\ell)} \geq 0$ and the constraint $\boldsymbol{\alpha}(\omega_*^{(\ell)})^T \mathbf{x}^{(\ell)} \geq c(\omega_*^{(\ell)})$ is violated at the beginning of iteration ℓ . \square

.2 Optimizing with the relative error

When considering the relative error in place of the absolute error, Problem (P_{general}) is replaced by

$$\min_{\substack{a_i \in \mathbb{R}, \\ i \in [0..n]}} \max_{t \in I} \left(\frac{|f(t) - p(t)| + |\tilde{p}(t) - p(t)|}{|f|} \right) \quad (P_{\text{general}}^{\text{rel}})$$

where f is assumed to be strictly positive over I . This is equivalent to multiplying the error variable \bar{a} by f in Problem (P), yielding the following new linear SIP formulation

$$\begin{aligned}
\min_{\mathbf{x} \in \mathbb{R}^d} \quad & \mathbf{z}^T \mathbf{x} \\
\text{s.t.} \quad & \tilde{\boldsymbol{\alpha}}(\omega)^T \mathbf{x} \geq c(\omega), \quad \omega \in \Omega,
\end{aligned} \quad (P^{\text{rel}})$$

with

$$\begin{aligned}
\mathbf{x} &= (\bar{a}, \mathbf{a}) \in \mathbb{R}^{n+2}, \quad \mathbf{z} = (1, 0, \dots, 0) \in \mathbb{R}^{n+2}, \\
\tilde{\boldsymbol{\alpha}}(t, \sigma_0, \dots, \sigma_k) &= (f(t), \sigma_0 \boldsymbol{\pi}_0^T(t) + \sum_{i=1}^k \sigma_i \boldsymbol{\pi}_i^T(t))^T \in \mathbb{R}^{n+2}, \\
\mathfrak{S} &= \{-1, 0, 1\}^{k+1}, \quad \omega = (t, \sigma_0, \dots, \sigma_k) \in \Omega := I \times \mathfrak{S}.
\end{aligned} \quad (16)$$

Therefore, replacing $\boldsymbol{\alpha}(\omega)$ by $\tilde{\boldsymbol{\alpha}}(\omega)$ in EVAL&APPROXOPTIMIZE and its subroutines provides an algorithm to compute the best degree- n polynomial w.r.t. both approximation and evaluation errors in relative setting.

.3 Airy function (continued)

The Airy function Ai is a special function frequently used in theoretical physics. In particular, some applications require to compute $\text{Ai}(t)$ for possibly large negative values of t , where the function exhibits a highly oscillatory

behavior (see Figure 6a). However, contrary to elementary functions, there exists no simple argument reduction for Ai . Therefore, one polynomial approximation is needed for each interval of the domain subdivision, and these intervals cannot be assumed to be small. Hence, controlling the evaluation error is essential.

We here consider the approximation of the Airy function over $I = [-4, 0]$ with polynomials evaluated using the Horner scheme with binary32 floating-point numbers, that is, single precision ($u = 2^{-24}$). Figure 6b shows the total error (approximation and evaluation) of the degree- n minimax polynomial (in blue), in function of n . Specifically, the total error starts to decrease when n increases, thanks to the improvement of the approximation error. However, at some point (here $n = 9$), the total error starts again to increase, due to the evaluation error of higher degree polynomials. On the contrary, the evaluation error of polynomials obtained with EVAL&APPROXOPTIMIZE (in yellow) continues decreasing to some asymptotic threshold (around $5 \cdot 10^{-6}$). In such a case, reducing the evaluation error by using a higher degree can be more efficient than increasing the floating-point precision.

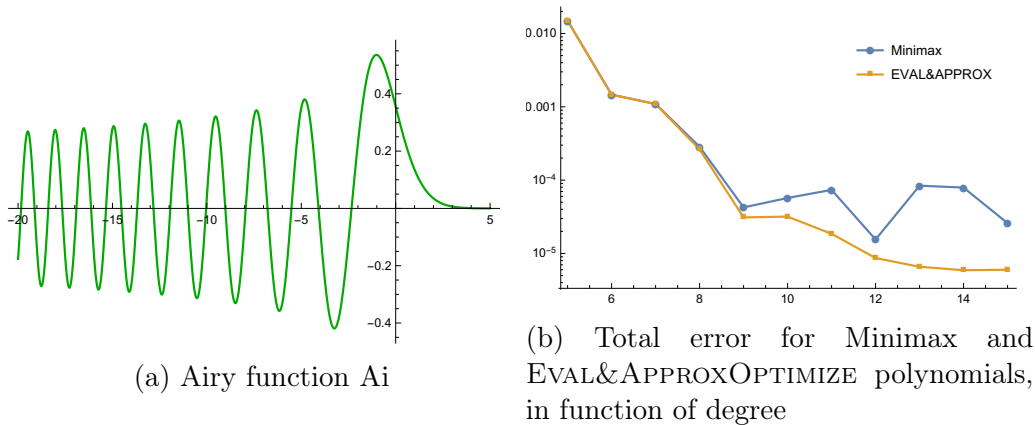


Figure 6: Approximation of Ai over $[-4, 0]$