



**HAL**  
open science

## LLL and polynomial approximation

Nicolas Brisebarre, Sylvain Chevillard

► **To cite this version:**

Nicolas Brisebarre, Sylvain Chevillard. LLL and polynomial approximation. LLL+25, Jun 2007, Caen, France. hal-04732267

**HAL Id: hal-04732267**

**<https://inria.hal.science/hal-04732267v1>**

Submitted on 16 Oct 2024

**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.



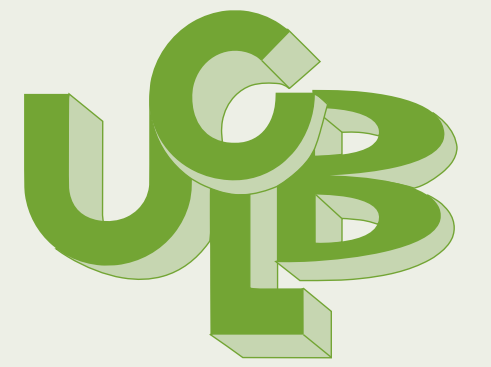
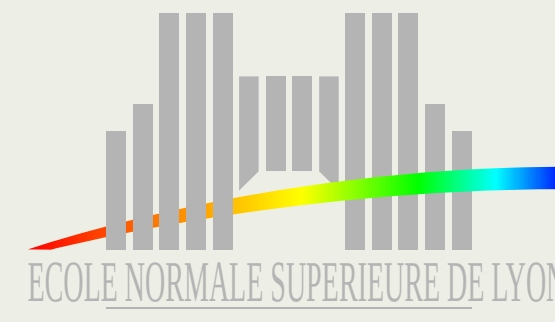
Distributed under a Creative Commons Attribution 4.0 International License

# LLL and polynomial approximation

Nicolas Brisebarre

Sylvain Chevillard

Arénaire project



## The problem

In practice, to implement mathematical functions, one uses a polynomial approximation to the function. The coefficients of the polynomial will be stored into floating-point variables.

A naive method is to compute the best polynomial with real coefficients  $p^*$  and to round each coefficient to the nearest floating-point number, thus leading to a polynomial  $\hat{p}$ . The polynomial  $\hat{p}$  may be quite inaccurate. On figure 1 in the *Example* box, the graph of the absolute errors  $\hat{p} - f$  and  $p^* - f$  in a real case are shown. The method described in this poster gives much better results (figure 2).

We address here the problem of finding a very good polynomial  $\bar{p}$  with floating-point coefficients, to approximate a continuous function on an interval.

**Lattices** are a *key tool* in the method we propose.

## IEEE-754 standard

The standard IEEE-754 defines how computers manipulate real numbers. The real numbers are approximated by the so-called *floating-point* numbers. Let us fix  $t \in \mathbb{N}$ .

The set of floating-point numbers of precision  $t$  is defined by  $\mathcal{F}_t = \{m \cdot 2^{e-t}, (m, e) \in \mathbb{Z}^2\}$   $m$  being written with exactly  $t$  bits.

In other words  $x \in \mathcal{F}_t$  is a number of the form  $x = 0.1m_2m_3 \dots m_t \cdot 2^e$  where  $m_i \in \{0, 1\}$ .

## Polynomial approximation

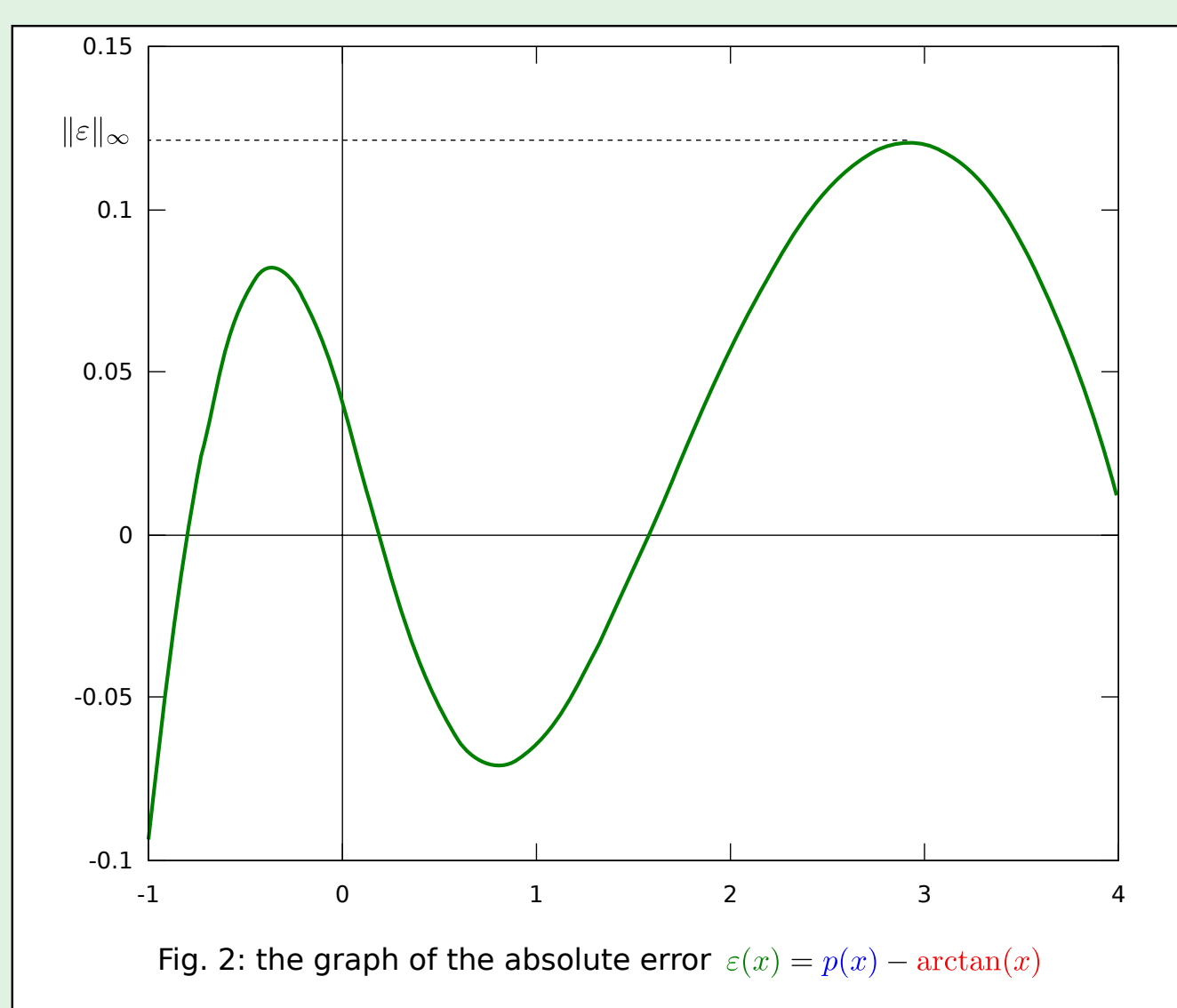
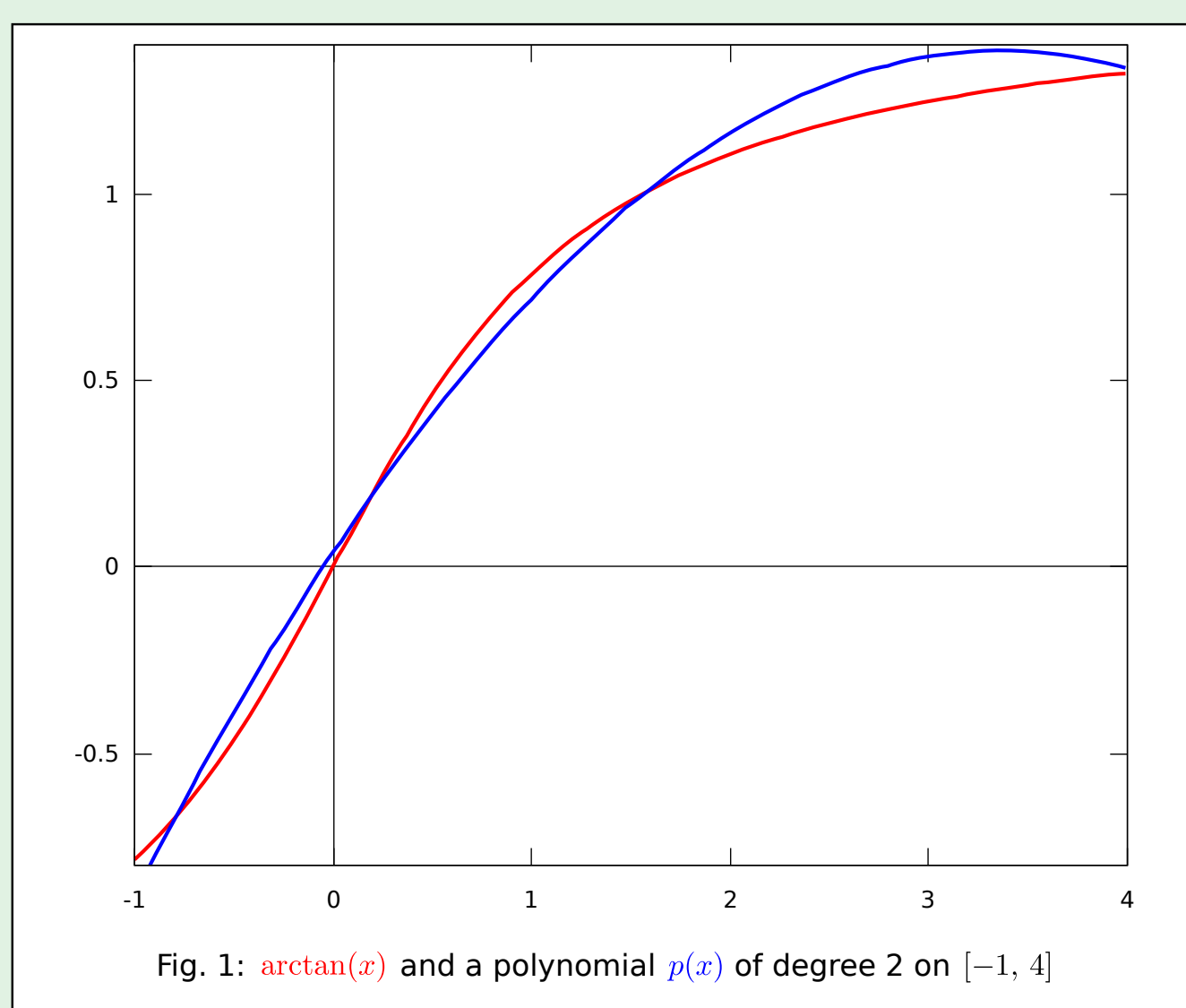
To evaluate a mathematical function  $f$  on several points of an interval, one generally replaces it by an approximation easier to evaluate.

A good choice consists in choosing a polynomial as an approximant because only additions and multiplications are required to evaluate it.

Using  $p(x)$  as an approximation to  $f(x)$  leads to an absolute error  $\varepsilon(x) = p(x) - f(x)$  (one may also consider the relative error defined as  $\varepsilon(x)/f(x)$ ).

The worst approximation quality corresponds to the maximum of the error function  $|\varepsilon|$  on the interval  $[a, b]$ . This is called the *infinite norm* of  $\varepsilon$  on  $[a, b]$ . The figures show the example of a polynomial of degree 2 approximating the function  $\arctan$  in  $[-1, 4]$ .

The *best approximation* problem consists in finding the polynomial with real coefficients leading to the smallest error in infinite norm. The existence of this polynomial and its uniqueness are ensured by theorems by Chebyshev and De la Vallée Poussin. Remez gave an iterative process that converges toward this polynomial.



## Formalization

Our goal: find  $p$  with floating-point coefficients that approximates  $f$  well.

Thus, each coefficient  $a$  of  $p$  must be chosen of the form  $a = m \cdot 2^e$  where  $m$  and  $e$  are two unknown integers.

$e$  is the exponent of the coefficient, thus corresponds to its order of magnitude. We assume that it is the same as the corresponding one in  $p^*$ .

We obtain now a problem that we will solve with lattice reduction techniques: find  $p$  that approximates  $f$  well and with the following form:

$$m_0 \cdot 2^{e_0} + m_1 \cdot 2^{e_1} X + \dots + m_n \cdot 2^{e_n} X^n.$$

## Resolution

Our goal: find  $p$  that approximates  $f$  well and with the following form:

$$m_0 \cdot 2^{e_0} + m_1 \cdot 2^{e_1} X + \dots + m_n \cdot 2^{e_n} X^n.$$

We use the idea of interpolation: we force the polynomial to be close to the function at some chosen points. Let  $(x_0, \dots, x_n)$  be  $n+1$  points in  $[a, b]$ . We search  $m_0, \dots, m_n$  such that for all  $i$

$$p(x_i) = m_0 \cdot 2^{e_0} + m_1 \cdot 2^{e_1} x_i + \dots + m_n \cdot 2^{e_n} x_i^n \simeq f(x_i).$$

Rewritten with vectors: 
$$m_0 \begin{pmatrix} 2^{e_0} \\ 2^{e_0} \\ \vdots \\ 2^{e_0} \end{pmatrix} + \dots + m_n \begin{pmatrix} 2^{e_n} \cdot x_0^n \\ 2^{e_n} \cdot x_1^n \\ \vdots \\ 2^{e_n} \cdot x_n^n \end{pmatrix} \simeq \begin{pmatrix} f(x_0) \\ f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}.$$

$\Gamma$  of the form  $\mathbb{Z}\vec{b}_0 + \mathbb{Z}\vec{b}_1 + \dots + \mathbb{Z}\vec{b}_n$        $\vec{t} \in \mathbb{R}^{n+1}$

We recognize a **closest vector problem!**

## Approximately solving CVP

The *closest vector problem (CVP)* is the following:

given a basis  $(\vec{b}_1, \dots, \vec{b}_n)$  of a lattice  $\mathcal{B}$  and a vector  $\vec{t}$ , let  $\mu = \text{dist}(\vec{t}, \mathcal{B})$ . The problem is to find a vector  $\vec{v} \in \mathcal{B}$  such that  $\|\vec{v} - \vec{t}\| = \mu$ .

The associated  $\alpha$ -approximation problem consists in finding a vector  $\vec{v}$  in the lattice such that  $\|\vec{v} - \vec{t}\| \leq \alpha\mu$ .

Babai's nearest plane algorithm uses an LLL-reduced basis of the lattice to solve the  $2^{n/2}$ -approximation of CVP:

```

Input:
  an LLL-reduced basis (b1, ..., bn)
  a vector t
Output:
  a 2^(n/2)-approximation of CVP

(c1, ..., cn) = GramSchmidt(b1, ..., bn);
v = t;
for (j=n; j > 0; j--) {
  v = v - [(v, cj)/(cj, cj)] * bj;
}
return t-v;
    
```

In this pseudo-code,  $[x]$  denotes the nearest integer to  $x$ ;  $(v, w)$  is the dot product of  $v$  and  $w$ .

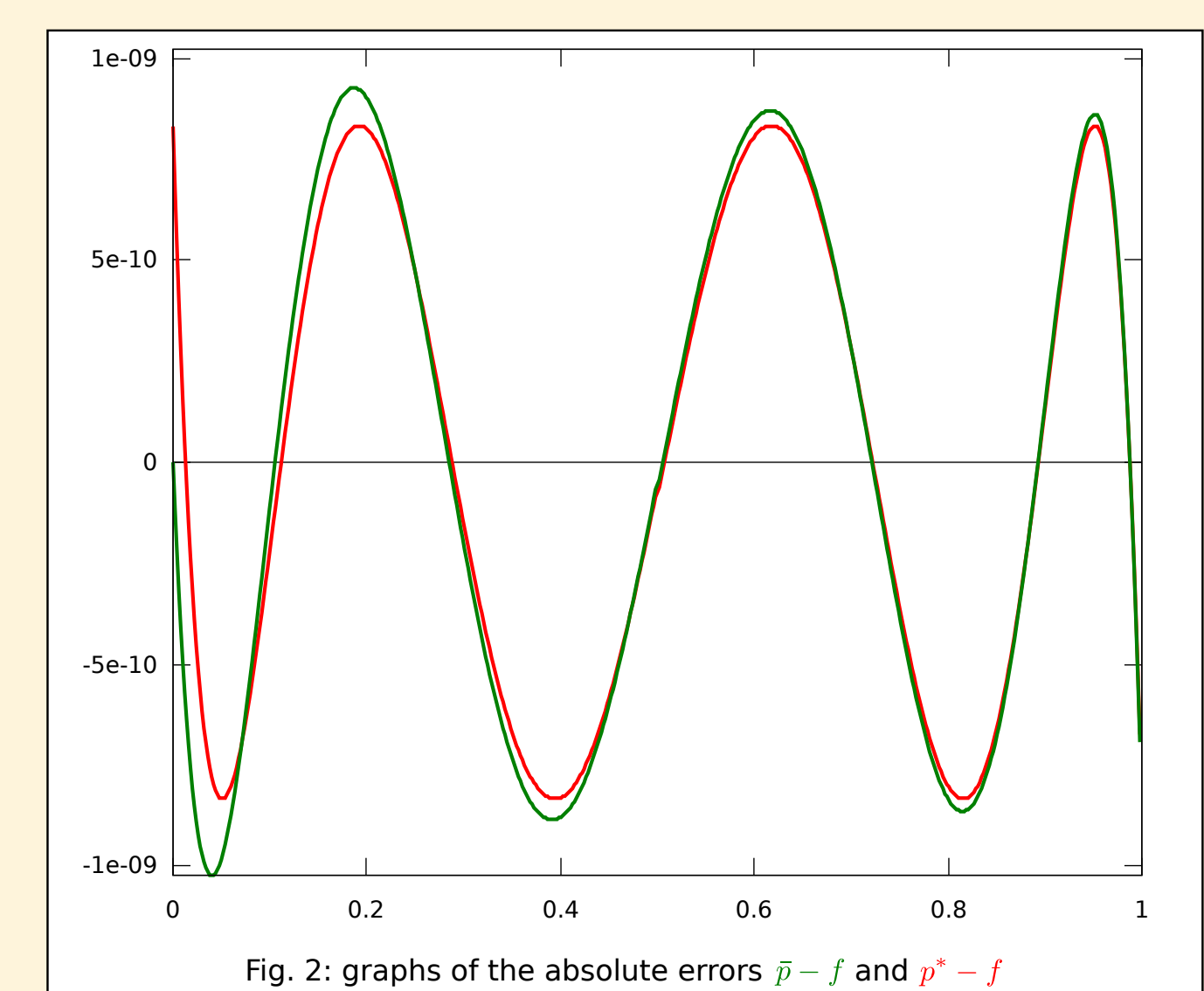
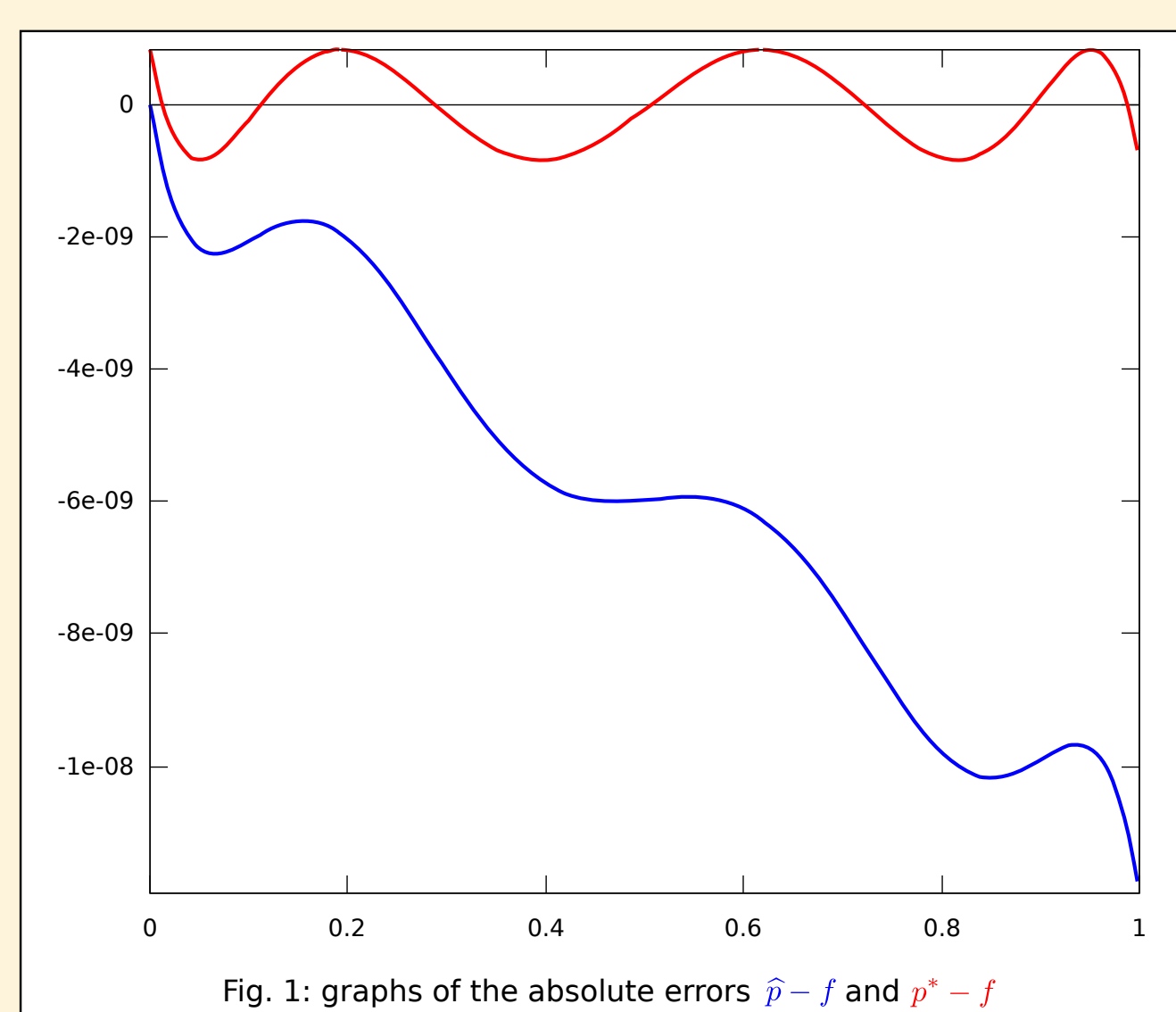
## Example

In this example, we try to approximate  $f : \log_2(1 + 2^{-x})$  on  $[0, 1]$  by a polynomial  $p$  of degree 6 with single precision coefficients (i.e. 24 bits of mantissa).

$p^*$  denotes the best polynomial with real coefficients;

$\hat{p}$  denotes the polynomial obtained by rounding each coefficient of  $p^*$  to the nearest floating-point number;

$\bar{p}$  denotes the polynomial that we obtain with our method.



As can be seen,  $\bar{p}$  is far closer to  $f$  than  $\hat{p}$ . Especially,  $\bar{p}$  is almost as accurate as the best polynomial  $p^*$ .

For the  $n+1 = 7$  points required by the method, we chose those where the red curve crosses the  $x$ -axis. A general theorem due to Chebyshev ensures that there are always at least  $n+1$  points where the best polynomial with real coefficients  $p^*$  equals the function, i.e. where  $p^* - f$  vanishes.