

A correctly rounded binary64 cube root

Robin Leroy (eggrobin)

REMOVE BEFORE FLIGHT 2021-04-36

This document describes the computations in `numerics/cbrt.cpp`.

On some abridged root-finding methods

We recall two families of root-finding methods from the late 17th century.

In [Fan91a], Thomas Fantet de Lagny first presents the iterations

$$a \mapsto \frac{1}{2}a + \sqrt{\frac{1}{4}a^2 + \frac{b}{3a}}, \quad (1)$$

hereafter the (*quadratic*) *irrational method*, and

$$a \mapsto a + \frac{ab}{3a^3 + b}, \quad (2)$$

the *rational method*, for the computation of the cube root $\sqrt[3]{a^3 + b}$, mentioning the existence of similar methods for arbitrarily higher powers. In [Fan91b] the above methods are again given, with an outline of the general method for higher powers, and a mention of their applicability to finding roots of polynomials other than $z^p - r$.

That general method is given in detail in [Fan92, p. 19]. Modernizing the notation, the general rule is as follows for finding a root of the monic polynomial of degree $p \geq 2$

$$f(z) := z^p + c_1 z^{p-1} + \dots + c_{p-1} z + c_p =: z^p - R(z)$$

with an initial approximation a .

Separate the binomial expansion of $(x + \frac{1}{2}a)^p$ into alternating sums of degree p and $p - 1$ in z ,

$$S_1 := \sum_{\substack{k=0 \\ 2 \nmid k}}^p \binom{p}{k} x^{p-k} \left(\frac{1}{2}a\right)^k \text{ and } S_2 := \sum_{\substack{k=0 \\ 2 \mid k}}^p \binom{p}{k} x^{p-k} \left(\frac{1}{2}a\right)^k,$$

and consider the following polynomials, of degree p and $p - 1$ in x for almost all a :

$$E_p := S_1 - \frac{1}{2}R\left(x + \frac{1}{2}a\right) \text{ and } E_{p-1} := S_2 - \frac{1}{2}R\left(x + \frac{1}{2}a\right).$$

Let E_{n-1} be the remainder of the polynomial division¹ of E_{n+1} by E_n ; its degree is $n - 1$ for almost all a . The iteration is $a \mapsto x + \frac{1}{2}a$, where x is a root of E_2 in the quadratic irrational method, and the root of E_1 in the rational method. Its order is p .

¹While the rest of the method is a straightforward translation, this step bears some explanations; its description in [Fan92] is

De ces deux égalitez, ou prifes féparément, ou comparées enfemble felon la methode des problèmes plus que déterminez tirez en une valeur d' x rationelle, ou fimplement d'un degré commode.

It is assumed that the reader is familiar with this “comparison according to the method of more-than-determined problems”. While the application of the root-finding method is described in painstaking detail in [Fan33], which outlines the treatment of overdetermined problems, it is perhaps this remark from [Fan97, p. 494] which lays it out most clearly:

Il n'y a rien de nouveau à remarquer fur les Problemes plus que déterminez du quatrième degré. La Regle générale est d'égalier tout à zero, & de divifer la plus haute équation par la moins élevée, ou l'également élevée l'une par l'autre, continuellement jufques à ce que l'on trouve le refte ou le divifeur le plus fimple.

Names and multiplicity of the irrational methods

Lagny does not require that the polynomial division be carried out all the way to E_2 , merely until one gets *une valeur d' x [...] d'un degré commode*, by which he likely means one that is constructible. When f is a cubic, he uses the term *formule irrationnelle* for $x + \frac{1}{2}a$ where x is a root of E_2 , but when it comes to computing the fifth root, the same term is used to refer to the case where x is a root of E_4 . In order to avoid confusion, we use the term *quadratic irrational method* when x is a root of E_2 in the general case, and we call the irrational formula from [Fan92, p. 43]² for $\sqrt[5]{a^5 + b}$

$$a \mapsto \frac{1}{2}a + \sqrt{\sqrt{\frac{1}{4}a^4 + \frac{b}{5a}} - \frac{1}{4}a^2}$$

Lagny's *quartic irrational method* for the fifth root, whereas the quadratic irrational method for the same fifth root would be³

$$a \mapsto \frac{a(7b - \sqrt{100a^{10} + 100a^5b - 7b^2})}{4b - 10a^5}.$$

When $p = 3$, there is however no ambiguity; assuming one does not wish to reduce the computation of the root of a cubic to that of another cubic, the only irrational method is the quadratic one.

Names of the rational method

Modern calculus allows us to give a more straightforward expression for the rational method than was available to Lagny; the proof of the following equivalence will be given at the end of this section.

Proposition. *The iteration of Lagny's rational method for a monic polynomial f of degree p is*

$$a \mapsto a + (p-1) \frac{(1/f)^{(p-2)}(a)}{(1/f)^{(p-1)}(a)}. \quad (3) \quad \square$$

The iteration (3) is a special case of the *Algorithmen* (A_ω^λ) defined by Schröder for an arbitrary polynomial f in [Sch70, pp. 349 sq.], equation (69); specifically, it is (A_{p-1}^0). As seen in the proof of the proposition, it is also a special case of Householder's equation (14) from [Hou70, p. 169], which generalizes it by substituting f/g for f , and letting f be an arbitrary analytic function. The case $g \equiv 1$ is mentioned in theorem 4.4.2, and that expression is given explicitly in [SG01].

For $p = 2$ and f an arbitrary polynomial, (3) is Newton's method, presented by Wallis in [Wal85, p. 338].

For $p = 3$ and f an arbitrary polynomial, it is Halley's rational method, given in [Hal94, pp. 142–143] in an effort to generalize⁴ Lagny's (2). It is usually simply known

²The formula has a misprint in [Fan92, p. 43], $-\frac{1}{2}a^2$ instead of $-\frac{1}{4}a^2$ under the radical. Halley remarks on it and gives the corrected formula in [Hal94, pp. 137, 140]. The misprint remains forty years later in [Fan33, p. 440 misnumbered 340].

³Both are of order 5, but the reader who wishes to compute a fifth root should note that leading term of the error of the quartic method is $\frac{2}{7}$ of that of the quadratic.

⁴Lagny's method is general, in that an iteration is given for any polynomial, albeit one whose order changes with the degree. However, while he refers to its results—and even corrects a misprint therein—, Halley did not have access to a copy of [Fan92].

Has Regulas, cum nondum librum videram, ab amico communicatas habui

and it appears that said friend communicated only the formulæ for the cube and fifth root, as opposed to the general method and its proof, as Halley writes

[...] *D. de Lagney* [...] qui cum totus fere sit in eliciendis Potestatum purarum radicibus, præfertim Cubicâ, pauca tantum eaque perplexa nec fatis demonstrata de affectarum radicum extractione fubjungit.

or, about the quartic irrational method for the fifth root, whereon Lagny does not elaborate as it is a direct application of the general method,

Author autem nullibi inveniendi methodum ejufve demonstrationem concedit, etiamfi maxime defiderari videatur [...].

Being unaware of this generality, Halley sets out to generalize (1) and (2) to arbitrary polynomials, and does so by keeping the order constant.

as Halley’s method, as the irrational method—which likewise generalizes Lagny’s irrational method for $p = 3$ while retaining constant order as the degree changes—has comparatively fallen into obscurity; see [ST95].

Considering, as remarked by [Sch70, p. 334], that a method can often be generalized from arbitrary polynomials or rational functions to arbitrary analytic functions, we call the iteration (3)

- Newton’s method when $p = 2$, for arbitrary f ;
- Lagny’s rational method when $p > 2$ and f is a polynomial of degree p ;
- Halley’s rational method when $p = 3$ and f is not a polynomial of degree 3;
- the Lagny–Schröder rational method of order p otherwise.

We do not simply call this last case “Schröder’s method”, as it is only a special case of the methods defined in [Sch70], so that the expression would be ambiguous.

Note that we avoid the name “Householder’s method” which appears in [SG01] and ulterior works (notably *MathWorld* and *Wikipedia*, both citing [SG01]), as it is variably used to refer to either (3) or to a method from a different family, namely φ_{p+1} from [Hou70, p. 168], equation (7), taking $\gamma_{p+1} \equiv 0$ in the resulting iteration; φ_3 is⁵ the iteration given in section 3.o.3 of [SG01]. As mentioned by Householder, both of those were described by Schröder a century prior anyway: Householder’s (7) is Schröder’s (18) from [Sch70, p. 327].

Bibliographic note

Our foray into the history of these methods was prompted by finding the “historical background” section of [ST95] while looking for a reference for Halley’s method: it is mentioned therein that this method, as applied to the cube root, is due to Lagny.

Searching for Lagny’s work led us to the historical note [Can61], wherein a note by the editors Terquem and Gerono reads

Naturellement, en mathématiques, séjour des propositions irréfragables, identiques en toute langue, en tout pays, ces rencontres ne peuvent manquer d’être assez fréquentes; nulle part les plagiats *effectifs* sont si rares, et les plagiats *apparents* si communs que dans la science exacte par excellence; mais les signaler est un devoir, un service rendu à l’histoire scientifique.

The editors then quote a letter by Prouhet, wherein he gives a reference to [Fan92].

Lagny’s work proved far more extensive than we expected: besides the above root finding methods for arbitrary polynomials, it contains an error analysis, and even a discussion of the principles of performance analysis based on a decomposition into elementary operations on—and writing of—decimal digits: a 17th century MIX.

Observing that the higher-order examples correspond to the well-known higher order method attributed to Householder in [SG01], we looked for its properties in [Hou70] so as to prove that observation, and found that Householder attributes them to Schröder. As mentioned in the translator’s note by Stewart in [SS93],

A. S. Householder used to claim you could evaluate a paper on root finding by looking for a citation of Schröder’s paper. If it was missing, the author had probably rediscovered something already known to Schröder.

It is quite possible that the irrational methods could be expressed using Schröder’s methods in one way or another, but while the result would probably be more convenient, it is unlikely to be something well-known, as irrational methods are far less popular nowadays—unjustifiedly so, as we shall see.

Prouhet’s letter in [Can61] ends with these words:

Tout cela est fort abrégé; mais qui nous délivrera des méthodes abrégées, qui n’en finissent pas?

⁵We are grateful to Peter Barfuss for this observation.

Proof of the proposition

We now prove the above proposition, which, substituting the definition of Lagny's rational method, is that

$$x + \frac{1}{2}a = a + (p-1) \frac{(1/f)^{(p-2)}(a)}{(1/f)^{(p-1)}(a)} =: \psi(a)$$

if x is the root of E_1 .

Proof. Let $E_p = d_0x^p + \dots + d_p$, $E_{p-1} = e_0x^{p-1} + \dots + e_{p-1}$. As shown in [Hou70, pp. 52–54], the polynomial remainders E_k are given up to a constant factor by [Hou70, p. 19] equation (23), i.e., for some α ,

$$\frac{E_k}{\alpha_k} = \det \begin{pmatrix} (E_p)_{p-1-k} \\ (E_{p-1})_{p-k} \end{pmatrix},$$

where the expression on the right-hand side is the *bigradient* defined in [Hou68] (3.2) or [Hou70, p. 19] (20),

$$\frac{E_k}{\alpha_k} = \det \begin{pmatrix} d_0 & d_1 & d_2 & \dots & d_{2(p-k)-3} & x^{p-k-2}E_p \\ & d_0 & d_1 & \dots & d_{2(p-k)-4} & x^{p-k-3}E_p \\ & & \ddots & & \vdots & \\ \mathbf{0} & & & d_0 & d_1 & \dots & d_{p-k-1} & x^0E_p \\ & & & e_0 & \dots & e_{p-k-2} & x^0E_{p-1} \\ & & & \vdots & & \vdots & \\ & & e_0 & \dots & e_{2(p-k)-4} & x^{p-k-2}E_p \\ & e_0 & e_1 & \dots & e_{2(p-k)-4} & x^{p-k-2}E_p \\ e_0 & e_1 & e_2 & \dots & e_{2(p-k)-3} & x^{p-k-1}E_p \end{pmatrix} =: \det \mathbf{E}_k,$$

where $d_n := 0$ for $n > p$, and $e_n := 0$ for $n > p-1$.

In particular, for $k = 1$,

$$\mathbf{E}_1 = \begin{pmatrix} d_0 & d_1 & d_2 & \dots & d_{p-3} & d_{p-2} & d_{p-1} & d_p & \mathbf{0} & x^{p-3}E_p \\ & d_0 & d_1 & \dots & d_{p-4} & d_{p-3} & d_{p-2} & d_{p-1} & d_p & x^{p-4}E_p \\ & & \ddots & & & & & & \vdots & \\ \mathbf{0} & & & d_0 & d_1 & d_2 & \dots & d_{p-2} & x^0E_p \\ & & & e_0 & e_1 & \dots & e_{p-3} & x^0E_{p-1} \\ & & & \vdots & & & \vdots & \\ & e_0 & \dots & e_{p-5} & e_{p-4} & e_{p-3} & e_{p-2} & e_{p-1} & x^{p-5}E_{p-1} \\ & e_0 & e_1 & \dots & e_{p-4} & e_{p-3} & e_{p-2} & e_{p-1} & x^{p-4}E_{p-1} \\ e_0 & e_1 & e_2 & \dots & e_{p-3} & e_{p-2} & e_{p-1} & \mathbf{0} & x^{p-3}E_{p-1} \end{pmatrix}.$$

Observe that, since the value of x used in the rational method is the root of E_1 , for that value of x , $\det \mathbf{E}_1 = 0$, i.e., \mathbf{E}_1 is singular.

Lemma. *The matrix \mathbf{E}_1 is singular if and only if $\mathbf{C}(x + \frac{1}{2}a)$ is singular, where*

$$\mathbf{C}(\psi) := \begin{pmatrix} \psi - a & c_0 & & \mathbf{0} \\ -1 & c_1 & c_0 & \\ 0 & c_2 & \ddots & \\ \vdots & \vdots & \ddots & c_0 \\ 0 & c_{p-1} & \dots & c_2 & c_1 \end{pmatrix}$$

and $c_0 := 1$.

Proof (Left as an exercise to the reviewer). Observe that:

- since $\det \mathbf{E}_1$ is a polynomial of degree 1, all terms divisible by x^2 must cancel out in the Laplace expansion of that determinant on the last column, the determinant is equal to

$$\pm \delta_1 x E_p \mp \delta_2 E_p \pm \delta_3 E_{p-1} \mp \delta_4 x E_{p-1},$$

where the δ_i are determinants of real matrices;

TODO(egg): Prove the lemma.

- by the same reasoning, only the linear and constant terms of the polynomials remain in the above expression, which simplifies to

$$\pm \delta_1 d_p x \mp \delta_2 (d_{p-1} x + d_p) \pm \delta_3 (e_{p-2} x + e_{p-1}) \mp \delta_4 e_{p-1} x;$$

- $E_p - E_{p-1} = S_1 - S_2$, which, by Lagny's *theoreme fondamentale* [Fan92, p. 17], is $\left(x - \frac{1}{2}a\right)^p$.

The proof is a calculation. \square

The proposition follows from the lemma and theorem 4.4.2 from [Hou70, p. 169]: $\psi(a)$ is Householder's (14) with $g \equiv 1$; for that value of g , theorem 4.4.2 states that (14) is the solution of (12) from the same page, which is $\det \mathcal{C}(\psi(a)) = 0$. By the lemma, for the value of x in Lagny's rational method, $x + \frac{1}{2}a$ solves that equation. \square

A faithfully rounded cube root

We now turn to the computation in `numerics/cbrt.cpp`.

Overview

Our general approach to computing a faithfully rounded cube root of $y > 0$ is the one described in [KB01]:

1. integer arithmetic is used to get an initial quick approximation q of $\sqrt[3]{y}$;
2. a root finding method is used to improve that to an approximation ξ with a third of the precision;
3. ξ is rounded to a third of the precision, resulting in the rounded approximation x whose cube x^3 can be computed exactly;
4. a single high order iteration of a root finding method is used to get the faithfully rounded result r_0 .

Notation

We define the fractional part as $\text{frac } a := a - \lfloor a \rfloor \in [0, 1[$, regardless of the sign of a .

The floating-point format used throughout is `binary64`; the quantities $p \in \mathbb{N}$ (precision in bits) and $\text{bias} \in \mathbb{N}$ are defined as in IEEE 754-2008, $p = 53$ and $\text{bias} = 1023$. Some of the individual methods discussed may be of general use; we thus give all constants to thirty-six decimal digits⁶, which should suffice for `decimal128`, `binary128`, and all smaller formats.

We use capital Latin letters for fixed-point numbers involved in the computation, and $A > 0$ for the normal floating-point number $a > 0$ reinterpreted as a binary fixed-point⁷ number with $p - 1$ bits after the binary point,

$$\begin{aligned} A &:= \text{bias} + \lfloor \log_2 a \rfloor + \text{frac}(2^{-\lfloor \log_2 a \rfloor} a) \\ &= \text{bias} + \lfloor \log_2 a \rfloor + 2^{-\lfloor \log_2 a \rfloor} a - 1, \end{aligned}$$

and *vice versa*,

$$a := 2^{\lfloor A \rfloor - \text{bias}} (1 + \text{frac } A).$$

⁶The constants used in the methods are rounded to the nearest decimal digit. A superscript sign after a final 5 serves as the sticky bit: the unrounded quantity is in excess of the rounded one if the sign is +, and in default if it is -. The error bounds (and their logarithms) are rounded towards positive infinity.

⁷The implementation uses integers (obtained by multiplying the fixed-point numbers by 2^{p-1}). For consistency with [KB01] we work with fixed-point numbers here. Since we do not multiply fixed point numbers together, the expressions are unchanged.

This corresponds to [KB01]’s $B + K + F$.

For both fixed- and floating-point numbers, given $\alpha \in \mathbb{R}$, we write $\llbracket \alpha \rrbracket$ for the nearest representable number (rounding ties to even). For fixed-point numbers, we write $\llbracket \alpha \rrbracket_0$ for directed rounding towards 0 to the fixed-point precision (as in division implemented with integer division). We write the unit roundoff $u := 2^{-p}$, and, after [Higo2, p. 63], $\gamma_n := \frac{nu}{1-nu}$. We discuss other rounding modes in appendix D.

To quote [Tre97], “If rounding errors vanished, 95% of numerical analysis would remain”. While we keep track of rounding errors throughout, they are of very little importance until the last step; when it is convenient to solely study the truncation error, we work with ideal quantities affected with a prime, which correspond to their primeless counterparts by removal of all intervening roundings.

The input y and all intervening floating-point numbers are taken to be normal; the rescaling performed to avoid overflows also avoids subnormals. We work only with correctly rounded addition, subtraction, multiplication, division, and square root; FMA is treated separately in appendix A.

1 Quick approximation

The quick approximation q is computed using fixed-point arithmetic as

$$Q := C + \left\llbracket \frac{Y}{3} \right\rrbracket_0,$$

where the fixed-point constant C is defined as⁸

$$C := \left\llbracket \frac{2 \text{bias} - \Gamma}{3} \right\rrbracket$$

for some $\Gamma \in \mathbb{R}$.

Let $\varepsilon_q := \frac{q}{\sqrt[3]{y}} - 1$, so that $\sqrt[3]{y}(1 + \varepsilon_q) = q$; the relative error of q as an approximation of $\sqrt[3]{y}$ is $|\varepsilon_q|$. Considering Y , Q , q , and ε_q as functions of y , we have

$$\begin{aligned} Y(8y) &= Y(y) + 3, \\ Q(8y) &= Q(y) + 1, \\ q(8y) &= 2q(y), \\ \varepsilon_q(8y) &= \varepsilon_q(y), \end{aligned}$$

so that the properties of ε_q need only be studied on some interval of the form $[\eta, 8\eta[$.

Pick $\eta := 2^{\lfloor \Gamma \rfloor}$, and $y \in [\eta, 8\eta[= [2^{\lfloor \Gamma \rfloor}, 2^{\lfloor \Gamma \rfloor + 3}[$, so that $\log_2 y \in [\lfloor \Gamma \rfloor, \lfloor \Gamma \rfloor + 3[$. Let $k := \lfloor \log_2 y \rfloor - \lfloor \Gamma \rfloor$; note that $k \in \{0, 1, 2\}$. Let $f := \text{frac}(2^{-\lfloor \log_2 y \rfloor} y) \in [0, 1[$. Up to at most 3 half-units in the last place from rounding (2 from the directed rounding of the division by three and 1 from the definition of C), we have

$$\begin{aligned} Q &\approx Q' := \text{bias} + \frac{\lfloor \log_2 y \rfloor}{3} + \frac{\text{frac}(2^{-\lfloor \log_2 y \rfloor} y) - \Gamma}{3}, \\ &= \text{bias} + \frac{\lfloor \Gamma \rfloor + k}{3} + \frac{f - \Gamma}{3}, \\ &= \text{bias} + \frac{k + f - \text{frac } \Gamma}{3}. \end{aligned}$$

Since $k \in [0, 2]$, the numerator $k + f - \text{frac } \Gamma$ lies in $] -1, 3[$. Further, it is negative only if $k = 0$, so that

$$\begin{aligned} \lfloor Q' \rfloor &= \begin{cases} \text{bias} - 1 & \text{if } k = 0 \text{ and } \text{frac } \Gamma > \text{frac}(2^{-\lfloor \Gamma \rfloor} y), \\ \text{bias} & \text{otherwise,} \end{cases} \text{ and} \\ \text{frac } Q' &= \begin{cases} 1 + \frac{f - \text{frac } \Gamma}{3} & \text{if } k = 0 \text{ and } \text{frac } \Gamma > f, \\ \frac{k + f - \text{frac } \Gamma}{3} & \text{otherwise.} \end{cases} \end{aligned}$$

⁸Note that there is a typo in the corresponding expression $C := (B - 0.1009678)/3$ in [KB01]; a factor of 2 is missing on the bias term.

Accordingly, for the quick approximation q , we have, again up to at most 3 half-units in the last place,

$$q \approx q' = \begin{cases} 1 + \frac{f - \text{frac } \Gamma}{3} & \text{if } k = 0 \text{ and } \text{frac } \Gamma > f, \\ 1 + \frac{k + f - \text{frac } \Gamma}{3} & \text{otherwise,} \end{cases}$$

With $\sqrt[3]{y} = 2^{\frac{[\Gamma] + k}{3}} \sqrt[3]{1 + f}$, we can define

$$\varepsilon'_q := \frac{q'}{\sqrt[3]{y}} - 1,$$

which we can express piecewise as a function of f and k . This gives us a bound on the relative error,

$$|\varepsilon|_q \leq |\varepsilon'_q|(1 + 3u).$$

The values $\Gamma = 0.1009678$ and $\varepsilon_q < 3.2\%$ from [KBo1] may be recovered by choosing Γ minimizing the maximum of $|\varepsilon'_q|$ over $y \in [\eta, 8\eta]$, or equivalently.

$$\Gamma_{\text{Kahan}} := \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_{y \in [\eta, 8\eta]} |\varepsilon'_q| = \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_{(f, k)} |\varepsilon'_q|$$

where the maximum is taken over $(f, k) \in [0, \text{frac } \Gamma[\times \{0\} \cup [0, 1[\times \{1, 2\}$,

$$= \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_{(f, k) \in \mathcal{E} \cup \mathcal{L}} |\varepsilon'_q|,$$

where $\mathcal{E} := \{(\text{frac } \Gamma, 0)\} \cup \{(0, k) \mid k \in \{0, 1, 2\}\}$ is the set of the endpoints of the intervals whereon q' is piecewise affine, and $\mathcal{L} := \left\{ \left(\frac{k - \text{frac } \Gamma}{2}, k \right) \mid k \in \{1, 2\} \right\}$ are the local extrema. We get more precisely⁹

$$\Gamma_{\text{Kahan}} \approx 0.10096\,78121\,55802\,88786\,36993\,42643\,55358\,1$$

with $\max_y |\varepsilon'_q| \approx 3.155\%$, yielding the constant

$$C_{\text{Kahan}} = {}_{16}2A9F\,7625\,3119\,D328 \cdot 2^{-52}$$

for IEEE 754-2008 binary64. However, as we will see in the next section, this value does not optimize the final error.

2 Getting to a third of the precision

We now consider multiple methods for the refinement of q to ξ . The rounding error in this step being both negligible and tedious to bound, its analysis is relegated to appendix B. Here we will study only the truncation error, and thus work only with the primed quantities.

Lagny's rational method

One way to compute ξ' is Lagny's rational method,

$$\xi' = q' + \frac{q'(y - q'^3)}{2q'^3 + y},$$

with the error

$$\varepsilon'_\xi := \frac{\xi'}{\sqrt[3]{y}} - 1.$$

With $q' = \sqrt[3]{y}(1 + \varepsilon'_q)$, we can express ε'_ξ using the transformation of the relative error error by one step of Lagny's rational method on the cube root,

$$\varepsilon'_\xi = \frac{2\varepsilon_q'^3 + \varepsilon_q'^4}{3 + 6\varepsilon_q' + 6\varepsilon_q'^2 + 2\varepsilon_q'^3} = \frac{2}{3}\varepsilon_q'^3 + \mathcal{O}(\varepsilon_q'^4).$$

⁹This value may be computed formally, but the expression is unwieldy.

If q' is computed using $\Gamma = \Gamma_{\text{Kahan}}$, we get $\max_y |\varepsilon'_\xi| \approx 21.96 \cdot 10^{-6}$, $\log_2 \max_y |\varepsilon'_\xi| \approx -15.47$. However, Γ_{Kahan} , which minimizes $\max_y |\varepsilon_q|$, does not minimize $\max_y |\varepsilon'_\xi|$. This is because while ε'_ξ is monotonic as a function of ε'_q , it is not odd: positive errors are reduced more than negative errors are, so that the minimum is attained for a different value of Γ . Specifically, we have

$$\begin{aligned} \Gamma_{L^1} &:= \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_y |\varepsilon'_\xi| \\ &\approx 0.09918\,74615\,29855\,99525\,66149\,20761\,31234\,35^- \end{aligned}$$

with $\max_y |\varepsilon'_q| \approx 3.2025\%$, but $\max_y |\varepsilon'_\xi| \approx 20.86 \cdot 10^{-6}$, $\log_2 \max_y |\varepsilon'_\xi| \approx -15.54$. The corresponding fixed-point constant is

$$C_{L^1} := {}_{16}2A9F\,7893\,782D\,A1CE \cdot 2^{-52}$$

for binary64.

While it is close to the 16 bits to which we will round in the next step, this error is still larger, and in any case is not comparatively negligible. As a result, it significantly contributes to the final error above $1u$, *i.e.*, to misrounding. Lagny's lesser-known irrational method provides us with a way to improve it.

Lagny's irrational method

As written in (1), Lagny's irrational method

$$\xi' = \frac{1}{2}q' + \sqrt{\frac{1}{4}q'^2 + \frac{y - q'^3}{3q'}}$$

seems prohibitively computationally expensive in comparison to the rational one: it adds a square root on the critical path, dependent on the result of a division. However, rewriting it as

$$\xi' = \frac{1}{2}q' + \frac{1/\sqrt{12}}{q'} \sqrt{4yq' - q'^4}, \quad (4)$$

one can evaluate it with similar¹⁰ performance to the rational method.

Its error is

$$\varepsilon'_\xi = \frac{-\varepsilon_q'^3}{3\left(\frac{1}{2} + \sqrt{\frac{1}{2} - 2\varepsilon_q'^2 - \frac{4}{3}\varepsilon_q'^3 - \frac{1}{3}\varepsilon_q'^4 - \varepsilon_q'^2}\right)} = -\frac{1}{3}\varepsilon_q'^3 + \mathcal{O}(\varepsilon_q'^4),$$

whose leading term is half that of the rational method; indeed we find that with $\Gamma = \Gamma_{\text{Kahan}}$, we have $\max_y |\varepsilon'_\xi| \approx 10.48 \cdot 10^{-6}$, $\log_2 \max_y |\varepsilon'_\xi| \approx -16.54$, gaining one bit with respect to the rational method. Here $\Gamma = \Gamma_{\text{Kahan}}$ is very close to optimal; with the optimal value

$$\Gamma_{L^2} \approx 0.10096\,82076\,65096\,37285\,40885\,52460\,33434\,6,$$

the error remains the same within the precision to which we have given it. However, we have other ways of improving the error at no cost to performance.

Canon optimization of Lagny's irrational method

The idea for this optimization comes from [Can18a], reproduced here with the author's permission:

A trick I've used for years and should write up: you can apply optimization to the iteration, not just the starting guess: $x' = xp(x)$, select $p(x)$ to be minimax error on bounded initial error in x . This yields a nice family of tunable approximations.

¹⁰We compare specific evaluations strategies for various architectures in appendix C; on some of those, other rewritings are superior to this one.

Everyone else seems to worry about starting estimate, but use standard iterations, which is appropriate for arbitrary precision, but silly with a fixed precision target.

Note that as p gets to be high-order, it converges quickly to the Taylor series for the correction, but there's a nice space with cheap initial approximations and order 2–5 or so, because we can evaluate these polynomials with lower latency [than] serially-dependent iterations.

Canon later elaborated on this in [Can18b]:

Quick version: we want to compute $1/\sqrt{y}$, we have an approximation x_0 , we want to improve it to $x_1 = x_0 p(x_0, y)$. For efficiency, we want p to be a polynomial correction.

handwavy motivation for brevity make p a polynomial in $x_0 x_0 y$, which is approximately 1.

Specifically, if x_0 has relative error e , $x_0 x_0 y$ is bounded by something like $1 \pm 2e$. So, we want to find p that minimizes $|x/x_0 - p(x_0 x_0 y)|$ on $[1 - 2e, 1 + 2e]$. NR¹¹ uses the $p = 1$ st order Taylor. We know that we can do better via usual approximation theory techniques.

We can also use higher-order approximations to hit any specific accuracy target in a single step. This isn't always better than iterating, but sometimes it is.

We do not use a polynomial—nor even a rational function—, nor do we express our refinement as a function of a quantity bounded by the error. However, we take advantage of Canon's key idea of “apply[ing] optimization to the iteration, not just the starting guess”; the latter is what we have so far done with Γ .

The constants $\frac{1}{2}$, $\frac{1}{4}$, and 3 in Lagny's irrational method may be modified with no effect on performance; altering the first two of these introduces rounding errors, but these need not concern us here. We thus write

$$\xi' = \kappa q' + \sqrt{\lambda q'^2 + \frac{y - q'^3}{\mu q'}}$$

and choose Γ , κ , λ , and μ minimizing relative error in the Чебышёв norm,

$$(\Gamma_{L^2C}, \kappa_{L^2C}, \lambda_{L^2C}, \mu_{L^2C}) := \operatorname{argmin}_{\Gamma, \kappa, \lambda, \mu} \max_y |\varepsilon'_\xi|.$$

Unfortunately, computing $\max_y |\varepsilon'_\xi|$ is not as easy as for the standard methods; the introduction of κ , λ , and μ breaks the monotonicity of $\varepsilon'_\xi(\varepsilon'_q)$, so that the local extrema of ε'_ξ are not found in the same place as those of ε'_q . Formally looking for zeros of the derivative of ε'_ξ with respect to f is impractical. Instead we find the local maxima by numerical maximization on the four pieces whereon q' is a smooth function of f .

That maximum can be minimized by a straightforward hill-climbing¹² starting from $\Gamma = \frac{1}{10}$, $\kappa = \frac{1}{2}$, $\lambda = \frac{1}{4}$, and $\mu = 3$. We obtain the values

$$\begin{aligned} \Gamma_{L^2C} &\approx 0.10007\,61614\,69941\,46538\,73178\,74111\,71965\,6, \\ \kappa_{L^2C} &\approx 0.49999\,99381\,08574\,04775\,14291\,72928\,30652\,9, \\ \lambda_{L^2C} &\approx 0.25000\,00000\,00145\,58487\,81104\,01052\,77249\,3, \\ \mu_{L^2C} &\approx 3.00074\,62871\,20756\,72280\,51404\,24030\,90920, \end{aligned}$$

for an error of $\max_y |\varepsilon'_\xi| \approx 2.6157 \cdot 10^{-6}$, $\log_2 \max_y |\varepsilon'_\xi| \approx -18.54$: this optimization gains two bits. The resulting ε'_ξ is remarkably equioscillating, as can be seen in figure ??.

¹¹Newton–Raphson, i.e., Newton's method. REMOVE BEFORE FLIGHT: Cite Raphson, cite Lagrange's remarks

¹²It is plausible that some variation of Pemež's algorithm could be used here, much like it can be adapted to rational functions; since the hill-climbing converged satisfactorily, and did so much faster than we were writing this document, we have not investigated this.

With the rewriting (4), the constants $1/\sqrt{12}$ and 4 should be replaced by

$$\sqrt{\frac{1 - \lambda_{L^2C}\mu_{L^2C}}{\mu_{L^2C}}} \approx 1/\sqrt{12.01194\,95117\,19793\,69720\,48452\,37177\,0703}$$

$$\approx 0.28853\,15115\,62316\,71905\,38451\,44194\,38406\,3$$

and

$$\frac{1}{1 - \lambda_{L^2C}\mu_{L^2C}} \approx 4.00298\,73779\,31697\,18250\,67433\,26901\,80421$$

respectively.

Note that a similar optimization could be applied to the rational method; however, it would not unconditionally be free: changing the 2 in the denominator turns an addition into a multiplication, and inserting additional constants adds more operations. Whether this hinders performance depends on the architecture. In any case, the optimization can scarcely gain more than two bits; such an optimized rational method would still have double the error of the optimized irrational method.

3 Rounded approximation

The number x is obtained from ξ by zeroing all but the most significant $\lfloor \frac{p}{3} \rfloor$ bits of its significand. The resulting relative error $\left| \frac{x}{\xi} - 1 \right|$ is greatest when the zeroed bits are all 1 and the remaining bits (except for the leading 1) are all 0; this is the case when the significand of ξ is $1 + 2^{-\lfloor \frac{p}{3} \rfloor + 1} - 2^{1-p}$, in which case that of x is 1, so that

$$\left| \frac{x}{\xi} - 1 \right| \leq 1 - \frac{1}{1 + 2^{-\lfloor \frac{p}{3} \rfloor + 1} - 2^{1-p}} < 2^{-\lfloor \frac{p}{3} \rfloor + 1} = 2^{-16}.$$

For the error of x as an approximation of the cube root,

$$\varepsilon_x := \frac{x}{\sqrt[3]{y}} - 1,$$

we have the bound $|\varepsilon_x| < (1 + |\varepsilon_\xi|)(1 + 2^{-16}) - 1$.

4 High order iteration

We use one iteration of the Lagny–Schröder rational method of order 5:

$$\left\| x - \frac{\left\| (x^3 - y) \left\| \left\| \left\| \left\| 10x^3 \right\| + 16y \right\| x^3 \right\| + \left\| y^2 \right\| \right\| \right\| \right\|}{\left\| 3x^2 \left\| \left\| \left\| 5x^3 \right\| + \left\| 17y \right\| x^3 \right\| + \left\| 5 \left\| y^2 \right\| \right\| \right\|} \right\| \right\|$$

where $3x^2$ and x^3 are exact thanks to the trailing 0s of x , $\llbracket x^6 \rrbracket$ is correctly rounded because it is computed as the square of x^3 , and $x^3 - y$ is exact by Sterbenz's lemma.

In infinite precision, this method is of such high order that if $|\varepsilon_x| < 14.5$, which is the case even if ξ is computed by the rational method, the relative error of the result is less than 2^{-75} . We will not seek to bound the truncation error more closely, nor to tweak the constants in the method to optimize it: as we will see, it is dominated by rounding.

Thanks to the exact cube and exact difference, the rounding analysis of the correction term is straightforward. All remaining sums being of positive terms, their relative error is readily bounded by the largest of those of their terms. This leads to bounds of γ_5 on the numerator and γ_5 on the denominator, overall $\frac{1+\gamma_6}{1-\gamma_5} - 1 < \frac{\gamma_{11}}{1-\gamma_5}$ on the correction term.

However, considering that our final bound on the excess above $1u$, and thus our final misrounding estimate, is proportional to this error, a more careful analysis is warranted. Observe that $x^3 = y(1 + \varepsilon_x)^3$, so that a sum

$$\Sigma' = \alpha x^m y^{p-m} (1 + \delta_1) + \beta x^n y^{p-n} (1 + \delta_2)$$

whose terms carry the errors δ_i may be rewritten as

$$\begin{aligned}\Sigma &= \alpha x^m y^{p-m} + \beta x^n y^{p-n} + \alpha x^m y^{p-m} \delta_1 + \beta x^n y^{p-n} \delta_2 \\ &= \alpha x^m y^{p-m} + \beta x^n y^{p-n} + \alpha y^p (1 + \varepsilon_x)^{3m} \delta_1 + \beta y^p (1 + \varepsilon_x)^{3n} \delta_2 \\ &= \alpha x^m y^{p-m} + \beta x^n y^{p-n} + y^p (\alpha \delta_1 + \beta \delta_2) + ?\end{aligned}$$

– REMOVE BEFORE FLIGHT, MOVE THE 6TH ORDER STUFF TO THE APPENDIX –

We use one iteration of the Lagny–Schröder rational method of order 6:

$$r_0 := \left[x - \frac{\left[\left[x(x^3 - y) \right] \left[\left[\left[5x^3 \right] + \left[17y \right] \right] x^3 \right] + \left[5 \left[y^2 \right] \right] \right]}{\left[\left[\left[7x^3 \right] + \left[42y \right] \right] \left[x^6 \right] \right] + \left[\left[\left[30x^3 \right] + 2y \right] \left[y^2 \right] \right]} \right],$$

where x^3 is exact thanks to the trailing 0s of x , $\llbracket x^6 \rrbracket$ is correctly rounded because it is computed as the square of x^3 , and $x^3 - y$ is exact by Sterbenz’s lemma.

In infinite precision, this method is of such high order that if x has a relative error below 2^{-14} , which is the case here, the relative error of the result is less than 2^{-100} ; we will not seek to bound the truncation error more closely, nor to tweak Γ nor the method itself to optimize it: as we will see, it is dominated by rounding.

Thanks to the exact cube and exact difference, the rounding analysis is straightforward; all remaining sums, being of positive terms, have a relative error bounded by the largest of those of their terms, so that we get no more than γ_4 on the second factor of the numerator, γ_6 on the whole numerator, and similarly γ_5 on the denominator, overall $\frac{1+\gamma_7}{1-\gamma_5} - 1 < \frac{\gamma_{12}}{1-\gamma_5}$ on the correction term.

Since the relative error of x is at most $(1 + \varepsilon_x)(1 + 2^{-16}) - 1$, the rounding-free correction term is no larger than $((1 + \varepsilon_x)(1 + 2^{-16}) - 1)\sqrt[3]{y}$, so that the absolute error of the computed correction term is bounded by

$$\frac{\gamma_{12}}{1 - \gamma_5} ((1 + \varepsilon_x)(1 + 2^{-16}) - 1)\sqrt[3]{y}.$$

This gives us our bound for the relative error of the result,

$$\left| \frac{r}{\sqrt[3]{y}} - 1 \right| < \left(1 + \frac{\gamma_{12}}{1 - \gamma_5} ((1 + \varepsilon_x)(1 + 2^{-16}) - 1) \right) (1 + u) - 1.$$

For binary64, this bound is $\left| \frac{r}{\sqrt[3]{y}} - 1 \right| < 1.0004336u$.

Correct rounding

A FMA

B Rounding error analysis for the second step

$$\xi := \left[q - \frac{\left[\left(\left[q^2 \right] q \right) - y \right) q \right]}{\left[2 \left[\left[q^2 \right] q \right] + y \right]} \right].$$

Note that the subtraction in the numerator is exact by Sterbenz’s lemma [Ste74, p. 138, theorem 4.3.1].

Let $\varepsilon_\xi := \frac{\xi}{\sqrt[3]{y}} - 1$ and

It is fairly clear that ε'_ξ dominates the rounding error in the approximation $\varepsilon_\xi \approx \varepsilon'_\xi$; for the sake of completeness we quantify this.

Since we have a cancellation in the expression for ξ , a little bit of care is needed to bound those errors. As mentioned above, q approximates q' to a relative error of at most $3u < \gamma_3$. The sum in the denominator

$$d := \llbracket 2\llbracket q^2 \rrbracket q \rrbracket + y$$

has positive terms, so its relative error with respect to $d' := 2q'^3 + y$ is readily bounded,

$$\frac{d}{d'} - 1 < \gamma_{3 \cdot 3 + 3} = \gamma_{12}.$$

The cancelling difference $b := \llbracket q^2 \rrbracket q \rrbracket - y$ differs from $b' := q'^3 - y$ by at most

$$\delta := q'^3 \gamma_{3 \cdot 3 + 2} = q'^3 \gamma_{11},$$

and the numerator $\llbracket bq \rrbracket$ from $b'q'$ by at most

$$((b' + \delta)q'(1 + \gamma_3))(1 + \gamma) - b'q' = (1 + \gamma_4)\delta q' + \gamma_4 b'q'.$$

We can bound the absolute error of the correction term as

$$\begin{aligned} \left| \left\llbracket \frac{\llbracket bq \rrbracket}{d} \right\rrbracket - \frac{b'q'}{d'} \right| &< \frac{b'q' + (1 + \gamma_4)\delta q' + \gamma_4 b'q'}{d'(1 - \gamma_{12})} - \frac{b'q'}{d'} = \frac{(1 + \gamma_4)\delta q' + (\gamma_4 + \gamma_{12})b'q'}{d'(1 - \gamma_{12})} \\ &< q' \frac{(1 + \gamma_4)\delta + \gamma_{16}b'}{d'(1 - \gamma_{12})}, \end{aligned}$$

and that of ξ as

$$|\xi - \xi'| < q' \left(\gamma_3 + \frac{(1 + \gamma_4)\delta + \gamma_{16}b'}{d'(1 - \gamma_{12})} \right).$$

Substituting the truncation errors and the definition of δ , we can bound the relative error arising from rounding:

$$\left| \frac{\xi}{\xi'} - 1 \right| < \frac{1 + \varepsilon'_q}{1 + \varepsilon'_\xi} \left(\gamma_3 + \frac{(1 + \gamma_4)\gamma_{11}(1 + \varepsilon'_q)^3 + \gamma_{16}((1 + \varepsilon'_q)^3 - 1)}{((1 + \varepsilon'_q)^3 + 1) + (1 - \gamma_{12})} \right).$$

Linearizing, this bound is $(6 + \frac{2}{3} + \mathcal{O}(\varepsilon'_q + \varepsilon'_\xi))u + \mathcal{O}(u^2)$. More palatably, for either choice of γ , we have

$$\left| \frac{\xi}{\xi'} - 1 \right| < 8u$$

provided that $p \geq 12$.

C Performance of Lagny's methods for the cube root

D Other rounding modes

E Comparison with other faithful implementations

References

- [Can18a] S. Canon. “A trick I’ve used for years and should write up”. Tweets at <https://twitter.com/stephentyrone/status/1016283784067665920> *sqq.*, with a note at <https://twitter.com/stephentyrone/status/1016328842296864770>. 9th July 2018.
- [Can18b] S. Canon. “A trick I’ve used for years and should write up: Quick version”. Tweets at <https://twitter.com/stephentyrone/status/1057788315699687424> *sqq.* 1st Nov. 2018.

- [Can61] M. Cantor. “Note historique sur l’extraction abrégée de la racine carrée”. In: *Nouvelles annales de mathématiques. Journal des candidats aux écoles polytechnique et normale* 1.20 (1861). Ed. by O. Terquem and C.-C. Gerono, pp. 46–47.
eprint: http://www.numdam.org/item/?id=NAM_1861_1_20__46_1.
- [Fan33] T. Fantet de Lagny. “Analyse Générale, ou Méthodes nouvelles pour résoudre les Problèmes de tous les Genres & de tous les Degrez à l’infini”. In: *Recueil des Mémoires de l’Académie Royale des Sciences depuis 1666 jusqu’à 1699*. Ed. by C. Richer. Vol. XI. Par la compagnie des libraires, 1733.
eprint: <https://books.google.fr/books?id=KwP-kH6gm1EC>.
- [Fan91a] T. Fantet de Lagny. “Nouvelle methode de Mr. T. F. de Lagny pour l’approximation des Racines cubiques”. In: *Le Journal des sçavans* 1691.17 (14th May 1691), pp. 200–203.
eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k56538h/f202.double>.
- [Fan91b] T. Fantet de Lagny. *Méthode nouvelle, infiniment générale et infiniment abrégée, Pour l’Extraction des Racines quarrées, cubiques, &c. & pour l’Approximation des mêmes Racines à l’infini dans toutes sortes d’égalitez. Proposée à examiner aux Mathématiciens de l’Europe*. De l’Imprimerie d’Antoine Lambin, ruë S. Jacques, au Miroir, 1691.
eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k1039787>.
- [Fan92] T. Fantet de Lagny. *Methodes nouvelles et abbregees pour l’extraction et l’approximation des racines. Et pour resoudre par le cercle et la ligne droite, plusieurs problèmes solides & sursolides ; comme la duplication du cube, l’invention de deux & de quatre moyennes proportionnelles, &c. dans toute la précision possible, & d’une maniere praticable. Avec une dissertation sur les methodes d’arithmetique & d’analyse ; où l’on établit des principes generaux pour en juger*. De l’Imprimerie de Jean Cusson, ruë saint Jacques, à l’Image de saint Jean Baptiste, 1692.
eprint: <https://nubis.univ-paris1.fr/ark:/15733/3415>.
- [Fan97] T. Fantet de Lagny. *Nouveaux elemens d’arithmetique et d’algebre, ou introduction aux mathematiques*. Chez Jean Jombert, près des Augustins, à l’Image Nôtre-Dame, 1697.
eprint: https://books.google.fr/books?id=IbTtzq_fixAC.
- [Hal09] E. Halley. “A new, exact, and easy Method of finding the Roots of any Equations generally, and that without any previous Reduction”. In: *The Philosophical Transactions of the Royal Society of London, from their commencement, in 1665, to the year 1800; Abridged, with notes and biographic illustrations*. Ed. by C. Hutton, G. Shaw and R. Pearson. Vol. III from 1683 to 1694. Translated from the Latin [Hal94]. 1809, pp. 640–649.
- [Hal94] E. Halley. “Methodus Nova Accurata & Facilis Inveniendi Radices Æquationum quarumcumque generaliter, sine prævia Reductione”. In: *Philosophical Transactions of the Royal Society* 18.210 (May 1694), pp. 136–148.
doi: 10.1098/rstl.1694.0029. English translation: [Hal09].
- [Higo2] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, 2002.
- [Hou68] A. S. Householder. “Bigradients and the Problem of Routh and Hurwitz”. In: *SIAM Review* 10.1 (Jan. 1968), pp. 56–66.
doi: 10.1137/1010003.
- [Hou70] A. S. Householder. *The Numerical Treatment of a Single Nonlinear Equation*. International Series in Pure and Applied Mathematics. McGraw-Hill, 1970.
- [KB01] W. Kahan and D. Bindel. “Computing a Real Cube Root”. 2001 retypesetting by Bindel of a purported 1991 version by Kahan, at <https://cscclub.uwaterloo.ca/~pbarfuss/qbrt.pdf>. 21st Apr. 2001.

- [Sch70] E. Schröder. “Ueber unendlich viele Algorithmen zur Auflösung der Gleichungen”. In: *Mathematische Annalen* 2 (June 1870), pp. 317–365. DOI: 10.1007/BF01444024. English translation: [SS93].
- [SG01] P. Sebah and X. Gourdon. “Newton’s method and high order iterations”. In: *Numbers, constants and computation*. 3rd Oct. 2001. eprint: <http://numbers.computation.free.fr/Constants/Algorithms/newton.html>.
- [SS93] E. Schröder and G. W. Stewart. *On Infinitely Many Algorithms for Solving Equations*. Tech. rep. UMIACS-TR-92-121. Translated from the German [Sch70]. Institute for Advanced Computer Studies, University of Maryland, College Park, Jan. 1993. eprint: <http://hdl.handle.net/1903/577>.
- [ST95] T. R. Scavo and J. B. Thoo. “On the Geometry of Halley’s Method”. In: *The American Mathematical Monthly* 102.5 (May 1995), pp. 417–426. DOI: 10.2307/2975033.
- [Ste74] P. H. Sterbenz. *Floating-point computation*. Prentice-Hall, 1974.
- [Tre97] L. N. Trefethen. *Maxims about numerical mathematics, science, computers, and life on Earth*. Maxims. Cornell University, 1997. eprint: <https://people.maths.ox.ac.uk/trefethen/maxims.html>.
- [Wal85] J. Wallis. *A Treatise of Algebra, both Historical and Practical. Shewing, the Original, Progress, and Advancement thereof, from time to time; and by what Steps it hath attained to the Heighth at which now it is*. Richard Davis, 1685. eprint: <https://books.google.fr/books?id=TXpmAAAACAAJ>.