# Rounding Error Analysis of Linear Recurrences
## Using Generating Series

Marc Mezzarobba

CNRS, École polytechnique

Diffferential Seminar @ RTCA6

November 30, 2023

# A Toy Example

# A toy example

[Boldo 2009]

$$c_{n+1} = 2\,c_n - c_{n-1} \qquad (c_0 = \diamond(1/3), c_{-1} = 0)$$

|  | **Floating-point arithmetic** | **Interval arithmetic** |
|---|---|---|
| $n = \;0$ | 0.333333333333333 | $[0.3333333333333333 \pm 1.49\mathrm{e}-17]$ |
| 5 | 2.00000000000000 | $[2.00000000000000 \pm 3.78\mathrm{e}-15]$ |
| 10 | 3.66666666666667 | $[3.666666666667 \pm 5.74\mathrm{e}-13]$ |
| 15 | 5.33333333333334 | $[5.3333333333 \pm 5.29\mathrm{e}-11]$ |
| 20 | 7.00000000000001 | $[7.00000000 \pm 1.60\mathrm{e}-9]$ |
| 25 | 8.66666666666668 | $[8.666667 \pm 4.65\mathrm{e}-7]$ |
| 30 | 10.3333333333333 | $[10.3333 \pm 4.41\mathrm{e}-5]$ |
| 35 | 12.0000000000000 | $[12.000 \pm 8.82\mathrm{e}-4]$ |
| 40 | 13.6666666666667 | $[1.4\mathrm{e}+1 \pm 0.406]$ |
| 45 | 15.33333333333334 | $[\pm 21.3]$ |
| 50 | 17.000000000000 | $[\pm 5.04\mathrm{e}+2]$ |

# Naïve error analysis

Model: $\diamond(x \, op \, y) = x \, op \, y + \varepsilon_{\text{op}}$ with $\varepsilon_{\text{op}} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ fixed-point arithmetic)

(multiplication by 2 is exact)

Exact rec.: $$c_{n+1} = 2 c_n - c_{n-1}$$

Model: $\diamond(x \, op \, y) = x \, op \, y + \varepsilon_{op}$ with $\varepsilon_{op} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ fixed-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2\tilde{c}_n - \tilde{c}_{n-1})$$

# Naïve error analysis

Model: $\diamond(x \, op \, y) = x \, op \, y + \varepsilon_{\mathrm{op}}$ with $\varepsilon_{\mathrm{op}} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ fixed-point arithmetic)

(multiplication by 2 is exact)

Exact rec.: 
$$c_{n+1} = 2c_n - c_{n-1}$$

Approx. rec.: 
$$\tilde{c}_{n+1} = \diamond(2\tilde{c}_n - \tilde{c}_{n-1})$$
$$= 2\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u}$$

Model: $\diamond(x \ op \ y) = x \ op \ y + \varepsilon_{op}$ with $\varepsilon_{op} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ fixed-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2\,c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond\big(2\,\tilde{c}_n - \tilde{c}_{n-1}\big)$$
$$= 2\,\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u}$$

$$|\tilde{c}_{n+1} - c_{n+1}| \leqslant 2\,|\tilde{c}_n - c_n| + |\tilde{c}_{n-1} - c_{n-1}| + \mathbf{u}$$

## Naïve error analysis

Model:   $\diamond(x \; op \; y) = x \; op \; y + \varepsilon_{op}$ with $\varepsilon_{op} \in [-\mathbf{u}, \mathbf{u}]$   ($\sim$ fixed-point arithmetic)

   (multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2 c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond\big(2\tilde{c}_n - \tilde{c}_{n-1}\big)$$
$$= 2\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u}$$

$$|\tilde{c}_{n+1} - c_{n+1}| \leqslant 2|\tilde{c}_n - c_n| + |\tilde{c}_{n-1} - c_{n-1}| + \mathbf{u}$$

Induction:
$$|\tilde{c}_n - c_n| \leqslant 3^n \, \mathbf{u}$$

# Naïve error analysis

Model: $\diamond(x \, op \, y) = x \, op \, y + \varepsilon_{op}$ with $\varepsilon_{op} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ fixed-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2 c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond\big(2 \tilde{c}_n - \tilde{c}_{n-1}\big)$$
$$= 2 \tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u}$$

$$|\tilde{c}_{n+1} - c_{n+1}| \leqslant 2 |\tilde{c}_n - c_n| + |\tilde{c}_{n-1} - c_{n-1}| + \mathbf{u}$$

Induction:
$$|\tilde{c}_n - c_n| \leqslant 3^n \mathbf{u}$$

🙁

# Naïve error analysis

Model: $\diamond(x \, op \, y) = x \, op \, y + \varepsilon_{op}$ with $\varepsilon_{op} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ fixed-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2 c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2 \tilde{c}_n - \tilde{c}_{n-1})$$
$$= 2 \tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with} \quad |\varepsilon_n| \leqslant \mathbf{u}$$

$$|\tilde{c}_{n+1} - c_{n+1}| \leqslant 2 |\tilde{c}_n - c_n| + |\tilde{c}_{n-1} - c_{n-1}| + \mathbf{u}$$

Induction:
$$|\tilde{c}_n - c_n| \leqslant 3^n \, \mathbf{u}$$
😦

Slightly better:
$$|\tilde{c}_n - c_n| \leqslant (\lambda_+ \alpha_+^n + \lambda_- \alpha_-^n - 4) \, \mathbf{u} \approx 2.4^n \, \mathbf{u} \qquad \alpha_\pm = 1 \pm \sqrt{2}$$
$$\approx 2.4^n \, \mathbf{u} \qquad \lambda_\pm = 4 \pm 3\sqrt{2}$$

# Naïve error analysis [3]

Model: $\diamond(x \, op \, y) = x \, op \, y + \varepsilon_{op}$ with $\varepsilon_{op} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ fixed-point arithmetic)

(multiplication by 2 is exact)

Exact rec.: $$c_{n+1} = 2c_n - c_{n-1}$$

Approx. rec.: $$\tilde{c}_{n+1} = \diamond(2\tilde{c}_n - \tilde{c}_{n-1})$$
$$= 2\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with} \quad |\varepsilon_n| \leqslant \mathbf{u}$$

$$|\tilde{c}_{n+1} - c_{n+1}| \leqslant 2|\tilde{c}_n - c_n| + |\tilde{c}_{n-1} - c_{n-1}| + \mathbf{u}$$

Induction: $$|\tilde{c}_n - c_n| \leqslant 3^n \, \mathbf{u}$$ 🙁

Slightly better: $$|\tilde{c}_n - c_n| \leqslant (\lambda_+ \alpha_+^n + \lambda_- \alpha_-^n - 4)\, \mathbf{u} \approx 2.4^n \, \mathbf{u} \qquad \alpha_\pm = 1 \pm \sqrt{2}$$
$$\approx 2.4^n \, \mathbf{u} \qquad\qquad\qquad\qquad \lambda_\pm = 4 \pm 3\sqrt{2}$$
🙁

This is what interval evaluation amounts to!

# Errors cancel out

Exact rec.:
$$c_{n+1} = 2c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2\tilde{c}_n - \tilde{c}_{n-1})$$
$$= 2\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant u$$

💡 The error satisfies "the same" recurrence as the computed sequence

# Errors cancel out

Exact rec.:
$$c_{n+1} = 2 c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond\left(2 \tilde{c}_n - \tilde{c}_{n-1}\right)$$
$$= 2 \tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant u$$

Let $\delta_n = \tilde{c}_n - c_n$ :

💡 The error satisfies "the same" recurrence as the computed sequence

# Errors cancel out

Exact rec.:
$$c_{n+1} = 2c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond\left(2\tilde{c}_n - \tilde{c}_{n-1}\right)$$
$$= 2\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant u$$

Let $\delta_n = \tilde{c}_n - c_n$ :
$$\delta_{n+1} = 2\delta_n - \delta_{n-1} + \varepsilon_n \qquad\qquad (\delta_0 = \delta_1 = 0)$$

💡 The error satisfies "the same" recurrence as the computed sequence

# Errors cancel out

Exact rec.: 
$$c_{n+1} = 2\,c_n - c_{n-1}$$

Approx. rec.: 
$$\tilde{c}_{n+1} = \diamond\!\left(2\,\tilde{c}_n - \tilde{c}_{n-1}\right)$$
$$= 2\,\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant u$$

Let $\delta_n = \tilde{c}_n - c_n$ :
$$\delta_{n+1} = 2\,\delta_n - \delta_{n-1} + \varepsilon_n \qquad\qquad (\delta_0 = \delta_1 = 0)$$
$$\curvearrowleft \textbf{ "local error"}$$

💡 The error satisfies "the same" recurrence as the computed sequence

# Errors cancel out

Exact rec.: $$c_{n+1} = 2\,c_n - c_{n-1}$$

Approx. rec.: $$\tilde{c}_{n+1} = \diamond\!\left(2\,\tilde{c}_n - \tilde{c}_{n-1}\right)$$
$$= 2\,\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant u$$

Let $\delta_n = \tilde{c}_n - c_n$ : $\qquad \delta_{n+1} = 2\,\delta_n - \delta_{n-1} + \varepsilon_n \qquad\qquad (\delta_0 = \delta_1 = 0)$

**"global error"** $\nearrow$ $\qquad\qquad\qquad\qquad$ $\nwarrow$ **"local error"**

💡 The error satisfies "the same" recurrence as the computed sequence

# Errors cancel out

Exact rec.:
$$c_{n+1} = 2c_n - c_{n-1}$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2\tilde{c}_n - \tilde{c}_{n-1})$$
$$= 2\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant u$$

Let $\delta_n = \tilde{c}_n - c_n$ :
$$\delta_{n+1} = 2\,\delta_n - \delta_{n-1} + \varepsilon_n \qquad (\delta_0 = \delta_1 = 0)$$

**"global error"** ↗       ↖ **"local error"**

Then
$$\delta_n = \sum_{k=1}^{n-1} k\,\varepsilon_{n-k}$$
$$|\delta_n| \leqslant \frac{n(n-1)}{2}u$$

😊

Calculations can become unwieldy (nested sums, determinants…)

💡 The error satisfies "the same" recurrence as the computed sequence

# Errors cancel out

Exact rec.: 
$$c_{n+1} = 2c_n - c_{n-1}$$

Approx. rec.: 
$$\tilde{c}_{n+1} = \diamond(2\tilde{c}_n - \tilde{c}_{n-1})$$
$$= 2\tilde{c}_n - \tilde{c}_{n-1} + \varepsilon_n \quad \text{with } |\varepsilon_n| \leqslant u$$

Let $\delta_n = \tilde{c}_n - c_n$ : 
$$\delta_{n+1} = 2\,\delta_n - \delta_{n-1} + \varepsilon_n \qquad (\delta_0 = \delta_1 = 0)$$
**"global error"** $\nearrow$ $\qquad\qquad\qquad$ $\nwarrow$ **"local error"**

Then 
$$\delta_n = \sum_{k=1}^{n-1} k\,\varepsilon_{n-k}$$
$$|\delta_n| \leqslant \frac{n(n-1)}{2}u \qquad\qquad \smiley$$

Calculations can become unwieldy (nested sums, determinants…)

💡 The error satisfies "the same" recurrence as the computed sequence

💡 Encode sequences by generating functions

$$\text{Sequence } (f_n)_{n \in \mathbb{Z}} \quad \longleftrightarrow \quad \text{Generating series } f(z) = \sum_{n=-\infty}^{\infty} f_n z^n$$

# Generating functions

💡 Encode sequences by generating functions

🖊 Sequence $(f_n)_{n\in\mathbb{Z}} \longleftrightarrow$ Generating series $f(z) = \sum_{n=-\infty}^{\infty} f_n z^n$

✚ Formulae for product, composition, …

✚ Analytic methods (Cauchy integrals)

✚ Fast algorithms

✚ Method of majorants

✚ …

# Generating functions

💡 Encode sequences by generating functions

---

Sequence $(f_n)_{n \in \mathbb{Z}}$ $\longleftrightarrow$ Generating series $f(z) = \sum_{n=-\infty}^{\infty} f_n z^n$

---

$$\delta_{n+1} = 2\,\delta_n - \delta_{n-1} + \varepsilon_n$$

$$\downarrow \sum_n \square\, z^n \qquad\qquad z\sum a_n z^n = \sum a_{n-1} z^n$$

$$z^{-1}\,\delta(z) = 2\,\delta(z) - z\,\delta(z) + \varepsilon(z)$$

✚ Formulae for product, composition, …

✚ Analytic methods (Cauchy integrals)

✚ Fast algorithms

✚ Method of majorants

✚ …

# Generating functions

💡 Encode sequences by generating functions

---

🖊 Sequence $(f_n)_{n \in \mathbb{Z}}$ $\longleftrightarrow$ Generating series $f(z) = \sum_{n=-\infty}^{\infty} f_n z^n$

---

Exact expression of the global error $\delta$:

$$\delta_{n+1} = 2\,\delta_n - \delta_{n-1} + \varepsilon_n$$

$$\downarrow \sum_n \square\, z^n \qquad\qquad z\sum a_n z^n = \sum a_{n-1} z^n$$

$$z^{-1}\,\delta(z) = 2\,\delta(z) - z\,\delta(z) + \varepsilon(z)$$

$$\delta(z) = \frac{z}{(1-z)^2}\,\varepsilon(z)$$

➕ Formulae for product, composition, …

➕ Analytic methods (Cauchy integrals)

➕ Fast algorithms

➕ Method of majorants

➕ …

# Generating functions

💡 Encode sequences by generating functions

✏️ Sequence $(f_n)_{n \in \mathbb{Z}}$ $\longleftrightarrow$ Generating series $f(z) = \sum_{n=-\infty}^{\infty} f_n z^n$

Exact expression of the global error $\delta$:

$$\delta_{n+1} = 2\,\delta_n - \delta_{n-1} + \varepsilon_n$$

$$\downarrow \sum_n \square\, z^n \qquad\qquad z \sum a_n z^n = \sum a_{n-1} z^n$$

$$z^{-1}\,\delta(z) = 2\,\delta(z) - z\,\delta(z) + \varepsilon(z)$$

$$\delta(z) = \frac{z}{(1-z)^2}\,\varepsilon(z)$$

✚ Formulae for product, composition, …
✚ Analytic methods (Cauchy integrals)
✚ Fast algorithms
✚ Method of majorants
✚ …

# Majorants

$$f(z) = \sum_n f_n z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n z^n$$

$$f \ll \hat{f} \text{ (``}\hat{f}\text{ majorizes } f\text{'')} \quad \Leftrightarrow \quad \forall n, |f_n| \leqslant \hat{f}_n$$

$$f(z) = \sum_n f_n z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n z^n$$

$$f \ll \hat{f} \text{ ("}\hat{f}\text{ majorizes f")} \quad \Leftrightarrow \quad \forall n, |f_n| \leqslant \hat{f}_n$$

Bound on the global error:
($\sim$ same proof, new notation)

$$\delta(z) = \frac{z}{(1-z)^2} \varepsilon(z)$$

# Majorants

$$f(z) = \sum_n f_n z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n z^n$$

$$f \ll \hat{f} \ (\text{"}\hat{f}\text{ majorizes }f\text{"}) \quad \Leftrightarrow \quad \forall n, |f_n| \leqslant \hat{f}_n$$

Bound on the global error:
($\sim$ same proof, new notation)

$$\delta(z) = \frac{z}{(1-z)^2}\, \varepsilon(z)$$

If $f(z) \ll \hat{f}(z)$ and $g(z) \ll \hat{g}(z)$,
then $f(z)\, g(z) \ll \hat{f}(z)\, \hat{g}(z)$

$$f(z) = \sum_n f_n z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n z^n$$

$$f \ll \hat{f} \text{ ("}\hat{f} \text{ majorizes } f\text{")} \quad \Leftrightarrow \quad \forall n, \; |f_n| \leqslant \hat{f}_n$$

Bound on the global error:
($\sim$ same proof, new notation)

$$\delta(z) \;=\; \frac{z}{(1-z)^2} \, \varepsilon(z)$$

If $f(z) \ll \hat{f}(z)$ and $g(z) \ll \hat{g}(z)$,
then $f(z)\,g(z) \ll \hat{f}(z)\,\hat{g}(z)$

Proof: 
$$\left| [z^n]\big(f(z)\,g(z)\big) \right| = \left| \sum_{i+j=n} f_i\, g_j \right|$$
$$\leqslant \sum_{i+j=n} \hat{f}_i\, \hat{g}_j$$
$$= [z^n]\big(\hat{f}(z)\,\hat{g}(z)\big)$$

# Majorants

$$f(z) = \sum_n f_n\, z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n\, z^n$$

$$f \ll \hat{f} \ (\text{``}\hat{f}\text{ majorizes }f\text{''}) \quad \Leftrightarrow \quad \forall n,\ |f_n| \leqslant \hat{f}_n$$

Bound on the global error:
($\sim$ same proof, new notation)

$$\delta(z) = \frac{z}{(1-z)^2}\, \varepsilon(z)$$

🧩 If $f(z) \ll \hat{f}(z)$ and $g(z) \ll \hat{g}(z)$,
then $f(z)\, g(z) \ll \hat{f}(z)\, \hat{g}(z)$

Proof:
$$\left| [z^n]\big(f(z)\, g(z)\big) \right| = \left| \sum_{i+j=n} f_i\, g_j \right|$$
$$\leqslant \sum_{i+j=n} \hat{f}_i\, \hat{g}_j$$
$$= [z^n]\big(\hat{f}(z)\, \hat{g}(z)\big)$$

🧩 $\displaystyle \varepsilon(z) = \sum_n \varepsilon_n\, z^n \ll \frac{u}{1-z}$

# Majorants

$$f(z) = \sum_n f_n z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n z^n$$

$$f \ll \hat{f} \text{ ("$\hat{f}$ majorizes $f$")} \quad \Leftrightarrow \quad \forall n, |f_n| \leqslant \hat{f}_n$$

Bound on the global error:
(~ same proof, new notation)

$$\delta(z) = \frac{z}{(1-z)^2}\, \varepsilon(z)$$

♣ If $f(z) \ll \hat{f}(z)$ and $g(z) \ll \hat{g}(z)$,
then $f(z)\, g(z) \ll \hat{f}(z)\, \hat{g}(z)$

Proof: $\left| [z^n]\left(f(z)\, g(z)\right) \right| = \left| \sum_{i+j=n} f_i\, g_j \right|$
$$\leqslant \sum_{i+j=n} \hat{f}_i\, \hat{g}_j$$
$$= [z^n]\left(\hat{f}(z)\, \hat{g}(z)\right)$$

♣ $\varepsilon(z) = \sum_n \varepsilon_n z^n \ll \dfrac{u}{1-z}$

♣ $\dfrac{z}{(1-z)^2} \in \mathbb{R}_+[[z]]$

$$f(z) = \sum_n f_n z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n z^n$$

$$f \ll \hat{f} \;(\text{"}\hat{f}\text{ majorizes f"}) \quad \Leftrightarrow \quad \forall n, |f_n| \leqslant \hat{f}_n$$

Bound on the global error:
($\sim$ same proof, new notation)

$$\begin{aligned}
\delta(z) &= \frac{z}{(1-z)^2}\, \varepsilon(z) \\
&\ll \frac{u\,z}{(1-z)^3}
\end{aligned}$$

If $f(z) \ll \hat{f}(z)$ and $g(z) \ll \hat{g}(z)$,
then $f(z)\,g(z) \ll \hat{f}(z)\,\hat{g}(z)$

Proof:
$$\begin{aligned}
\left| [z^n]\,\big(f(z)\,g(z)\big) \right| &= \left| \sum_{i+j=n} f_i\, g_j \right| \\
&\leqslant \sum_{i+j=n} \hat{f}_i\, \hat{g}_j \\
&= [z^n]\,\big(\hat{f}(z)\,\hat{g}(z)\big)
\end{aligned}$$

$$\varepsilon(z) = \sum_n \varepsilon_n z^n \ll \frac{u}{1-z}$$

$$\frac{z}{(1-z)^2} \in \mathbb{R}_+[[z]]$$

# Majorants

$$f(z) = \sum_n f_n z^n, \quad \hat{f}(z) = \sum_n \hat{f}_n z^n$$

$$f \ll \hat{f} \text{ ("$\hat{f}$ majorizes $f$")} \quad \Leftrightarrow \quad \forall n, \ |f_n| \leqslant \hat{f}_n$$

Bound on the global error:
($\sim$ same proof, new notation)

$$\begin{aligned}
\delta(z) &= \frac{z}{(1-z)^2} \, \varepsilon(z) \\
&\ll \frac{u \, z}{(1-z)^3}
\end{aligned}$$

$$|\delta_n| \leqslant \frac{n \, (n-1)}{2} \, u \quad \text{☺}$$

If $f(z) \ll \hat{f}(z)$ and $g(z) \ll \hat{g}(z)$,
then $f(z) \, g(z) \ll \hat{f}(z) \, \hat{g}(z)$

Proof: $\left| [z^n] \left( f(z) \, g(z) \right) \right| = \left| \sum_{i+j=n} f_i \, g_j \right|$
$$\leqslant \sum_{i+j=n} \hat{f}_i \, \hat{g}_j$$
$$= [z^n] \left( \hat{f}(z) \, \hat{g}(z) \right)$$

$$\varepsilon(z) = \sum_n \varepsilon_n z^n \ll \frac{u}{1-z}$$

$$\frac{z}{(1-z)^2} \in \mathbb{R}_+[[z]]$$

# Related work

> in the backward direction. There has been less attention devoted to computation which utilizes the difference equation in the forward direction, not because a forward algorithm is more difficult to analyze, but <mark>rather for the opposite reason—that its analysis was considered straightforward.</mark> Of the above

[Wimp 1972]

**'Linear' error propagation**

| | |
|---|---|
| Henrici 1962 | finite difference schemes for ODE |
| Oliver 1967 | linear recurrences |

**Explicit bounds** (not necessarily easy to compute)

| | |
|---|---|
| von Neumann & Goldstine 1947, Turing 1948 | triangular system solving |
| Elliott 1968 | sums of generalized Fourier series |
| Wimp 1972 | order 2 |
| Barrio & Melendo & Serrano 2003 | order $n$, $O(u^2)$ |

**Transfer functions of digital filters**

| | |
|---|---|
| Liu & Kaneko 1969 | random errors |
| Hilaire & Lopez 2013 | error bounds |

# Relative error propagation

Model: $\diamond(x \ op \ y) = (x \ op \ y)(1 + \varepsilon_{op})$ with $\varepsilon_{op} \in [-u, u]$ ($\sim$ floating-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2 c_n - c_{n-1} \qquad\qquad \times (1 + \varepsilon_n)$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2 \tilde{c}_n - \tilde{c}_{n-1})$$
$$= (2 \tilde{c}_n - \tilde{c}_{n-1})(1 + \varepsilon_n) \quad \text{with } |\varepsilon_n| \leqslant u$$

With $\delta_n = \tilde{c}_n - c_n$ :

# Relative error propagation

Model: $\diamond(x \; op \; y) = (x \; op \; y)(1 + \varepsilon_{\mathrm{op}})$ with $\varepsilon_{\mathrm{op}} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ floating-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2 c_n - c_{n-1} \qquad \times(1 + \varepsilon_n)$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2 \tilde{c}_n - \tilde{c}_{n-1})$$
$$= (2 \tilde{c}_n - \tilde{c}_{n-1})(1 + \boldsymbol{\varepsilon_n}) \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u} \qquad \times(-1)$$

With $\delta_n = \tilde{c}_n - c_n$:
$$\delta_{n+1} - c_{n+1}\varepsilon_n = (2\delta_n - \delta_{n-1})(1 + \varepsilon_n)$$

# Relative error propagation

Model: $\diamond(x \; op \; y) = (x \; op \; y)(1 + \varepsilon_{\mathrm{op}})$ with $\varepsilon_{\mathrm{op}} \in [-\mathbf{u}, \mathbf{u}]$ ($\sim$ floating-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2\,c_n - c_{n-1} \qquad\qquad \times(1 + \varepsilon_n)$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2\,\tilde{c}_n - \tilde{c}_{n-1})$$
$$= (2\,\tilde{c}_n - \tilde{c}_{n-1})\,(1 + \boldsymbol{\varepsilon_n}) \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u} \qquad \times(-1)$$

With $\delta_n = \tilde{c}_n - c_n$:
$$\delta_{n+1} - c_{n+1}\,\varepsilon_n = (2\,\delta_n - \delta_{n-1})\,(1 + \varepsilon_n)$$
$$\delta_{n+1} - 2\,\delta_n + \delta_{n-1} = \varepsilon_n\,(c_{n+1} + 2\,\delta_n - \delta_{n-1})$$

# Relative error propagation

Model: $\diamond(x \; op \; y) = (x \; op \; y)(1 + \varepsilon_{op})$ with $\varepsilon_{op} \in [-u, u]$ ($\sim$ floating-point arithmetic)

(multiplication by 2 is exact)

Exact rec.:
$$c_{n+1} = 2c_n - c_{n-1} \qquad \times(1 + \varepsilon_n)$$

Approx. rec.:
$$\tilde{c}_{n+1} = \diamond(2\tilde{c}_n - \tilde{c}_{n-1})$$
$$= (2\tilde{c}_n - \tilde{c}_{n-1})(1 + \varepsilon_n) \quad \text{with } |\varepsilon_n| \leqslant u \qquad \times(-1)$$

With $\delta_n = \tilde{c}_n - c_n$:
$$\delta_{n+1} - c_{n+1}\varepsilon_n = (2\delta_n - \delta_{n-1})(1 + \varepsilon_n)$$
$$\delta_{n+1} - 2\delta_n + \delta_{n-1} = \varepsilon_n(c_{n+1} + 2\delta_n - \delta_{n-1})$$

Translate:
$$\Big\downarrow \sum_n \square z^n$$
$$(z^{-1} - 2 + z)\,\delta(z) = \varepsilon(z) \odot \big(z^{-1}c(z) + (2 - z)\,\delta(z)\big)$$

# Relative error propagation

Model: $\diamond(x \; op \; y) = (x \; op \; y)(1 + \varepsilon_{\mathrm{op}})$ with $\varepsilon_{\mathrm{op}} \in [-\mathbf{u}, \mathbf{u}]$  ($\sim$ floating-point arithmetic)

(multiplication by 2 is exact)

Exact rec.: 
$$c_{n+1} = 2 c_n - c_{n-1} \qquad\qquad \times(1 + \varepsilon_n)$$

Approx. rec.: 
$$\tilde{c}_{n+1} = \diamond(2 \tilde{c}_n - \tilde{c}_{n-1})$$
$$= (2 \tilde{c}_n - \tilde{c}_{n-1})(1 + \varepsilon_n) \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u} \qquad \times(-1)$$

With $\delta_n = \tilde{c}_n - c_n$ : 
$$\delta_{n+1} - c_{n+1} \varepsilon_n = (2 \delta_n - \delta_{n-1})(1 + \varepsilon_n)$$
$$\delta_{n+1} - 2 \delta_n + \delta_{n-1} = \varepsilon_n (c_{n+1} + 2 \delta_n - \delta_{n-1})$$

Translate:
$$\Big\downarrow \; \sum_n \square \, z^n$$

$$(z^{-1} - 2 + z) \, \delta(z) = \varepsilon(z) \odot \big( z^{-1} c(z) + (2 - z) \, \delta(z) \big)$$

Solve for $\delta$ ?? 
$$\delta(z) = \frac{(z \, \varepsilon(z)) \odot \big( c(z) + z \, (2 - z) \, \delta(z) \big)}{(1 - z)^2} \qquad \text{☹}$$

# Relative error propagation

Model: $\diamond(x\ op\ y) = (x\ op\ y)(1 + \varepsilon_{\mathrm{op}})$ with $\varepsilon_{\mathrm{op}} \in [-\mathbf{u}, \mathbf{u}]$    ($\sim$ floating-point arithmetic)

(multiplication by 2 is exact)

Exact rec.: $$c_{n+1} = 2\,c_n - c_{n-1} \qquad\qquad\qquad \times(1 + \varepsilon_n)$$

Approx. rec.: $$\tilde{c}_{n+1} = \diamond(2\,\tilde{c}_n - \tilde{c}_{n-1})$$
$$= (2\,\tilde{c}_n - \tilde{c}_{n-1})\,(1 + \boldsymbol{\varepsilon_n}) \quad \text{with } |\varepsilon_n| \leqslant \mathbf{u} \qquad \times(-1)$$

With $\delta_n = \tilde{c}_n - c_n$: $$\delta_{n+1} - c_{n+1}\,\varepsilon_n = (2\,\delta_n - \delta_{n-1})\,(1 + \varepsilon_n)$$
$$\delta_{n+1} - 2\,\delta_n + \delta_{n-1} = \varepsilon_n\,(c_{n+1} + 2\,\delta_n - \delta_{n-1})$$

Translate:

$$\Big\downarrow \sum_n \square\, z^n$$

$$(z^{-1} - 2 + z)\,\delta(z) = \varepsilon(z) \odot \big(z^{-1}\,c(z) + (2 - z)\,\delta(z)\big)$$

Solve for $\delta$ ?? $$\delta(z) = \frac{(z\,\varepsilon(z)) \odot \big(c(z) + z\,(2 - z)\,\delta(z)\big)}{(1 - z)^2}$$ :(

$$^{\sharp}f(z) = \sum_n |f_n|\, z^n \qquad\qquad \delta(z) \ll \frac{z\,(2 + z)\,\mathbf{u}}{(1 - z)^2}\, {}^{\sharp}\delta(z) + \frac{{}^{\sharp}c(z)\,\mathbf{u}}{(1 - z)^2}$$ :(

# Majorizing equations

💡 Obtain the bound as a solution of a "similar" equation
$$\delta(z) \ll \frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}{}^{\sharp}\delta(z) + \frac{{}^{\sharp}\mathbf{c}(z)\,\mathbf{u}}{(1-z)^2}$$

---

**Lemma.** [$\sim$ Cauchy]

Let $\hat{a}(z), \hat{b}(z) \in \mathbb{R}_+[[z]]$ with $\hat{a}(0) = 0$. Suppose $y \in \mathbb{R}_+[[z]]$ satisfies

$$y(z) \ll \hat{a}(z)\,y(z) + \hat{b}(z).$$

Then $y(z)$ is majorized by the solution of $\hat{y}(z) = \hat{a}(z)\,\hat{y}(z) + \hat{b}(z)$, i.e.,

$$y(z) \ll \hat{y}(z) = \frac{\hat{b}(z)}{1 - \hat{a}(z)}.$$

---

# Majorizing equations

💡 Obtain the bound as a solution of a "similar" equation

$$\delta(z) \ll \frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}{}^\sharp\delta(z) + \frac{{}^\sharp\mathbf{c}(z)\,\mathbf{u}}{(1-z)^2}$$

---

**Lemma.** [$\sim$ Cauchy]

Let $\hat{a}(z), \hat{b}(z) \in \mathbb{R}_+[[z]]$ with $\hat{a}(0) = 0$. Suppose $y \in \mathbb{R}_+[[z]]$ satisfies

$$y(z) \ll \hat{a}(z)\,y(z) + \hat{b}(z).$$

Then $y(z)$ is majorized by the solution of $\hat{y}(z) = \hat{a}(z)\,\hat{y}(z) + \hat{b}(z)$, i.e.,

$$y(z) \ll \hat{y}(z) = \frac{\hat{b}(z)}{1 - \hat{a}(z)}.$$

---

**Proof.**

• $y_0 \leqslant \hat{b}_0 = \hat{y}_0$

• $y_n \leqslant \sum_{i=0}^{n} \hat{a}_i\,y_{n-i} + \hat{b}_n$ □

# Majorizing equations

💡 Obtain the bound as a solution of a "similar" equation $\qquad \delta(z) \ll \dfrac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}{}^\sharp\delta(z) + \dfrac{{}^\sharp\mathbf{c}(z)\,\mathbf{u}}{(1-z)^2}$

---

**Lemma.** [$\sim$ Cauchy]

Let $\hat{a}(z), \hat{b}(z) \in \mathbb{R}_+[[z]]$ with $\hat{a}(0) = 0$. Suppose $y \in \mathbb{R}_+[[z]]$ satisfies

$$y(z) \ll \hat{a}(z)\,y(z) + \hat{b}(z).$$

Then $y(z)$ is majorized by the solution of $\hat{y}(z) = \hat{a}(z)\,\hat{y}(z) + \hat{b}(z)$, i.e.,

$$y(z) \ll \hat{y}(z) = \frac{\hat{b}(z)}{1 - \hat{a}(z)}.$$

---

**Proof.**

• $y_0 \leqslant \hat{b}_0 = \hat{y}_0$

• $y_n \leqslant \displaystyle\sum_{i=1}^{n} \hat{a}_i\,y_{n-i} + \hat{b}_n$ $\hfill \square$

# Majorizing equations

💡 Obtain the bound as a solution of a "similar" equation
$$\delta(z) \ll \frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}{}^\sharp\delta(z) + \frac{{}^\sharp c(z)\,\mathbf{u}}{(1-z)^2}$$

---

**Lemma.** [$\sim$ Cauchy]

Let $\hat{a}(z), \hat{b}(z) \in \mathbb{R}_+[[z]]$ with $\hat{a}(0) = 0$. Suppose $y \in \mathbb{R}_+[[z]]$ satisfies

$$y(z) \ll \hat{a}(z)\,y(z) + \hat{b}(z).$$

Then $y(z)$ is majorized by the solution of $\hat{y}(z) = \hat{a}(z)\,\hat{y}(z) + \hat{b}(z)$, i.e.,

$$y(z) \ll \hat{y}(z) = \frac{\hat{b}(z)}{1 - \hat{a}(z)}.$$

---

**Proof.**

- $y_0 \leqslant \hat{b}_0 = \hat{y}_0$

- $y_n \leqslant \sum_{i=1}^{n} \hat{a}_i\,y_{n-i} + \hat{b}_n \leqslant \sum_{i=1}^{n} \hat{a}_i\,\hat{y}_{n-i} + \hat{b}_n = \hat{y}_n$ $\qquad\qquad\qquad\square$

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathfrak{u}}{(1-z)^2}}_{\hat{\mathfrak{a}}(z)} \sharp\delta(z) + \underbrace{\frac{\sharp\mathfrak{c}(z)\,\mathfrak{u}}{(1-z)^2}}_{\hat{\mathfrak{b}}(z)}$$

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\,\sharp\delta(z) + \underbrace{\frac{\sharp c(z)\,\mathbf{u}}{(1-z)^2}}_{\hat{b}(z)}$$

$$c(z) = \frac{c_0}{(1-z)^2}$$
$$\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z)$$

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\,\sharp\delta(z) + \underbrace{\frac{\hat{c}(z)\ \ \mathbf{u}}{(1-z)^2}}_{\hat{b}(z)}$$

$$c(z) = \frac{c_0}{(1-z)^2}$$
$$\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z)$$

# Bound on the floating-point error

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\,\sharp\delta(z) + \underbrace{\frac{\hat{c}(z)\;\mathbf{u}}{(1-z)^2}}_{\hat{b}(z)}$$

$$c(z) = \frac{c_0}{(1-z)^2}$$
$$\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z)$$

By the lemma:

$$\delta(z) \ll \frac{\hat{b}(z)}{1-\hat{a}(z)}$$

# Bound on the floating-point error

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\,\sharp\delta(z) + \underbrace{\frac{\hat{c}(z)\ \mathbf{u}}{(1-z)^2}}_{\hat{b}(z)}$$

$$c(z) = \frac{c_0}{(1-z)^2}$$

$$\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z)$$

By the lemma:

$$\delta(z) \ll \frac{\hat{b}(z)}{1-\hat{a}(z)}$$

$$= \frac{\hat{c}(z)\,\mathbf{u}}{1 - 2\,(1+\mathbf{u})\,z + (1-\mathbf{u})\,z^2}$$

# Bound on the floating-point error

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\,\sharp\delta(z) + \underbrace{\frac{\hat{c}(z)\;\mathbf{u}}{(1-z)^2}}_{\hat{b}(z)}$$

$$c(z) = \frac{c_0}{(1-z)^2}$$
$$\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z)$$

By the lemma:
$$\delta(z) \ll \frac{\hat{b}(z)}{1-\hat{a}(z)}$$
$$= \frac{\hat{c}(z)\,\mathbf{u}}{1-2\,(1+\mathbf{u})\,z + (1-\mathbf{u})\,z^2}$$
$$= \frac{\hat{c}(z)\,\mathbf{u}}{(1-\boldsymbol{\alpha}\,z)\,(1-\boldsymbol{\beta}\,z)}$$

$$\boldsymbol{\alpha} = 1 + 2\sqrt{\mathbf{u}} + O(\mathbf{u})$$

# Bound on the floating-point error

$$\overset{\sharp}{\delta}(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\,\overset{\sharp}{\delta}(z) + \underbrace{\frac{\hat{c}(z)\;\mathbf{u}}{(1-z)^2}}_{\hat{b}(z)} \qquad\qquad \begin{aligned} c(z) &= \frac{c_0}{(1-z)^2} \\ &\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z) \end{aligned}$$

By the lemma:
$$\begin{aligned} \delta(z) &\ll \frac{\hat{b}(z)}{1-\hat{a}(z)} \\ &= \frac{\hat{c}(z)\,\mathbf{u}}{1 - 2\,(1+\mathbf{u})\,z + (1-\mathbf{u})\,z^2} \\ &= \frac{\hat{c}(z)\,\mathbf{u}}{(1-\alpha\,z)\,(1-\beta\,z)} \qquad\qquad \alpha = 1 + 2\,\sqrt{\mathbf{u}} + O(\mathbf{u}) \\ &\ll \frac{|c_0|\,\mathbf{u}}{(1-\alpha\,z)^4} \end{aligned}$$

# Bound on the floating-point error

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\, \sharp\delta(z) + \underbrace{\frac{\hat{c}(z)\ \mathbf{u}}{(1-z)^2}}_{\hat{b}(z)}$$

$$c(z) = \frac{c_0}{(1-z)^2}$$
$$\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z)$$

By the lemma:
$$\delta(z) \ll \frac{\hat{b}(z)}{1-\hat{a}(z)}$$
$$= \frac{\hat{c}(z)\,\mathbf{u}}{1 - 2\,(1+\mathbf{u})\,z + (1-\mathbf{u})\,z^2}$$
$$= \frac{\hat{c}(z)\,\mathbf{u}}{(1-\boldsymbol{\alpha}\,z)\,(1-\boldsymbol{\beta}\,z)} \qquad \boldsymbol{\alpha} = 1 + 2\sqrt{\mathbf{u}} + O(\mathbf{u})$$
$$\ll \frac{|c_0|\,\mathbf{u}}{(1-\boldsymbol{\alpha}\,z)^4}$$

Absolute error on $c_n$:
$$|\delta_n| \leqslant \frac{|c_0|}{6}\,(n+3)^3\,\boldsymbol{\alpha}^n\,\mathbf{u}$$

# Bound on the floating-point error

$$\sharp\delta(z) \ll \underbrace{\frac{z\,(2+z)\,\mathbf{u}}{(1-z)^2}}_{\hat{a}(z)}\,\sharp\delta(z) + \underbrace{\frac{\hat{c}(z)\,\mathbf{u}}{(1-z)^2}}_{\hat{b}(z)}$$

$$c(z) = \frac{c_0}{(1-z)^2}$$
$$\ll \frac{|c_0|}{(1-z)^2} =: \hat{c}(z)$$

By the lemma:

$$\begin{aligned}
\delta(z) &\ll \frac{\hat{b}(z)}{1-\hat{a}(z)} \\
&= \frac{\hat{c}(z)\,\mathbf{u}}{1 - 2\,(1+\mathbf{u})\,z + (1-\mathbf{u})\,z^2} \\
&= \frac{\hat{c}(z)\,\mathbf{u}}{(1-\boldsymbol{\alpha}\,z)\,(1-\boldsymbol{\beta}\,z)} \qquad\qquad \boldsymbol{\alpha} = 1 + 2\sqrt{\mathbf{u}} + O(\mathbf{u}) \\
&\ll \frac{|c_0|\,\mathbf{u}}{(1-\boldsymbol{\alpha}\,z)^4}
\end{aligned}$$

Absolute error on $c_n$: $\qquad |\delta_n| \leqslant \frac{|c_0|}{6}\,(n+3)^3\,\alpha^n\,\mathbf{u}$

Exponential for fixed $\mathbf{u}$, but $O(n^3\,\mathbf{u})$ if $n = O(\mathbf{u}^{-1/2})$

# Differential equations

# Evaluation of Legendre polynomials

```
Algorithm 1. Evaluation of Legendre polynomials in GMP fixed-point arithmetic.
Input: An integer x and t ≥ 0 such that |2⁻ᵗx| ≤ 1, and n ≥ 1
Output: p, q such that |2⁻ᵗp − Pₙ₋₁(2⁻ᵗx)|, |2⁻ᵗq − Pₙ(2⁻ᵗx)| ≤ (0.75 (n+1)(n+2)+1) 2⁻ᵗ
1: void legendre(mpz_t p, mpz_t q, int n, const mpz_t x, int t) {
2:     mpz_t tmp; int k; mpz_init(tmp);                      ▷ Comments use the notation of
3:     mp_limb_t denlo, den = 1;                             ▷ the proof of Corollary 6
4:     mpz_set_ui(p, 1); mpz_mul_2exp(p, p, t);                            ▷ p₀ = 2ᵗ
5:     mpz_set(q, x);                                                      ▷ q₀ = x̂
6:     for (k = 1; k < n; k++) {
7:         mpz_mul(tmp, q, x); mpz_tdiv_q_2exp(tmp, tmp, t);   ▷ ⌈x̂ q_{k−1} 2⁻ᵗ⌋
8:         mpz_mul_si(p, p, -k*k);
9:         mpz_addmul_ui(p, tmp, 2*k+1);                  ▷ −k²p_{k−1} + (2k + 1) tmp
10:        mpz_swap(p, q);
11:        if (mpn_mul_1(&denlo, &den, 1, k+1)) {   ▷ If multiplication overflows
12:            mpz_tdiv_q_ui(p, p, den);                           ▷ ⌈p/d_{k−1}⌋
13:            mpz_tdiv_q_ui(q, q, den);
14:            den = k+1;                                          ▷ d_k = k + 1
15:        } else den = denlo;                             ▷ d_k = (k + 1) d_{k−1}
16:    }
17:    mpz_tdiv_q_ui(p, p, den/n); mpz_tdiv_q_ui(q, q, den);
18:    mpz_clear(tmp);
19: }
```

Context: rigorous arbitrary-precision
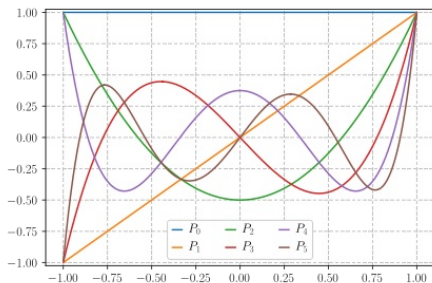   Gauss-Legendre quadrature
   [Johansson 2018]

$$P_{n+1}(x) = \frac{1}{n+1}\left[ (2n+1)\, x\, P_n(x) - n\, P_{n-1}(x) \right]$$

$\tilde{p}_n = P_n(x)$ evaluated using this recurrence
in t-bit fixed-point arithmetic

Bound $|\tilde{p}_n - P_n(x)|$.

# Legendre polynomials: error analysis

Exact rec.:
$$p_{n+1} = \frac{1}{n+1}\big[(2n+1)\,x\,p_n - n\,p_{n-1}\big] \qquad\qquad p_n := P_n(x)$$

Approx. rec.:
$$\tilde{p}_{n+1} = \frac{1}{n+1}\big[(2n+1)\,x\,\tilde{p}_n - n\,\tilde{p}_{n-1}\big] + \varepsilon_{n+1} \qquad \text{with } \varepsilon_n \leqslant 3\,u$$

# Legendre polynomials: error analysis

Exact rec.:
$$p_{n+1} = \frac{1}{n+1}\big[(2n+1)x\,p_n - n\,p_{n-1}\big] \qquad\qquad p_n := P_n(x)$$

Approx. rec.:
$$\tilde{p}_{n+1} = \frac{1}{n+1}\big[(2n+1)x\,\tilde{p}_n - n\,\tilde{p}_{n-1}\big] + \varepsilon_{n+1} \qquad \text{with } \varepsilon_n \leqslant 3\,u$$

Global error:

$\delta_n = \tilde{p}_n - p_n$ $\qquad\qquad (n+1)\,\delta_{n+1} = (2n+1)x\,\delta_n - n\,\delta_{n-1} + (n+1)\,\varepsilon_{n+1}$

# Legendre polynomials: error analysis

Exact rec.:
$$p_{n+1} = \frac{1}{n+1}\big[(2n+1)x\,p_n - n\,p_{n-1}\big] \qquad\qquad p_n := P_n(x)$$

Approx. rec.:
$$\tilde{p}_{n+1} = \frac{1}{n+1}\big[(2n+1)x\,\tilde{p}_n - n\,\tilde{p}_{n-1}\big] + \varepsilon_{n+1} \qquad \text{with}\,\varepsilon_n \leqslant 3\,\mathbf{u}$$

Global error:

$$\delta_n = \tilde{p}_n - p_n \qquad\qquad (n+1)\,\delta_{n+1} = (2n+1)x\,\delta_n - n\,\delta_{n-1} + (n+1)\,\varepsilon_{n+1}$$

Translate:

$$\sum_n \square\, z^n$$

$$(1 - 2xz + z^2)\,\delta'(z) = z\,(x-z)\,\delta(z) + \varepsilon'(z)$$

# Legendre polynomials: error analysis

Exact rec.:
$$p_{n+1} = \frac{1}{n+1}\big[(2n+1)\,x\,p_n - n\,p_{n-1}\big] \qquad p_n := P_n(x)$$

Approx. rec.:
$$\tilde{p}_{n+1} = \frac{1}{n+1}\big[(2n+1)\,x\,\tilde{p}_n - n\,\tilde{p}_{n-1}\big] + \varepsilon_{n+1} \qquad \text{with } \varepsilon_n \leqslant 3\,\mathbf{u}$$

Global error:

$$\delta_n = \tilde{p}_n - p_n \qquad\qquad (n+1)\,\delta_{n+1} = (2n+1)\,x\,\delta_n - n\,\delta_{n-1} + (n+1)\,\varepsilon_{n+1}$$

Translate:
$$\Big\downarrow \sum_n \square\, z^n$$

$$(1 - 2\,x\,z + z^2)\,\delta'(z) = z\,(x - z)\,\delta(z) + \varepsilon'(z)$$

Solve:
$$\delta(z) = \frac{1}{\sqrt{1 - 2\,x\,z + z^2}} \int_0^z \frac{\varepsilon'(w)}{\sqrt{1 - 2\,x\,w + w^2}}\,\mathrm{d}w$$

# Legendre polynomials: bound

$$\delta(z) \;=\; \frac{1}{\sqrt{1 - 2\,x\,z + z^2}} \int_0^z \frac{\varepsilon'(w)}{\sqrt{1 - 2\,x\,w + w^2}}\, \mathrm{d}w \qquad\qquad \varepsilon(z) \ll \frac{3\,\mathbf{u}}{1 - z}$$

$$\delta(z) \; = \; \frac{1}{\sqrt{1-2\,x\,z+z^2}} \int_0^z \frac{\varepsilon'(w)}{\sqrt{1-2\,x\,w+w^2}}\,\mathrm{d}w \qquad\qquad \varepsilon(z) \ll \frac{3\,\mathbf{u}}{1-z}$$

$$\frac{3\,\mathbf{u}}{(1-z)^2}$$

$$\delta(z) \;=\; \frac{1}{\sqrt{1-2\,\mathsf{x}\,z+z^2}} \int_0^z \frac{\varepsilon'(w)}{\sqrt{1-2\,\mathsf{x}\,w+w^2}}\,\mathrm{d}w \qquad\qquad \varepsilon(z) \ll \frac{3\,\mathfrak{u}}{1-z}$$

$$\frac{1}{1-z} \qquad \frac{1}{1-z}\,\frac{3\,\mathfrak{u}}{(1-z)^2}$$

$$\delta(z) = \frac{1}{\sqrt{1 - 2x z + z^2}} \int_0^z \frac{\varepsilon'(w)}{\sqrt{1 - 2x w + w^2}} \, dw \qquad \varepsilon(z) \ll \frac{3u}{1-z}$$

$$\ll \frac{1}{1-z} \int \frac{1}{1-z} \frac{3u}{(1-z)^2}$$

$$\delta(z) = \frac{1}{\sqrt{1-2\,x\,z+z^2}} \int_0^z \frac{\varepsilon'(w)}{\sqrt{1-2\,x\,w+w^2}}\,\mathrm{d}w \qquad \varepsilon(z) \ll \frac{3\,\mathbf{u}}{1-z}$$

$$\ll \frac{1}{1-z} \int \frac{1}{1-z} \frac{3\,\mathbf{u}}{(1-z)^2}$$

$$= \frac{3}{2} \frac{1}{(1-z)^3}\,\mathbf{u}$$

---

**Proposition.** [Johansson & M.]

For all $x \in [-1, 1]$ and $n \in \mathbb{N}$, the error in the recursive fixed-point computation of Legendre polynomials satisfies

$$|\tilde{p}_n - P_n(x)| \leqslant \frac{3}{4}\,(n+1)\,(n+2)\,\mathbf{u}.$$

$$\begin{aligned}
\delta(z) &= \frac{1}{\sqrt{1-2\,x\,z+z^2}} \int_0^z \frac{\varepsilon'(w)}{\sqrt{1-2\,x\,w+w^2}}\, \mathrm{d}w & \varepsilon(z) \ll \frac{3\,\mathbf{u}}{1-z} \\
&\ll \frac{1}{1-z} \int \frac{1}{1-z} \frac{3\,\mathbf{u}}{(1-z)^2} \\
&= \frac{3}{2} \frac{1}{(1-z)^3}\, \mathbf{u}
\end{aligned}$$

---

**Proposition.** [Johansson & M.]

For all $x \in [-1,1]$ and $n \in \mathbb{N}$, the error in the recursive fixed-point computation of Legendre polynomials satisfies

$$|\tilde{p}_n - P_n(x)| \leqslant \frac{3}{4}\,(n+1)\,(n+2)\,\mathbf{u}.$$

---

We were lucky that the equation could be solved explicitly

# Partial sums of differentially finite series

$$L(u) = a_r u^{(r)} + \cdots + a_1 u' + a_0 u = 0, \quad a_i \in \mathbb{C}[z]$$

Given | the operator $L$,     compute an **enclosure** of $\displaystyle\sum_{n=0}^{N-1} u_n \zeta^n$.
the initial values $u_0, \ldots, u_{r-1}$
an evaluation point $\zeta$
a truncation order $N$

**Assumptions**

ordinary point    $a_r(0) \neq 0$

"obvious" geometric convergence    $|\zeta| < \min\{|\xi| : a_r(\xi) = 0\}$

**Strategy**

- Compute a recurrence on the $u_n$
- Compute and sum the $u_n \zeta^n$ iteratively          $\Rightarrow$    need to avoid interval blow-up

# D-finite series: error propagation

Exact rec.:
$$u_n = \frac{-1}{b_s(n)}\big[b_{s-1}(n)\,u_{n-1} + \cdots + b_1(n)\,u_{n-s+1} + b_0(n)\,u_{n-s}\big]$$

Approx. rec.:
$$\tilde{u}_n = \frac{-1}{b_s(n)}\big[b_{s-1}(n)\,\tilde{u}_{n-1} + \cdots + b_1(n)\,\tilde{u}_{n-s+1} + b_0(n)\,\tilde{u}_{n-s}\big] + \varepsilon_n$$

local error bound $|\varepsilon_n| \leqslant \hat{\varepsilon}_n$ computed on the fly

('running' error analysis)

## D-finite series: error propagation

Exact rec.: $\quad u_n = \dfrac{-1}{b_s(n)} \big[ b_{s-1}(n)\, u_{n-1} + \cdots + b_1(n)\, u_{n-s+1} + b_0(n)\, u_{n-s} \big]$

Approx. rec.: $\quad \tilde{u}_n = \dfrac{-1}{b_s(n)} \big[ b_{s-1}(n)\, \tilde{u}_{n-1} + \cdots + b_1(n)\, \tilde{u}_{n-s+1} + b_0(n)\, \tilde{u}_{n-s} \big] + \varepsilon_n$

local error bound $|\varepsilon_n| \leqslant \hat{\varepsilon}_n$ computed on the fly

('running' error analysis)

The global error $\delta_n = \tilde{u}_n - u_n$ satisfies

$$b_s(n)\, \delta_n + b_{s-1}(n)\, \delta_{n-1} + \cdots + b_0(n)\, \delta_{n-s} = b_s(n)\, \varepsilon_n$$

$$\Big\downarrow \ \textstyle\sum_n \square\, z^n$$

$$a_r(z)\, \delta^{(r)}(z) + \cdots + a_1(z)\, \delta'(z) + a_0(z)\, \delta(z) = Q(\theta) \cdot \varepsilon(z) \qquad\qquad \theta = z\dfrac{\mathrm{d}}{\mathrm{d}z}$$

$Q(\theta) = b_s(0)\, \theta\, (\theta - 1) \cdots (\theta - s + 1)$ (ordinary point)

Compute a bound on $\delta_n$ given one on $\varepsilon_n$?

$$a_r(z)\ \delta^{(r)}(z)\ +\cdots+a_1(z)\ \delta'(z)\ +a_0(z)\ \delta(z)\ =Q(\theta)\cdot\varepsilon(z)$$

**Lemma.** [$\sim$ Cauchy]

Let $a_0,...,a_r\in\mathbb{C}[z]$. Suppose $y\in\mathbb{C}[[z]]$ satisfies

$$a_r(z)\,y^{(r)}(z)+\cdots+a_0(z)\,y(z)=Q(\theta)\cdot\varepsilon(z).$$

Suppose $\varepsilon(z)\ll\hat\varepsilon(z)$.

One can compute a rational series $\hat a(z)\in\mathbb{R}_+[[z]]$ such that $y(z)$ is majorized by any solution of

$$\hat y'(z)=\hat a(z)\,\hat y(z)+\hat\varepsilon(z)$$

such that $|y_0|,...,|y_{r-1}|\leqslant\hat y_0$.

# D-finite series: error bound

$$a_r(z)\, \delta^{(r)}(z) + \cdots + a_1(z)\, \delta'(z) + a_0(z)\, \delta(z) = Q(\theta) \cdot \varepsilon(z)$$

---

**Lemma.** [$\sim$ Cauchy]

Let $a_0, ..., a_r \in \mathbb{C}[z]$. Suppose $y \in \mathbb{C}[[z]]$ satisfies

$$a_r(z)\, y^{(r)}(z) + \cdots + a_0(z)\, y(z) = Q(\theta) \cdot \varepsilon(z).$$

Suppose $\varepsilon(z) \ll \hat{\varepsilon}(z)$.

One can compute a rational series $\hat{a}(z) \in \mathbb{R}_+[[z]]$ such that $y(z)$ is majorized by any solution of

$$\hat{y}'(z) = \hat{a}(z)\, \hat{y}(z) + \hat{\varepsilon}(z)$$

such that $|y_0|, ..., |y_{r-1}| \leqslant \hat{y}_0$.

---

Solve:
$$\hat{\delta}(z) \;=\; \hat{h}(z)\left( \mathrm{cst} + \int_0^z \frac{\hat{\varepsilon}(w)}{\hat{h}(w)}\, dw \right), \qquad \hat{h}(z) = \exp \int_0^z \hat{a}(w)\, dw$$

# D-finite series: error bound

$$a_r(z)\, \delta^{(r)}(z) + \cdots + a_1(z)\, \delta'(z) + a_0(z)\, \delta(z) = Q(\theta) \cdot \varepsilon(z)$$

---

**Lemma.** [$\sim$ Cauchy]

Let $a_0, ..., a_r \in \mathbb{C}[z]$. Suppose $y \in \mathbb{C}[[z]]$ satisfies

$$a_r(z)\, y^{(r)}(z) + \cdots + a_0(z)\, y(z) = Q(\theta) \cdot \varepsilon(z).$$

Suppose $\varepsilon(z) \ll \hat{\varepsilon}(z)$.

One can compute a rational series $\hat{a}(z) \in \mathbb{R}_+[[z]]$ such that $y(z)$ is majorized by any solution of

$$\hat{y}'(z) = \hat{a}(z)\, \hat{y}(z) + \hat{\varepsilon}(z)$$

such that $|y_0|, ..., |y_{r-1}| \leqslant \hat{y}_0$.

---

Solve:
$$\hat{\delta}(z) \;=\; \hat{h}(z)\left( \mathrm{cst} + \int_0^z \frac{\hat{\varepsilon}(w)}{\hat{h}(w)}\, dw \right), \qquad \hat{h}(z) = \exp \int_0^z \hat{a}(w)\, dw$$

# D-finite series: error bound

$$a_r(z) \; \delta^{(r)}(z) \; + \cdots + a_1(z) \; \delta'(z) \; + a_0(z) \; \delta(z) \; = Q(\theta) \cdot \varepsilon(z)$$

**Lemma.** [$\sim$ Cauchy]

Let $a_0, ..., a_r \in \mathbb{C}[z]$. Suppose $y \in \mathbb{C}[[z]]$ satisfies

$$a_r(z) \, y^{(r)}(z) + \cdots + a_0(z) \, y(z) = Q(\theta) \cdot \varepsilon(z).$$

Suppose $\varepsilon(z) \ll \hat{\varepsilon}(z)$.

One can compute a rational series $\hat{a}(z) \in \mathbb{R}_+[[z]]$ such that $y(z)$ is majorized by any solution of

$$\hat{y}'(z) = \hat{a}(z) \, \hat{y}(z) + \hat{\varepsilon}(z)$$

such that $|y_0|, ..., |y_{r-1}| \leqslant \hat{y}_0$.

Solve:
$$\hat{\delta}(z) \;=\; \hat{h}(z) \left( \text{cst} + \int_0^z \frac{\hat{\varepsilon}(w)}{\hat{h}(w)} \, dw \right), \qquad \hat{h}(z) = \exp \int_0^z \hat{a}(w) \, dw$$

Choose $\hat{\varepsilon}_n = \bar{\varepsilon} \, \hat{h}_n$: $\qquad \qquad = \; \bar{\varepsilon} \, z \, \hat{h}(z)$

$$a_r(z)\ \delta^{(r)}(z)\ + \cdots + a_1(z)\ \delta'(z)\ + a_0(z)\ \delta(z)\ = Q(\theta) \cdot \varepsilon(z)$$

**Lemma.** [$\sim$ Cauchy]

Let $a_0, ..., a_r \in \mathbb{C}[z]$. Suppose $y \in \mathbb{C}[[z]]$ satisfies

$$a_r(z)\, y^{(r)}(z) + \cdots + a_0(z)\, y(z) = Q(\theta) \cdot \varepsilon(z).$$

Suppose $\varepsilon(z) \ll \hat{\varepsilon}(z)$.

One can compute a rational series $\hat{a}(z) \in \mathbb{R}_+[[z]]$ such that $y(z)$ is majorized by any solution of

$$\hat{y}'(z) = \hat{a}(z)\,\hat{y}(z) + \hat{\varepsilon}(z)$$

such that $|y_0|, ..., |y_{r-1}| \leqslant \hat{y}_0$.

Solve:
$$\hat{\delta}(z)\ =\ \hat{h}(z)\left( \text{cst} + \int_0^z \frac{\hat{\varepsilon}(w)}{\hat{h}(w)}\, dw \right), \qquad \hat{h}(z) = \exp \int_0^z \hat{a}(w)\, dw$$

Choose $\hat{\varepsilon}_n = \bar{\varepsilon}\, \hat{h}_n$:
$$=\ \bar{\varepsilon}\, z\, \hat{h}(z)$$

$$a_r(z) \; \delta^{(r)}(z) \; + \cdots + a_1(z) \; \delta'(z) \; + a_0(z) \; \delta(z) \; = Q(\theta) \cdot \varepsilon(z)$$

**Lemma.** [$\sim$ Cauchy]

Let $a_0, ..., a_r \in \mathbb{C}[z]$. Suppose $y \in \mathbb{C}[[z]]$ satisfies

$$a_r(z) \, y^{(r)}(z) + \cdots + a_0(z) \, y(z) = Q(\theta) \cdot \varepsilon(z).$$

Suppose $\varepsilon(z) \ll \hat{\varepsilon}(z)$.

One can compute a rational series $\hat{a}(z) \in \mathbb{R}_+[[z]]$ such that $y(z)$ is majorized by any solution of

$$\hat{y}'(z) = \hat{a}(z) \, \hat{y}(z) + \hat{\varepsilon}(z)$$

such that $|y_0|, ..., |y_{r-1}| \leqslant \hat{y}_0$.

Solve:
$$\hat{\delta}(z) \;=\; \hat{h}(z) \left( \text{cst} + \int_0^z \frac{\hat{\varepsilon}(w)}{\hat{h}(w)} \, dw \right), \qquad \hat{h}(z) = \exp \int_0^z \hat{a}(w) \, dw$$

Choose $\hat{\varepsilon}_n = \bar{\varepsilon} \, \hat{h}_n$:
$$= \; \bar{\varepsilon} \, z \, \hat{h}(z)$$

♻ Compute a bound on the truncation error at the same time

# Bernoulli Numbers

# Scaled Bernoulli numbers

$$B_n = 1, \frac{-1}{2}, \frac{1}{6}, 0, \frac{-1}{30}, 0, \frac{1}{42}, 0, \frac{-1}{30}, 0, \frac{5}{66}, 0, \frac{-691}{2730}, 0, \frac{7}{6}, 0, \frac{-3617}{510}, ...$$

$$|B_{2k}| \sim \frac{2\,(2\,k)!}{(2\,\pi)^{2k}}$$

$$b_k = \frac{B_{2k}}{(2\,k)!} \qquad\qquad b(z) = \sum_{k=0}^{\infty} b_k\, z^k = \frac{\sqrt{z}/2}{\tanh(\sqrt{z}/2)}$$

**Algorithm.** [Brent 1980, based on a suggestion of Reinsch]

$$b_k = \frac{1}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{b_j}{(2\,k+1-2\,j)!\,4^{k-j}}$$

be used with sufficient guard digits, or a more stable recurrence must be used. If we multiply both sides of (30) by $\sinh(x/2)/x$ and equate coefficients, we get the recurrence

$$C_k + \frac{C_{k-1}}{3!\,4} + \cdots + \frac{C_1}{(2k-1)!\,4^{k-1}} = \frac{2k}{(2k+1)!\,4^k} \qquad (36)$$

If (36) is used to evaluate $C_k$, using precision $n$ arithmetic, the error is only $O(k^2 2^{-n})$. Thus,

[Brent 1980]

$$b_k = \frac{1}{(2k)!\,4^k} - \sum_{j=0}^{k-1} \frac{b_j}{(2k+1-2j)!\,4^{k-j}}, \quad \tilde{b}_k = \text{computed values}$$

**Exercise 4.35** Prove (or give a plausibility argument for) the statements made in §4.7 that: (a) if a recurrence based on (4.59) is used to evaluate the scaled Bernoulli number $C_k$, using precision $n$ arithmetic, then the relative error is of order $4^k 2^{-n}$; and (b) if a recurrence based on (4.60) is used, then the relative error is $O(k^2 2^{-n})$.

[Brent & Zimmermann 2010]

**Conjecture.** [Brent, Zimmermann]

The computed values $\tilde{b}_k$ satisfy $\tilde{b}_k = b_k\,(1 + \eta_k)$ where $\eta_k = O(k \cdot \mathbf{u})$.

$\mathbf{u} = $ unit roundoff

**Remark.** To be understood as $\eta_k = O(k \cdot \mathbf{u})$ when $k = O(\mathbf{u}^{-1})$

or $|\eta_k| \leqslant C_k\, \mathbf{u}$ as $\mathbf{u} \to 0$ with $C_k = O(k)$ (resp. $O(k^2)$)

$$b_k = \frac{1}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{b_j}{(2\,k+1-2\,j)!\,4^{k-j}}, \quad \tilde{b}_k = \text{computed values}$$

**Local error analysis.**

$$\tilde{b}_k = \frac{1+s_k}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{\tilde{b}_j\,(1+t_{k,j})}{(2\,k+1-2\,j)!\,4^{k-j}}$$

$$|s_k| \leqslant \hat{\theta}_{2k}$$
$$|t_{k,j}| \leqslant \hat{\theta}_{3(k-j)+2}$$
$$\text{where } \hat{\theta}_n = (1+u)^n - 1$$

$$b_k = \frac{1}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{b_j}{(2\,k+1-2\,j)!\,4^{k-j}}, \quad \tilde{b}_k = \text{computed values}$$

**Local error analysis.**

$$\tilde{b}_k = \frac{1+\mathbf{s_k}}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{\tilde{b}_j\,(1+\mathbf{t_{k,j}})}{(2\,k+1-2\,j)!\,4^{k-j}}$$

$$|\mathbf{s_k}| \leqslant \hat{\theta}_{2k}$$
$$|\mathbf{t_{k,j}}| \leqslant \hat{\theta}_{3(k-j)+2}$$
$$\text{where } \hat{\theta}_n = (1+\mathbf{u})^n - 1$$

**Linearity.** $\delta_k := \tilde{b}_k - b_k = \text{global error}$

$$\delta_k = \frac{\mathbf{s_k}}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{\delta_j + (\,\mathbf{b_j}\, + \,\delta_j\,)\,\mathbf{t_{k,j}}}{(2\,k+1-2\,j)!\,4^{k-j}}$$

# Error analysis

$$b_k = \frac{1}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{b_j}{(2\,k+1-2\,j)!\,4^{k-j}}, \quad \tilde{b}_k = \text{computed values}$$

**Local error analysis.**

$$\tilde{b}_k = \frac{1 + s_k}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{\tilde{b}_j\,(1 + t_{k,j})}{(2\,k+1-2\,j)!\,4^{k-j}}$$

$$|s_k| \leqslant \hat{\theta}_{2k}$$
$$|t_{k,j}| \leqslant \hat{\theta}_{3(k-j)+2}$$
$$\text{where } \hat{\theta}_n = (1 + u)^n - 1$$

**Linearity.** $\delta_k := \tilde{b}_k - b_k = \text{global error}$

$$\delta_k = \frac{s_k}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{\delta_j + (\,b_j + \delta_j\,)\,t_{k,j}}{(2\,k+1-2\,j)!\,4^{k-j}}$$

**Inequation on the global error.**

$$\delta(z) \ll \check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp\delta(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp b(z)$$

$${}^\sharp f(z) = \sum_k |f_k|\,z^k$$

where
$$C(z) = \cosh\,(\sqrt{z}/2), \qquad S(z) = (\sqrt{z}/2)^{-1} \sinh\,(\sqrt{z}/2), \qquad \check{S}(z) = \frac{(\sqrt{z}/2)}{\sin\,(\sqrt{z}/2)},$$

$$\tilde{C}(z) = C(a^2\,z) - C(z), \qquad \tilde{S}(z) = S(a^4\,z) - S(z) - (a^2 - 1) \qquad \text{with } a = 1 + u$$

# Error analysis

$$b_k = \frac{1}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{b_j}{(2\,k+1-2\,j)!\,4^{k-j}}, \quad \tilde{b}_k = \text{computed values}$$

**Local error analysis.**

$$\tilde{b}_k = \frac{1 + s_k}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{\tilde{b}_j\,(1 + t_{k,j})}{(2\,k+1-2\,j)!\,4^{k-j}}$$

$$|s_k| \leqslant \hat{\theta}_{2k}$$
$$|t_{k,j}| \leqslant \hat{\theta}_{3(k-j)+2}$$
$$\text{where } \hat{\theta}_n = (1+u)^n - 1$$

**Linearity.** $\delta_k := \tilde{b}_k - b_k = \text{global error}$

$$\delta_k = \frac{s_k}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{\delta_j + (\,b_j\, + \,\delta_j\,)\,t_{k,j}}{(2\,k+1-2\,j)!\,4^{k-j}}$$

**Inequation on the global error.**

$$\delta(z) \ll \check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,^\sharp\delta(z) + \check{S}(z)\,\tilde{S}(z)\,^\sharp b(z) \qquad {}^\sharp f(z) = \sum_k |f_k|\,z^k$$

where
$$C(z) = \cosh(\sqrt{z}/2), \qquad S(z) = (\sqrt{z}/2)^{-1}\sinh(\sqrt{z}/2), \qquad \check{S}(z) = \frac{(\sqrt{z}/2)}{\sin(\sqrt{z}/2)},$$
$$\tilde{C}(z) = C(a^2\,z) - C(z), \qquad \tilde{S}(z) = S(a^4\,z) - S(z) - (a^2 - 1) \qquad \text{with } a = 1 + u$$

$$\delta(z) \ll \check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp b(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp \delta(z)$$

**'Explicit' majorant.** By the first lemma on majorizing equations

$$\delta(z) \ll \frac{\check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp b(z)}{1 - \check{S}(z)\,\tilde{S}(z)} \quad =: \hat{\delta}(z)$$

$$\delta(z) \ll \check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp b(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp\delta(z)$$

**'Explicit' majorant.** By the first lemma on majorizing equations

$$\delta(z) \ll \frac{\check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,{}^\sharp b(z)}{1 - \check{S}(z)\,\tilde{S}(z)} \quad =: \hat{\delta}(z)$$

**Asymptotic behavior.**                    💡 Series notation → computer algebra

$$\hat{\delta}(z) = \left( \frac{2\,(1 - \cosh w)\cos(w)}{w^{-2}\sin(w)^2} + \frac{4\,(\cosh w - 1) + w\sinh w}{w^{-1}\sin w} \right)\mathbf{u} + O(\mathbf{u}^2)$$

$$w = \sqrt{z}/2$$

Unique dominant pole at $z = 4\,\pi^2$,
multiplicity (w.r.t. $z$) $= 2$
$\Rightarrow \hat{\delta}_k = O(k\,(2\,\pi)^{-2k}) \cdot \mathbf{u} + O(\mathbf{u}^2)$
$\Rightarrow \eta_k = \text{"}O(k \cdot \mathbf{u})\text{"}$

$$\delta(z) = \hat{\delta}(z) \ll \frac{\check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,^\sharp b(z)}{1 - \check{S}(z)\,\tilde{S}(z)}$$

**Controlling the dominant pole.**

Suppose $\mathbf{u} \leqslant 2^{-16}$.

Then $\hat{\delta}(z)$ has a pole at $\gamma = \left(\dfrac{2\,\pi}{1 + \varphi(\mathbf{u})}\right)^2$ where $0 \leqslant \varphi(\mathbf{u}) \leqslant 2\,(\cosh \pi - 1)\,\mathbf{u}$.

This is the only pole with $|z| < 153.7 \approx (3.9\,\pi)^2$.

(A little analysis + comparison with the limiting case using Rouché's theorem.)

$$\delta(z) = \hat{\delta}(z) \ll \frac{\check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\, {}^{\sharp}b(z)}{1 - \check{S}(z)\,\tilde{S}(z)}$$

**Controlling the dominant pole.**

Suppose $u \leqslant 2^{-16}$.

Then $\hat{\delta}(z)$ has a pole at $\gamma = \left( \dfrac{2\pi}{1 + \varphi(u)} \right)^2$ where $0 \leqslant \varphi(u) \leqslant 2\,(\cosh \pi - 1)\,u$.

This is the only pole with $|z| < 153.7 \approx (3.9\,\pi)^2$.

(A little analysis + comparison with the limiting case using Rouché's theorem.)

**Symbolic-numeric estimate.**

$$\hat{\delta}(z) = \frac{2 \; \text{explicit } R(u)}{1 - z/\gamma} - \frac{2}{1 - z/(2\pi)^2} + \text{analytic for } |z| < 153.7$$

$$\hat{\delta}(z) \ll \frac{2\,|\,R(u)\,-1|}{1 - z/\gamma} + \frac{\text{explicit and } O(u)}{(1 - z/\gamma)^2} + \frac{\sup_{|z=\lambda\gamma|} \text{ analytic}}{1 - z/(\lambda\,\gamma)}$$

Cauchy's formula + interval arithmetic.

$$\delta(z) = \hat{\delta}(z) \ll \frac{\check{S}(z)\,\tilde{C}(z) + \check{S}(z)\,\tilde{S}(z)\,{}^{\sharp}b(z)}{1 - \check{S}(z)\,\tilde{S}(z)}$$

**Controlling the dominant pole.**

Suppose $\mathbf{u} \leqslant 2^{-16}$.

Then $\hat{\delta}(z)$ has a pole at $\gamma = \left(\dfrac{2\,\pi}{1 + \varphi(\mathbf{u})}\right)^2$ where $0 \leqslant \varphi(\mathbf{u}) \leqslant 2\,(\cosh \pi - 1)\,\mathbf{u}$.

This is the only pole with $|z| < 153.7 \approx (3.9\,\pi)^2$.

(A little analysis + comparison with the limiting case using Rouché's theorem.)

**Symbolic-numeric estimate.**

$$\hat{\delta}(z) = \frac{2\ \text{explicit } \mathbf{R}(\mathbf{u})}{1 - z/\gamma} - \frac{2}{1 - z/(2\,\pi)^2} + \ \text{analytic for } |z| < 153.7$$

$$\hat{\delta}(z) \ll \frac{2\,|\ \mathbf{R}(\mathbf{u})\ -1|}{1 - z/\gamma} + \frac{\text{explicit and } O(\mathbf{u})}{(1 - z/\gamma)^2} + \frac{\sup_{|z=\lambda\gamma|}\ \text{analytic}}{1 - z/(\lambda\,\gamma)}$$

Cauchy's formula + interval arithmetic.

$$b(z) = \frac{\sqrt{z}/2}{\tanh(\sqrt{z}/2)} \qquad\qquad b_k = \frac{1}{(2\,k)!\,4^k} - \sum_{j=0}^{k-1} \frac{b_j}{(2\,k+1-2\,j)!\,4^{k-j}}$$

$$\tilde{b}_k = b_k\,(1+\eta_k)$$

**Theorem.** The total relative error satisfies
$$|\eta_k| \leqslant (1+21.2\,\mathbf{u})^k\,(1.1\,k+446)\,\mathbf{u}$$

**Corollary.** Assuming $\mathbf{u} < 2^{-16}$ and $43\,k\,\mathbf{u} \leqslant 1$, one has $|\eta_k| \leqslant (3\,k+1213)\,\mathbf{u}$.

# Conclusion

Error analyses of linear recurrences can (should!) use generating series

- Local errors $\rightarrow$ global errors via exact expressions or equations
- Cauchy majorants
- Analytic methods

- Legendre polynomials
- General D-finite functions
- Bernoulli numbers

- Other algorithms for D-finite functions, e.g., $O(n\,M(d)/d)$
- Tighter bounds in practice
- Backward recurrence schemes
- Orthogonal polynomials, numerical integration schemes, …

# Conclusion

Error analyses of linear recurrences can (should!) use generating series

- Local errors $\rightarrow$ global errors via exact expressions or equations
- Cauchy majorants
- Analytic methods

- Legendre polynomials
- General D-finite functions
- Bernoulli numbers

- Other algorithms for D-finite functions, e.g., $O(n\,M(d)/d)$
- Tighter bounds in practice
- Backward recurrence schemes
- Orthogonal polynomials, numerical integration schemes, …

**Thank you!**

# Credits