A few observations about summation algorithms

Ilse C.F. Ipsen

Joint work with Eric Hallman also thanks to Johnathan Rhyne

North Carolina State University Raleigh, NC, USA

Research supported by NSF DMS

This talk

Forward error when summing *n* real numbers x_1, \ldots, x_n

$$s_n = x_1 + \cdots + x_n$$

in floating point arithmetic

Overview

- The traditional method
- More accurate methods (without higher precision or better hardware)
 - Inspired by computer architecture and formal methods: Shifted summation
 - Kahan's 1965 method: Compensated summation
- Summary

The traditional method

Sequential (recursive) summation

[Higham: Accuracy and Stability of Numerical Algorithms, Chapter 4]

• Exact arithmetic

$$s_1 = x_1, \qquad s_k = s_{k-1} + x_k, \quad 2 \le k \le n$$

• Floating point arithmetic¹

$$\hat{s}_1 = x_1, \qquad \hat{s}_k = (\hat{s}_{k-1} + x_k)(1 + \delta_k), \quad 2 \le k \le n$$

where $|\delta_k| \leq u$ are n-1 roundoffs

• First order forward error

$$\frac{\hat{s}_n - s_n}{s_n} \bigg| \le n u \underbrace{\frac{\sum_{k=1}^n |x_k|}{|s_n|}}_{\ge 1} + \mathcal{O}(u^2)$$

¹Assume: x_1, \ldots, x_n are floating point numbers

Tighter forward error bounds for sequential summation

• Exact expression [Hallman 2020]

$$\hat{s}_n - s_n = \sum_{k=2}^n s_k \, \delta_k \, \prod_{j=k}^n (1+\delta_j)$$

Probabilistic bounds

[Higham, Mary 2019, 2020], [Rhyne 2020], [Hallman 2020]

Treat roundoffs δ_k as zero-mean bounded random variables (independent, or mean-independent) Error $\approx O(\sqrt{nu})$ with high probability

• Next: More accurate algorithms

Shifted summation

Motivation for shifted summation

• Computer architecture, formal methods for program verification:

[Solovyev et al 2015], [Dahlquist, Salvia, Constatinides 2019], [Lohar, Prokop, Darulova 2019], [Constantinides, Dahlquist, Rakamaric, Salvia 2021]

Assume x_1, \ldots, x_n drawn from a distribution Compute statistics for accumulated errors Determine probability of errors in certain interval

• Probabilistic bounds for random data

[Higham, Mary 2020], [Hallman 2020] Idea: Sequential summation accurate if x_k tightly clustered

Sequential summation of random data [Hallman 2020]

Assume:

- Roundoffs δ_k are independent zero-mean random variables, and $|\delta_k| \leq u$
- Summands x_k are random variables with mean μ and 'variance' $\max_{1 \le k \le n} |x_k \mu| \le \sigma$

Then for any 0 $<\delta<1$ with probability at least $1-\delta$

$$|\hat{s}_n - s_n| \leq (1 + \gamma)(\lambda n^{3/2} |\mu| + \lambda^2 n \sigma) u$$

where $\lambda = \sqrt{2 \ln (6/\delta)}$ and $\gamma = \exp \left(\lambda \left(\sqrt{n}u + nu^2\right)/(1-u)\right) - 1$

Sequential summation accurate if x_k tightly clustered around zero

Application to non-random data [Higham, Mary 2020]

Example: Assume $x_k = 10^4 + y_k$ where $|y_k| \le 1$

- Sequential summation $|\hat{s}_n - s_n| \le nu \sum_{k=1}^n |x_k| + \mathcal{O}(u^2) = n^2 u (10^4 + 1) + \mathcal{O}(u^2)$
- 'Center' the data to reduce their 'mean'

Center: $y_k = x_k - 10^4$, $1 \le k \le n$ Sum centered data $t_1 = y_1$, $t_k = t_{k-1} + y_k$, $1 \le k \le n$ Uncenter: $s = t_n + n \cdot 10^4$

$$|\hat{t}_n - t_n| \leq nu \sum_{k=1}^n |y_k| + \mathcal{O}(u^2) \leq n^2 u + \mathcal{O}(u^2)$$

Centered sum captures 'tail' bits of $x_k \Longrightarrow$ smaller error

Shifted sequential summation [Higham, Mary 2020]

Input: summands
$$x_1, \ldots, x_n$$
, shift c
 $t_0 = 0$
for $k = 1 : n$ do
 $\hat{y}_k = (x_k - c)(1 + \epsilon_k)$ { ϵ_k are centering roundoffs}}
 $\hat{t}_k = (\hat{t}_{k-1} + \hat{y}_k)(1 + \delta_k)$ { δ_k are summation roundoffs}
end for
 $\hat{y}_{n+1} = cn(1 + \epsilon_{n+1})$
 $\hat{s}_n = (\hat{t}_n + \hat{y}_{n+1})(1 + \delta_{n+1})$ {uncentering}

Error [Hallman, II 2021]

$$\hat{s}_n - s_n = \underbrace{\sum_{k=1}^{n+1} t_k \delta_k}_{\text{Summation}} \prod_{\ell=k+1}^{n+1} (1+\delta_\ell) + \underbrace{\sum_{k=1}^{n+1} y_k \epsilon_k}_{\text{Centering}} \prod_{\ell=k}^{n+1} (1+\delta_\ell)$$

Error bounds for shifted summation [Hallman, II 2021]

Error depends only on shifted quantities

• Deterministic

$$|\hat{s}_n - s_n| \le u(1+u)^n \left(\sum_{k=2}^n \underbrace{|s_k - kc|}_{|t_k|} + \sum_{k=1}^n \underbrace{|x_k - c|}_{|y_k|} + |s| + |nc|
ight)$$

Probabilistic

Let δ_j and ϵ_j be independent zero-mean random variables Then for any $0 < \delta < 1$ with probability at least $1 - \delta$,

$$|\hat{s}_n-s_n|\leq \max_{1\leq k\leq n+1}\left(|s_k-kc|+|x_k-c|
ight)eta\sqrt{2\ln(2/\delta)}$$

where $\beta = \sqrt{\frac{u}{2}\gamma_{2n}} \le \sqrt{\frac{(n+2)u^2}{1-2(n+2)u}} \approx \sqrt{n+2} u$ (provided $n \ll 1/u$)

Numerical experiments: Sequential summation with shifting vs no shifting

- Number of summands $n = 10^6$
- Shift $c = (\max_k x_k + \min_k x_k)/2$
- Working precision: Julia Float64 $u = 2^{-53} \approx 1.11 \cdot 10^{-16}$
- 'Exact' computation: Julia Float256
- Plots: relative errors $|\hat{s}_n s_n|/|s_n|$ versus n
- Probabilistic 'bound'

$$u\sqrt{n+2} \max_{1 \le k \le n+1} \frac{|s_k - kc| + |x_k - c|}{|s_n|} \underbrace{\sqrt{2\ln(2/\delta)}}_{3.26} \qquad \delta = 10^{-2}$$

Two different types of summands x_k , $1 \le k \le n$

Figure 1: $x_k = 10^4 + y_k$ where y_k is uniform[0,1] Figure 2: x_k is normal(0,1) Summands tightly clustered at 10^4 $x_k = 10^4 + y_k$ where y_k is uniform[0,1]



Shifting increases accuracy compared to plain summation Bound accurate within factor of 100

Summands tightly clustered at zero x_k is normal(0,1)



Shifting hurts accuracy $(-.2 \le c \le 1.2)$ over plain summation Bound accurate within factor of 10-100

Summary: Shifted sequential summation

- Error bounds depend only on shifted quantities and hold to all orders of *u*
- Probabilistic bound: Error $\approx \mathcal{O}(\sqrt{n}u)$
- Shifting improves accuracy if shift *c* decreases magnitude of partial sums

$$\max_{k}\left(|s_{k}-k c|+|x_{k}-c|\right) \ll \max_{k}\left(|s_{k}|+|x_{k}|\right)$$

See advice in [Higham book, Section 4.2]

• Extension to general summation algorithms [Hallman, II 2021]

Roundoffs treated as (mean-) independent random variables Probabilistic bounds depend on $\sqrt{\text{height of computational tree}}$ Bounds valid to all orders of u

• Warning: Shifting can worsen accuracy for naturally centered data, such as *normal*(0,1)

Kahan's 1965 method: Compensated sequential summation

Compensated sequential summation

[Goldberg 1991], [Higham book], [Kahan 1973]

Input: summands x_1, \ldots, x_n $s_1 = x_1, c_1 = 0$ for k = 2 : n do $y_k = x_k - c_{k-1}$ $s_k = s_{k-1} + y_k$ $c_k = (s_k - s_{k-1}) - y_k$ end for return s_n Summands tightly clustered at 10^4 $x_k = 10^4 + y_k$ where y_k is uniform[0,1]



Compensated summation as or more accurate than shifted summation with $c = (\min_k x_k + \max_k x_k)/2$

Summands clustered at zero x_k is normal(0,1)



Compensated summation more accurate than shifted summation with $c = (\min_k x_k + \max_k x_k)/2$

Roundoff in sequential compensated summation

[Goldberg 1991], [Kahan 1973]

$$\begin{split} \hat{s}_{1} &= s_{1} = x_{1}, \ \hat{c}_{1} = 0\\ \text{for } k &= 2:n \text{ do}\\ \hat{y}_{k} &= (x_{k} - \hat{c}_{k-1})(1 + \eta_{k})\\ \hat{s}_{k} &= (\hat{s}_{k-1} + \hat{y}_{k})(1 + \sigma_{k})\\ \hat{c}_{k} &= ((\hat{s}_{k} - \hat{s}_{k-1})(1 + \delta_{k}) - \hat{y}_{k})(1 + \beta_{k})\\ \text{end for} \end{split}$$

Forward error

$$|\hat{s}_n - s_n| \leq \left(2u + \mathcal{O}(nu^2)\right) \sum_{k=1}^n |x_k|$$

Recursion for error and correction [Hallman, II 2021]

$$e_{k} = \hat{s}_{k} - s_{k}$$

$$\hat{c}_{k} = ((\hat{s}_{k} - \hat{s}_{k-1})(1 + \delta_{k}) - \underbrace{(x_{k} - \hat{c}_{k-1})(1 + \eta_{k})}_{\hat{y}_{k}})(1 + \beta_{k})$$

satisfy the recurrence

$$\begin{bmatrix} e_k \\ \hat{c}_k \end{bmatrix} = P_k \begin{bmatrix} e_{k-1} \\ \hat{c}_{k-1} \end{bmatrix} + P_k \begin{bmatrix} s_{k-1} \\ -x_k \end{bmatrix} + \begin{bmatrix} -s_k \\ 0 \end{bmatrix} \qquad 2 \le k \le n$$

where

$$\mathsf{P}_k \equiv \begin{bmatrix} 1+\sigma_k & -(1+\eta_k)(1+\sigma_k) \\ \sigma_k(1+\delta_k)(1+\beta_k) & (1+\eta_k)(1-(1+\sigma_k)(1+\delta_k))(1+\beta_k) \end{bmatrix}$$

Explicit expressions for error and correction [Hallman, II 2021]

$$\begin{bmatrix} e_n \\ \hat{c}_n \end{bmatrix} = (P_n \cdots P_2) \begin{bmatrix} x_1 \\ 0 \end{bmatrix} + \sum_{j=2}^n (P_n \cdots P_j) \begin{bmatrix} 0 \\ -x_j \end{bmatrix} + \begin{bmatrix} -s_n \\ 0 \end{bmatrix}$$

A second explicit expression

$$e_n = s_n \sigma_n + \sum_{j=4}^n x_j \eta_j \prod_{k=j}^n (1 + \sigma_k) \\ + \sum_{j=2}^{n-1} (s_j \sigma_j - \hat{c}_j (1 + \eta_{j+1})) \prod_{k=j+1}^n (1 + \sigma_k)$$

A third explicit expression [Hallman, II 2021]

$$e_n = \hat{s}_n - s_n = \underline{s_n \sigma_n} + \underline{X_n} + \underbrace{\underline{X_n \sigma_n} + \underline{E_{n-1}(1 + \sigma_n)}}_{\mathcal{O}(u^2)}$$

where

$$X_{n} \equiv x_{n}\eta_{n}(1+\beta_{n}) + \sum_{j=2}^{n-1} x_{j}(\eta_{j}-\delta_{j}) \prod_{\ell=j}^{n} (1+\beta_{\ell})(1+\eta_{\ell+1})$$

$$\mathcal{O}(u)$$

$$E_{n-1} \equiv e_{n-1}\Theta_{n-1} + \sum_{j=2}^{n-1} e_{j}(\Theta_{j}+\delta_{j+1}) \prod_{\ell=j+1}^{n} (1+\beta_{\ell})(1+\eta_{\ell+1})$$

$$\Theta_{k} = 1 - (1+\delta_{k})(1+\beta_{k})(1+\eta_{k+1}) \quad 2 \le k \le n-1$$

First order error [Hallman, II 2021]

The previous expressions imply

$$\hat{s}_n - s_n = s_n \sigma_n + x_n \eta_n + \sum_{k=2}^{n-1} x_k (\eta_k - \delta_k) + \mathcal{O}(u^2)$$

Error dominated by correction and final summation roundoffs

$$egin{aligned} \hat{y}_k &= (x_k - \hat{c}_{k-1})(1 + \eta_k), \ \hat{c}_k &= ((\hat{s}_k - \hat{s}_{k-1})(1 + \delta_k) - \hat{y}_k) \, (1 + eta_k) \ \hat{s}_n &= (\hat{s}_{n-1} + \hat{y}_n)(1 + \sigma_n) \end{aligned}$$

• First order bound $|\hat{s}_n - s_n| \le 3u \sum_{k=1}^n |x_k| + O(u^2)$ inconsistent with previous bounds $|\hat{s}_n - s_n| \le 2u \sum_{k=1}^n |x_k| + O(u^2)$

Second order error [Hallman, II 2021]

Let
$$\mu_k \equiv \eta_k - \delta_k$$
, $1 \le k \le n - 1$, $\mu_n = \eta_n$

-

$$\begin{split} \hat{s}_{n} - s_{n} &= s_{n}\sigma_{n} + \sum_{k=2}^{n} x_{k}\mu_{k} \left(1 + \sigma_{n}\right) \\ &- \sum_{k=2}^{n-1} \left(s_{k}\sigma_{k}(\mu_{k+1} + \beta_{k} + \delta_{k}) + x_{k}\delta_{k}(\mu_{k+1} + \beta_{k} + \eta_{k})\right) + \mathcal{O}(u^{3}) \end{split}$$

Deterministic bound

$$|\hat{s}_n - s_n| \leq (3u + 4n u^2) \sum_{k=1}^n |x_k| + \mathcal{O}(u^3)$$

Probabilistic bounds [Hallman, II 2021]

Vector of summands $\mathbf{x} = \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix}^T$ Independent zero-mean random variables δ_j and ϵ_j Then for any $0 < \delta < 1$ with probability at least δ

Second order bound

$$\begin{aligned} |\hat{s}_n - s_n| &\leq u\left((2+6u) \|\mathbf{x}\|_2 + \sqrt{|s_n|^2 + 16u^2 \sum_{k=1}^{n-1} |s_k|^2}\right) \sqrt{2\ln(2/\delta)} + \mathcal{O}(u^3) \\ &\leq u\left((2+6u) \|\mathbf{x}\|_2 + \sqrt{1+16(n-2)u^2} \|\mathbf{x}\|_1\right) \sqrt{2\ln(2/\delta)} + \mathcal{O}(u^3) \end{aligned}$$

First order bound

$$|\hat{s}_n - s_n| \le u(2||\mathsf{x}||_2 + |s_n|)\sqrt{2\ln(2/\delta)} + \mathcal{O}(u^2)$$

Numerical experiments: Compensated summation in half precision

- Working precision: Julia Float16 $u = 2^{-11} \approx 4.88 \cdot 10^{-4}$
- 'Exact' computation: Julia Float64
- Plots: relative errors $|\hat{s}_n s_n|/|s_n|$ versus n
- First order probabilistic bound

$$u \frac{2||\mathbf{x}||_2 + |s_n|}{|s_n|} \underbrace{\sqrt{2\ln(2/\delta)}}_{3.26} \quad \text{where} \quad \delta = 10^{-2}$$

Two different types of summands x_k , $1 \le k \le n$

Figure 1: x_k is uniform[0,1], $n = 6 \cdot 10^4$ Figure 2: x_k is normal(0,1), $n = 10^6$

Summands have same sign: x_k is uniform[0,1]



Compensated summation accurate to machine precision Ordinary summation not accurate (well-conditioned problem) Bound accurate within factor of 10

Summands have different signs: x_k is normal(0,1)



Bound accurate within factor of 10

Summary

Forward error in summation of *n* real numbers without recourse to higher precision or better arithmetic

• Shifted sequential summation Explicit expression and probabilistic bounds valid to all orders Extension to general summation algorithms Centering more accurate if it decreases magnitude of partial sums But centering can hurt accuracy Mixed precision versions do not appear to be effective

• Compensated sequential summation Three explicit expressions for error First and second order deterministic and probabilistic error bounds First order bound differs by *u* from existing bounds First order probabilistic bound accurate within factor of 10 Accurate but more expensive Compromise: FABsum [Blanchard, Higham, Mary 2020]

E. Hallman and I. C. F. Ipsen: Deterministic and Probabilistic Error Bounds for Floating Point Summation Algorithms, arXiv:2107.01604