

浮動小数点演算による内積の丸め誤差解析

樋口 裕幸* 尾崎 克久†

* 芝浦工業大学大学院理工学研究科システム理工学専攻

† 芝浦工業大学システム理工学部数理科学科

概要. 数値計算に幅広く採用されている浮動小数点演算は高速である一方で誤差の問題を抱えている. そのために, 計算結果と真値の距離の上限, すなわち絶対誤差の上限を求める丸め誤差解析の研究が進んでいる. 本稿では, 内積に対する丸め誤差解析の最近の結果をその特徴を含めて紹介しつつ, 著者らにより改善できた誤差評価式とその特徴を紹介する.

Notes on Error Estimates for Dot Product by Floating-Point Arithmetic

Hiroyuki Higuchi* Katsuhisa Ozaki†

*Division of Systems Engineering and Science,
Graduate School of Engineering and Science,
Shibaura Institute of Technology

†Department of Mathematical Science,
College of Systems Engineering and Science,
Shibaura Institute of Technology

Abstract. This paper is concerned with a rounding error analysis of dot product by floating-point arithmetic. Floating-point arithmetic is fast performed on recent architectures. However, since the number of total bits in floating-point format is finite, rounding errors may occur in arithmetic operations. Recently, new analysis of rounding errors for dot product was developed. We review the results of the error analysis of dot product and revise them.

1. はじめに

本稿では, 浮動小数点演算による内積の丸め誤差解析について最近の成果を紹介するとともに, 先行研究の改善を行う. 浮動小数点数とその演算は IEEE 754 規格 [3] によって定められた 2 進数フォーマットを使用する. オーバーフロー・アンダーフローという用語の定義は IEEE 754 規格 [3] に従うものとする. \mathbb{F} を固定された精度における浮動小数点数の正規化数・非正規化数・零の和集合とする. $\text{fl}(\cdot)$ や $\text{float}(\cdot)$ という表記は文献 [4] と同じ使い方をする. すなわち, $\text{fl}(\cdot)$ は括弧内におけるすべての 2 項演算を最近点への丸め

モード（偶数丸め：roundTiesToEven を採用）による浮動小数点演算で評価した結果を意味する．2 数の積に対して「 \cdot 」記号を基本的に用いない．ただし，浮動小数点演算中に誤差が発生する可能性がある場合や，「 \cdot 」記号を入れたほうが読みやすい部分には使用する． $a, b, c \in \mathbb{F}$ に対して， $\text{fl}(\text{fl}(a + b) + c)$ は $\text{fl}((a + b) + c)$ のように表記を簡略化する． $\text{float}(\cdot)$ は括弧内を浮動小数点演算で計算した結果を意味し，括弧内の計算順序は任意とする．例えば $p \in \mathbb{F}^4$ に対して $\text{float}(\sum_{i=1}^4 p_i)$ は

$$\text{fl}(((p_1 + p_2) + p_3) + p_4), \quad \text{fl}((p_1 + p_2) + (p_3 + p_4)), \quad \text{fl}((p_1 + p_3) + (p_2 + p_4))$$

など，どれでも良い．本稿全体を通して， $\text{fl}(\cdot)$ ， $\text{float}(\cdot)$ の評価中にオーバーフローは発生しないと仮定する． \mathbf{u} を相対丸めとする．IEEE 754 規格により $\mathbf{u} \in \{2^{-24}, 2^{-53}, 2^{-113}\}$ となり，倍精度浮動小数点数（binary64）に対しては $\mathbf{u} = 2^{-53}$ である． $a \in \mathbb{R}$ に対する unit in the first place（先頭桁の単位）を表す実関数 $\text{ufp}(a)$ を下記のように導入する．

$$(1.1) \quad \text{ufp}(a) := \begin{cases} 0, & \text{if } a = 0, \\ 2^{\lfloor \log_2 |a| \rfloor}, & \text{otherwise} \end{cases}$$

$a \neq 0$ の場合， $\text{ufp}(a)$ は $|a|$ 以下の最大の 2 のべき乗数を表す． ufp に関して，

$$(1.2) \quad \text{ufp}(a) \leq |a| < 2\text{ufp}(a)$$

が成立する．この不等式は $\text{ufp}(a) = |a|$ のときには自明であり，そうでない場合は

$$\text{ufp}(a) = 2^{\lfloor \log_2 |a| \rfloor} < 2^{\log_2 |a|} = |a| < 2^{\lfloor \log_2 |a| \rfloor + 1} = 2\text{ufp}(a)$$

から示される．

これから，成分がすべて浮動小数点数であるベクトルの総和を例に挙げながら，本論文で着目する誤差評価式の特徴を整理したい．浮動小数点演算によるベクトル $p \in \mathbb{F}^n$ の総和については，文献 [2] の第 3 章や第 4 章でよく紹介されており，

$$(1.3) \quad |\text{float}(\sum_{i=1}^n p_i) - \sum_{i=1}^n p_i| \leq ((1 + \mathbf{u})^{n-1} - 1) \sum_{i=1}^n |p_i|$$

や， $(n-1)\mathbf{u} < 1$ のとき

$$(1.4) \quad |\text{float}(\sum_{i=1}^n p_i) - \sum_{i=1}^n p_i| \leq \gamma_{n-1} \sum_{i=1}^n |p_i|, \quad \gamma_n = \frac{n\mathbf{u}}{1 - n\mathbf{u}}$$

が成立する．ここで，丸め誤差解析に関して得られる誤差評価式に対して着目する点を

ポイント (a) 任意の入力サイズ $n \in \mathbb{N}$ に対して成立するか？

ポイント (b) 誤差上限に関する係数に対する計算回数・計算コストは？

ポイント (c) 浮動小数点演算でそのまま誤差の上限が評価可能か？

ポイント (d) 誤差の上限は最適か？または過大評価を抑えているか？

ポイント (e) アンダーフローに対応できているか？

とする．ポイント (b) の誤差に関する係数とは，総和の場合は $\sum_{i=1}^n |p_i|$ に関する項に乘ずる \mathbf{u} と n を含む値を意味する．すなわち，誤差評価式 (1.3) では $(1 + \mathbf{u})^{n-1} - 1$ が，また誤差評価式 (1.4) では γ_{n-1} が係数である．例えば誤差評価式 (1.3) はポイント (a) について達成されており，誤差評価式 (1.4) は n の大きさに制限がある．ポイント (b) を考えれば， $n-1$ 乗の計算を含む誤差評価式 (1.3) より誤差評価式 (1.4) のほうが大きな n に対して計算回数が少ない．ポイント (c) については，誤差の上限が浮動小数点演算の結果ではないため，誤差評価式 (1.3) と (1.4) のどちらも満たされていない．ポイント (d) については，どちらも係数は最適ではないものの，係数の大きさは誤差評価式 (1.3) のほうが誤差評価式 (1.4) よりも良い．浮動小数点数の総和については，非正規化数の範囲内における和の演算は正確に行われるため，ここではポイント (e) は議論しなくて良い．これらの特徴をすべて満たす誤差評価式を導出することは非常に難しく，目的によって優先する特徴を設定することになる．

近年，ベクトルの総和の丸め誤差解析に対して進展があった．Rump は浮動小数点数の総和について，

$$(1.5) \quad \left| \text{float}\left(\sum_{i=1}^n p_i\right) - \sum_{i=1}^n p_i \right| \leq (n-1)\mathbf{u} \cdot \text{ufp}\left(\text{float}\left(\sum_{i=1}^n |p_i|\right)\right)$$

を示した [6,7]^{*1}．ここでは，誤差評価式 (1.5) の左辺と右辺にある $\text{float}(\cdot)$ 内の計算順序はともに同じとする．ベクトルを $p = (1, \mathbf{u}, \mathbf{u}, \dots, \mathbf{u})^T \in \mathbb{F}^n$ とし，すべての n において前方から順に足しこめば $\text{fl}(((p_1 + p_2) + p_3) + \dots + p_n) = 1$ となる．また誤差が $(n-1)\mathbf{u}$ となることから，零ベクトル以外の入力に対しても不等式における等号の例がすべての n に対して見つかる．よって，係数 $(n-1)\mathbf{u}$ はこれ以上小さくできないため，ポイント (d) において最適な評価と言える．また， $\alpha \in \mathbb{F}$ に対する $\text{ufp}(\alpha)$ は3回の浮動小数点演算で計算可能であることが文献 [6] に示されているため， ufp の評価は簡単である．以上より，誤差評価式 (1.5) はポイント (a), (b), (d) が達成された優れた誤差評価式である．また，誤差評価式 (1.5) の右辺を浮動小数点演算ですべて評価可能にした

$$(1.6) \quad n\mathbf{u} \leq 1 \implies \left| \text{fl}\left(\sum_{i=1}^n p_i\right) - \sum_{i=1}^n p_i \right| \leq \text{fl}((n-1) \cdot (\mathbf{u} \cdot \text{ufp}(\text{fl}\left(\sum_{i=1}^n |p_i|\right))))$$

も [6] にて示されており，誤差評価式 (1.6) に対してはポイント (b), (c), (d) が達成されている．ここで誤差評価式 (1.6) における $\text{fl}\left(\sum_{i=1}^n p_i\right)$ は，前方から順に浮動小数点演算で足し

^{*1} 計算順序を固定した誤差評価式が [6] にてまず示され，その後任意の計算順序について成立することが [7] にて示された．

こんだ総和の結果を表す．また，誤差評価式 (1.4) の改良として，

$$(1.7) \quad |\text{float}(\sum_{i=1}^n p_i) - \sum_{i=1}^n p_i| \leq (n-1)\mathbf{u} \sum_{i=1}^n |p_i|$$

が，Jeannerod と Rump によって示された [4]．誤差評価式 (1.7) は n に制限がなく，かつ誤差上限の中に \mathbf{u}^2 以降に相当する項がないことからポイント (a)，(b) を満たすきれいな評価式であり，誤差評価式 (1.3) と (1.4) に対してもポイント (d) の点で優れている．以上が総和に関する丸め誤差解析の紹介であった．

本稿の構成は以下の通りである．1 章において総和の丸め誤差解析について紹介し，本論文で着目する誤差評価式に関する特徴を紹介した．2 章では提案手法に用いられる丸め誤差解析のモデルと 2 分木を用いた帰納法の概要について紹介する．3 章では内積の丸め誤差解析に関する先行研究を紹介し，著者らによって改良できた誤差評価式について紹介を行う．また近年の CPU では標準的に使用できることも多い Fused Multiply-Add を用いた誤差評価式について紹介する．

2. 誤差解析のモデル

本章では，誤差評価式の導出に必要な浮動小数点演算に対する丸め誤差解析モデルを紹介する．

定理 2.1 (Jeannerod and Rump [5], 2014) $a, b \in \mathbb{F}$ に対して，

$$\begin{aligned} \text{fl}(a \pm b) &= (a \pm b)(1 + \delta_1), \quad |\delta_1| \leq \frac{\mathbf{u}}{1 + \mathbf{u}} < \mathbf{u} \\ \text{fl}(a \pm b) &= (a \pm b) + \delta_2, \quad |\delta_2| \leq \mathbf{u} \cdot \text{ufp}(a \pm b) \end{aligned}$$

が成り立つ．さらに，一般に

$$(2.1) \quad \text{ufp}(a \pm b) \leq \text{ufp}(\text{fl}(a \pm b))$$

となるため，

$$(2.2) \quad |\delta_2| \leq \mathbf{u} \cdot \text{ufp}(a \pm b) \leq \mathbf{u} \cdot \text{ufp}(\text{fl}(a \pm b))$$

が成立する．

定理 2.1 における $|\delta_1|$ の上限は，[2] を含む多くの文献では $|\delta_1| \leq \mathbf{u}$ で抑えられていた．その後， $|\delta_1| \leq \frac{\mathbf{u}}{1 + \mathbf{u}}$ で抑えられ，これが最適であることが Jeannerod と Rump らによって示された [5]．

次に，浮動小数点数の積に対する丸め誤差解析モデル [6] を紹介する．以後， \mathbf{u}_S は非正規化数を含めた浮動小数点数における正の最小数とし，倍精度浮動小数点数であれば

$\mathbf{u}_S = 2^{-1074}$ である. また, \mathbf{u}_N は, 正の正規化数の最小値とする. 倍精度浮動小数点数に対しては $\mathbf{u}_N = 2^{-1022}$ である. ここで, $\mathbf{u}, \mathbf{u}_N, \mathbf{u}_S$ の関係は

$$(2.3) \quad 2\mathbf{u}_N\mathbf{u} = \mathbf{u}_S$$

である.

定理 2.2 (Rump [6], 2012) $a, b \in \mathbb{F}$ に対して,

$$\text{fl}(a \cdot b) = ab + \delta + \eta, \quad |\delta| \leq \mathbf{u} \cdot \text{ufp}(ab), \quad |\eta| \leq \frac{\mathbf{u}_S}{2}, \quad \delta\eta = 0$$

が成り立つ.

定理 2.2 においてアンダーフローが発生するときは $\delta = 0$ と, そうでない場合は $\eta = 0$ としてよい.

$a, b, c \in \mathbb{F}$ に対して, $a \cdot b + c$ に対する浮動小数点演算では, 積と和の計算をあわせて 2 回の丸め誤差が発生する可能性がある. 実数 $ab + c$ を最近点の浮動小数点数に丸めた結果を返す FMA (Fused Multiply-Add) と言われる命令がある. $\text{fl}_F(a, b, c)$ を $a \cdot b + c$ を FMA により計算した結果とすると, 次の定理 2.3 が成立する.

定理 2.3 $a, b, c \in \mathbb{F}$ に対して,

$$\text{fl}_F(a, b, c) = ab + c + \delta + \eta, \quad |\delta| \leq \mathbf{u} \cdot \text{ufp}(ab + c), \quad |\eta| \leq \frac{\mathbf{u}_S}{2}, \quad \delta\eta = 0$$

が成り立つ.

また, ある実数 $a \in \mathbb{R}$ が浮動小数点数であるかを判定する定理を紹介する.

定理 2.4 $a \in \mathbb{R}$ がある. ただし, $|a|$ は \mathbb{F} に属する浮動小数点数の最大値以下であり, $a \in \mathbf{u}_S\mathbb{Z}$ とする. このとき, $c \geq |a| \in \mathbf{uc}\mathbb{Z}$ となる 2 のべき乗数 $c = 2^k$, $k \in \mathbb{Z}$ が存在すれば $a \in \mathbb{F}$ である.

証明 仮定より, 実数 a は絶対値の大きさの意味では浮動小数点数で表現するために適切な範囲にあり, 問題は仮数部に情報が収まるか否かである. 「 $c = |a|$ 」と「 $c > |a|$ 」の場合に分けて考える. $c = |a|$ の場合は a は 2 のべき乗数であるために $a \in \mathbb{F}$ となる. $c > |a|$ の場合, c が 2 のべき乗数であるから $\frac{1}{2}c \geq \text{ufp}(a)$ となる. よって最悪の場合でも $\frac{1}{2}c$ を先頭ビットの単位とし, \mathbf{uc} を最終ビットの単位とする仮数部の範囲に情報が収まるために $a \in \mathbb{F}$ となる.

□

本論文に使用される等式を紹介する.

補題 2.5 $n \in \mathbb{N}$ に対して, $n \leq \mathbf{u}^{-1}$ ならば $\text{fl}(n\mathbf{u}) = n\mathbf{u}$ である.

この補題は定理 2.4 において $c = 1$ とすることにより直ちに証明される. 次章より n の範囲に制約がある場合において, $\text{fl}((n+2)\mathbf{u}) = (n+2)\mathbf{u}$ という関係が用いられる.

次に, 内積の丸め誤差解析に必要な 2 分木を用いた帰納法の概要を説明する. $x, y \in \mathbb{F}^n$ に対する浮動小数点演算による内積計算は, $\tilde{p}_i = \text{fl}(x_i \cdot y_i)$ と計算した後に $\sum_{i=1}^n \tilde{p}_i$ を浮動小数点演算で評価する. ここで, この総和に対する計算順序は 2 分木によって整理できる. 例として $x, y \in \mathbb{F}^8$ における

$$\text{fl}(((((((\tilde{p}_1 + \tilde{p}_2) + \tilde{p}_3) + \tilde{p}_4) + \tilde{p}_5) + \tilde{p}_6) + \tilde{p}_7) + \tilde{p}_8))$$

に対する 2 分木を Fig. 1 に示した. また,

$$\text{fl}(((\tilde{p}_1 + \tilde{p}_2) + (\tilde{p}_3 + \tilde{p}_4)) + ((\tilde{p}_5 + \tilde{p}_6) + (\tilde{p}_7 + \tilde{p}_8)))$$

に対応する 2 分木を Fig. 2 に示した.

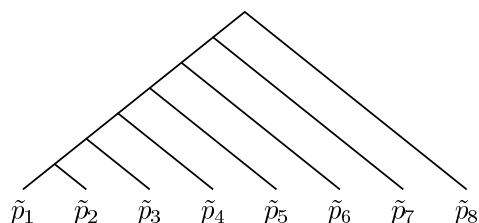


Fig. 1. Binary tree for the recursive order of the summation

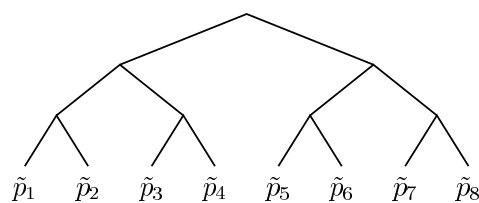
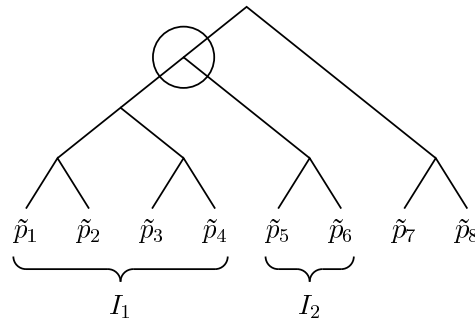


Fig. 2. Binary tree for the pairwise order of the summation

ここで, 2 分木を用いた誤差評価式に関する帰納法の概要について解説する. まず 2 分木のすべての葉について誤差評価式が成立することを確認する. あるノードについて葉から左の子までに計算した \tilde{p}_i のインデックスの集合を I_1 とし, 葉から右の子までに計算した \tilde{p}_i のインデックスの集合を I_2 とする ($I_1 \cap I_2 = \emptyset$). Fig. 3 における丸印のノードに対しては

$$I_1 = \{1, 2, 3, 4\}, \quad I_2 = \{5, 6\}$$

となる. 集合内の要素の個数を $|I_1| = n_1$, $|I_2| = n_2$ とする. 2 分木の葉から左の子ノードと右の子ノードまでの計算について誤差評価式が成立すると仮定する. すなわち誤差評価式が I_1 と I_2 に対して成立すると仮定する. 最後に $I_1 \cup I_2$ に対して誤差評価式が成立することを確認すれば, 帰納的に証明が進み, 最後は 2 分木のルートに関しての誤差評価式が成立する. 詳細は文献 [4] を参照するとよい.

Fig. 3. I_1 and I_2 of the node in binary tree

3. 内積に関する誤差解析

本章のはじめに $x, y \in \mathbb{R}^n$ に対する内積 $x^T y$ の誤差評価式を紹介する. $\text{float}(x^T y)$ はすべての $\text{fl}(x_i \cdot y_i)$ を計算した後, 総和を任意の計算順序にて計算した結果である. よって, $\text{float}(x^T y)$ は内積計算に対する Winograd の方法 [8] など, 計算方法を根本的に変えて得られた結果を含まないとする.

これからベクトルの内積計算に関する誤差評価式について先行研究を紹介する. 演算中にアンダーフローが発生しないと仮定すれば,

$$(3.1) \quad |\text{float}(x^T y) - x^T y| \leq ((1 + \mathbf{u})^n - 1)|x|^T |y|$$

という誤差評価式が知られている (文献 [2] の 3 章を参照). これは第 1 章で紹介したポイント (a) を満たしている. ここで, ベクトル $z \in \mathbb{R}^n$ に対して $|z| = (|z_1|, \dots, |z_n|)^T$ とする. ベクトルの長さに関して $n\mathbf{u} < 1$ という制約をつければ

$$(3.2) \quad |\text{float}(x^T y) - x^T y| \leq \gamma_n |x|^T |y|$$

が知られている (文献 [2] の 3 章を参照). 誤差評価式 (3.2) は, ポイント (b) の意味では誤差評価式 (3.1) より優れている. 近年, 誤差評価式 (3.1) と (3.2) の改良として

$$(3.3) \quad |\text{float}(x^T y) - x^T y| \leq n\mathbf{u} |x|^T |y|$$

が示された [4]. 誤差評価式 (3.3) はポイント (a) が達成されており, 誤差評価式 (3.1) と (3.2) に対してポイント (b), (d) の点でも優れている. もし, 浮動小数点演算の際にアンダーフローを考慮するならば

$$(3.4) \quad |\text{float}(x^T y) - x^T y| \leq n\mathbf{u} |x|^T |y| + \frac{n}{2} \mathbf{u}_s$$

が示される.

また、式 (1.1) にて定義された \mathbf{ufp} を用いた内積の誤差評価についても研究が行われた。 $(n+2)\mathbf{u} \leq 1$ のとき

$$(3.5) \quad |\mathbf{fl}(x^T y) - x^T y| \leq (n+2)\mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}(|x|^T |y|)) + \frac{n}{2}\mathbf{u}_S$$

が導出されている [6]。ここで、 $\mathbf{fl}(x^T y)$ は $\mathbf{fl}(x_i \cdot y_i)$, $1 \leq i \leq n$ を計算した後、前から順に足しこんで得られた計算結果、すなわち

$$\mathbf{fl}((((x_1 \cdot y_1 + x_2 \cdot y_2) + x_3 \cdot y_3) + x_4 \cdot y_4) \dots)$$

とする。ポイント (c) を満たす例として、 $2(n+2)\mathbf{u} \leq 1$ のとき、

$$(3.6) \quad |\mathbf{fl}(x^T y) - x^T y| \leq \mathbf{fl}(((n+2)\mathbf{u}) \cdot \mathbf{ufp}(\mathbf{fl}(|x|^T |y|)) + \mathbf{u}_N)$$

と

$$(3.7) \quad |\mathbf{fl}(x^T y) - x^T y| \leq \mathbf{fl}((n+2) \cdot (\mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}(|x|^T |y|))) + \mathbf{u}_N)$$

が文献 [6] に示されている。

ここで、先行研究 [6] における内積の誤差評価式について考察する。総和に関しては誤差評価式 (1.5) のように任意の計算順に対応している。ただし、内積については誤差評価式 (3.5) のように対応をしていない。また、誤差評価式 (3.5) に $n = 1$ を代入すれば、

$$|\mathbf{fl}(x_1 y_1) - x_1 y_1| \leq 3\mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}(|x_1| |y_1|)) + \frac{1}{2}\mathbf{u}_S$$

となり、定理 2.2 に比べて過大評価なことがすぐにわかる。また、 \mathbf{ufp} を用いた総和の誤差評価式 (1.5) に対しては n の制約はなく、 \mathbf{ufp} を用いた内積の誤差評価式 (3.5) については n の制約がある。一般にアンダーフローに対応するには、小さな定数を最後に足すことで対応しており、他の方法は議論されていない。さらに FMA 命令を用いた場合の内積の丸め誤差解析に関しては先行研究が見当たらない。以上をふまえ、

- 誤差評価式 (3.5) と (3.6) の任意の計算順序化
- $n = 2$ のときの最適な誤差評価式、すなわちポイント (d) の追求
- n が小さいとき、より良い誤差上限の係数を与えられないか？すなわちポイント (d) の改善
- 任意の n に対応した誤差上限の係数の改善、すなわちポイント (a) の達成とポイント (d) の改善
- ポイント (e) に関する新しい対策
- FMA を用いた内積の誤差解析

について研究した成果を順次述べる。

3.1 任意の計算順序に対する ufp を用いたベクトルの内積の誤差評価

ここでは、誤差評価式 (3.5) と (3.6) に相当する誤差評価を任意の計算順序に拡張する.

定理 3.1 $x, y \in \mathbb{F}^n$ とする. $\text{fl}(|x_i \cdot y_i|) \geq \mathbf{u}_N$ のとき $d_i = 1$ とし, そうでない場合には $d_i = 0$ とする. また $d = \sum_{i=1}^n d_i$ とする ($0 \leq d \leq n$). このとき, 次の不等式が成立する.

$$(3.8) \quad |\text{float}(x^T y) - x^T y| \leq (n + 1 + n\mathbf{u} - \mathbf{u})\mathbf{u}A + \frac{n-d}{2}\mathbf{u}_S, \quad A := \text{ufp}(\text{float}(|x|^T |y|))$$

ここでは、不等式左辺と右辺の $\text{float}(x^T y)$ の計算順は同じとする. さらに $n\mathbf{u} \leq 1$ ならば

$$(3.9) \quad |\text{float}(x^T y) - x^T y| \leq (n + 2)\mathbf{u}A + \mathbf{u}_N$$

が成立する. さらに, $2(n + 1)\mathbf{u} \leq 1$ ならば

$$(3.10) \quad |\text{float}(x^T y) - x^T y| \leq \text{fl}(((n + 2)\mathbf{u}) \cdot A + \mathbf{u}_N)$$

と浮動小数点演算のみで評価できる誤差上限を得る.

証明 $p_i = x_i y_i$, $\tilde{p}_i = \text{fl}(x_i \cdot y_i)$ とする. 内積の誤差評価について,

$$(3.11) \quad \begin{aligned} |\text{float}(x^T y) - x^T y| &= |\text{float}(\sum_{i=1}^n \tilde{p}_i) - \sum_{i=1}^n \tilde{p}_i + \sum_{i=1}^n \tilde{p}_i - \sum_{i=1}^n p_i| \\ &\leq |\text{float}(\sum_{i=1}^n \tilde{p}_i) - \sum_{i=1}^n \tilde{p}_i| + \sum_{i=1}^n |\tilde{p}_i - p_i| \end{aligned}$$

と, 積に関する誤差と総和の誤差に分ける. さらに, 総和の誤差については誤差評価式 (1.5) を適用すると

$$(3.12) \quad |\text{float}(x^T y) - x^T y| \leq (n - 1)\mathbf{u}A + \sum_{i=1}^n |\tilde{p}_i - p_i|$$

を得る. ここで, \tilde{p}_i に関する誤差の和 $\sum_{i=1}^n |\tilde{p}_i - p_i|$ の上限について

$$(3.13) \quad \begin{aligned} \sum_{i=1}^n |\tilde{p}_i - p_i| &\leq \mathbf{u} \sum_{i=1}^n |\tilde{p}_i| + \frac{n-d}{2}\mathbf{u}_S \leq \mathbf{u}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|) + (n - 1)\mathbf{u}A) + \frac{n-d}{2}\mathbf{u}_S \\ &\leq \mathbf{u}(2A + (n - 1)\mathbf{u}A) + \frac{n-d}{2}\mathbf{u}_S = (2 + n\mathbf{u} - \mathbf{u})\mathbf{u}A + \frac{n-d}{2}\mathbf{u}_S \end{aligned}$$

が導かれる．この不等式は，はじめに定理 2.2 をすべての i に使用し，次に $\sum_{i=1}^n |\tilde{p}_i|$ の上限について誤差評価式 (1.5) を $\sum_{i=1}^n |\tilde{p}_i|$ について解いた式を用いる．さらに不等式 (1.2) の関係を用いて $\text{float}(\sum_{i=1}^n |\tilde{p}_i|)$ の上限をとることにより得られる．

式 (3.13) における上限を式 (3.12) の $\sum_{i=1}^n |\tilde{p}_i - p_i|$ に代入し，整理すると誤差評価式 (3.8) を得る．さらに， $n\mathbf{u} \leq 1$ ならば $n\mathbf{u} - \mathbf{u} < 1$ であり，さらに $n\mathbf{u} \leq 1$ と等式 (2.3) の関係より $\frac{n-d}{2}\mathbf{u}_S \leq \mathbf{u}_N$ となるため，誤差評価式 (3.9) が得られる．

ここから誤差評価式 (3.10) の証明を行う． $2(n+1)\mathbf{u} \leq 1$ を仮定する．このとき，補題 2.5 より $(n+2)\mathbf{u} = \text{fl}((n+2)\mathbf{u})$ となる．これより「 $(n+2)\mathbf{u}A < \mathbf{u}_N$ 」, 「 $\mathbf{u}_N \leq (n+2)\mathbf{u}A < \frac{1}{2}\mathbf{u}^{-1}\mathbf{u}_N$ 」, 「 $\frac{1}{2}\mathbf{u}^{-1}\mathbf{u}_N \leq (n+2)\mathbf{u}A$ 」の場合に分けて証明を行う．

1. まず $(n+2)\mathbf{u}A < \mathbf{u}_N$ の場合を考える．このとき，定理 2.2 より

$$(3.14) \quad (n+2)\mathbf{u}A = \text{fl}((n+2)\mathbf{u})A = \text{fl}(((n+2)\mathbf{u}) \cdot A) + \eta, \quad |\eta| \leq \frac{1}{2}\mathbf{u}_S$$

が成立する．また，仮定の $(n+2)\mathbf{u}A < \mathbf{u}_N$ より

$$(3.15) \quad \text{fl}(((n+2)\mathbf{u}) \cdot A) \leq \mathbf{u}_N$$

が成立する． n の制限から $n\mathbf{u} - \mathbf{u} < 1$ であり，さらに等式 (2.3) の関係から $\frac{n+1}{2}\mathbf{u}_S < \frac{3}{4}\mathbf{u}_N$ が得られる．不等式 (3.15) より， $\text{fl}(((n+2)\mathbf{u}) \cdot A)$ は最大でも \mathbf{u}_N までとなる．最大で \mathbf{u}_N までの数と非正規化数の和は，定理 2.4 において $\mathbf{u}c = \mathbf{u}_S$ と置くことにより無誤差であることが示される．よって

$$(3.16) \quad \text{fl}(((n+2)\mathbf{u}) \cdot A) + \frac{3}{4}\mathbf{u}_N = \text{fl}(((n+2)\mathbf{u}) \cdot A + \frac{3}{4}\mathbf{u}_N)$$

となる．誤差評価式 (3.8) の上限は，

$$(3.17) \quad (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}A + \frac{n-d}{2}\mathbf{u}_S$$

$$(3.18) \quad \leq (n+2)\mathbf{u}A + \frac{n}{2}\mathbf{u}_S = \text{fl}((n+2)\mathbf{u})A + \frac{n}{2}\mathbf{u}_S$$

$$(3.19) \quad = \text{fl}(((n+2)\mathbf{u}) \cdot A) + \eta + \frac{n}{2}\mathbf{u}_S \leq \text{fl}(((n+2)\mathbf{u}) \cdot A) + \frac{n+1}{2}\mathbf{u}_S$$

$$(3.20) \quad < \text{fl}(((n+2)\mathbf{u}) \cdot A) + \frac{3}{4}\mathbf{u}_N$$

$$(3.21) \quad = \text{fl}(((n+2)\mathbf{u}) \cdot A + \frac{3}{4}\mathbf{u}_N) < \text{fl}(((n+2)\mathbf{u}) \cdot A + \mathbf{u}_N)$$

となる. 上記の変形は, (3.17) から (3.18) の変形は $n\mathbf{u} - \mathbf{u} < 1$ を, (3.18) から (3.19) の変形は等式 (3.14) を, (3.19) から (3.20) の変形は $\frac{n+1}{2}\mathbf{u}_S < \frac{3}{4}\mathbf{u}_N$ を用いている. 最後に, (3.20) から (3.21) の等式変形には等式 (3.16) を用いた.

2. 次に $\mathbf{u}_N \leq (n+2)\mathbf{u}A < \frac{1}{2}\mathbf{u}^{-1}\mathbf{u}_N$ の場合を考える. 定理 2.2 と n の制限により得られる $(n+2)\mathbf{u} < 1$ より

$$(3.22) \quad (n+2)\mathbf{u}A = \text{fl}((n+2)\mathbf{u})A = \text{fl}(((n+2)\mathbf{u}) \cdot A) + \delta$$

$$(3.23) \quad |\delta| \leq \mathbf{u} \cdot \text{ufp}((n+2)\mathbf{u} \cdot A) \leq \frac{1}{2}\mathbf{u}A$$

を得る. また, n の制限により $\frac{n}{2}\mathbf{u}_S < \frac{1}{2}\mathbf{u}_N$, $n\mathbf{u} - \mathbf{u} - \frac{1}{2} < 0$ である. 誤差評価式 (3.8) の上限は,

$$\begin{aligned} & (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}A + \frac{n-d}{2}\mathbf{u}_S \\ (3.24) \quad & = (n+2)\mathbf{u}A + (n\mathbf{u}-\mathbf{u}-1)\mathbf{u}A + \frac{n-d}{2}\mathbf{u}_S \\ (3.25) \quad & \leq \text{fl}(((n+2)\mathbf{u}) \cdot A) + \delta + (n\mathbf{u}-\mathbf{u}-1)\mathbf{u}A + \frac{n}{2}\mathbf{u}_S \\ (3.26) \quad & < \text{fl}(((n+2)\mathbf{u}) \cdot A) + (n\mathbf{u}-\mathbf{u}-\frac{1}{2})\mathbf{u}A + \frac{1}{2}\mathbf{u}_N \\ (3.27) \quad & = \text{fl}(((n+2)\mathbf{u}) \cdot A) + \mathbf{u}_N + (n\mathbf{u}-\mathbf{u}-\frac{1}{2})\mathbf{u}A - \frac{1}{2}\mathbf{u}_N \\ (3.28) \quad & \leq \text{fl}(((n+2)\mathbf{u}) \cdot A) + \mathbf{u}_N - \frac{1}{2}\mathbf{u}_N \\ (3.29) \quad & \leq \text{fl}(((n+2)\mathbf{u}) \cdot A + \mathbf{u}_N) \end{aligned}$$

となる. これより, 上式の導出について説明する. (3.24) から (3.25) の変形には不等式 (3.22) を, (3.25) から (3.26) の変形には不等式 (3.23) と $\frac{n}{2}\mathbf{u}_S < \frac{1}{2}\mathbf{u}_N$ を用いた. (3.27) から (3.28) は $n\mathbf{u} - \mathbf{u} - \frac{1}{2}$ が負であることから, この項をなくすことにより上限を得た. 最後の (3.28) から (3.29) の変形では, $(n+2)\mathbf{u}A$ に対する仮定から得られる $\text{fl}(((n+2)\mathbf{u}) \cdot A) \leq \frac{1}{2}\mathbf{u}^{-1}\mathbf{u}_N$ の条件下で成立する

$$\text{fl}(((n+2)\mathbf{u}) \cdot A) + \mathbf{u}_N = \text{fl}(((n+2)\mathbf{u}) \cdot A + \mathbf{u}_N) + \delta_1, \quad |\delta_1| \leq \frac{1}{2}\mathbf{u}_N$$

を使用した.

3. 最後に $(n+2)\mathbf{u}A \geq \frac{1}{2}\mathbf{u}^{-1}\mathbf{u}_N$ の場合を考える. ここで, n に対すると制約から得る n の最大値と, 仮定 $(n+2)\mathbf{u}A \geq \frac{1}{2}\mathbf{u}^{-1}\mathbf{u}_N$ から得る $\mathbf{u}A$ の最小値を考えれば, 式 (3.26) において

$$(n\mathbf{u}-\mathbf{u}-1)\mathbf{u}A + \frac{1}{2}\mathbf{u}_N \leq 0$$

を得るため、誤差評価式 (3.10) を得る。

□

3.2 長さ 2 のベクトルの内積に関する最適な誤差評価

本節では、長さ 2 のベクトルの内積に関する最適な誤差評価について述べる。ここで最適な誤差評価とは誤差上限がこれ以上小さくできないものとし、誤差上限を表す式が唯一であることは保証しない。この長さ 2 のベクトルの内積は 2 次の正方行列に対する行列式の計算や、複素数の積において現れる。

補題 3.2 $0 \leq a, b \in \mathbb{F}$ に対して

$$(3.30) \quad \text{ufp}(a) + \text{ufp}(b) \leq 1.5\text{ufp}(a + b)$$

が成立する。 $a \neq 0$ のとき $a = m_a \cdot 2^{e_a}$, $b \neq 0$ のとき $b = m_b \cdot 2^{e_b}$ と表す。ここでは $1 \leq m_a, m_b < 2$, $m_a, m_b \in \mathbb{F}$, $e_a, e_b \in \mathbb{Z}$ である。 $e_a = e_b$ または $ab = 0$ のとき、

$$(3.31) \quad \text{ufp}(a) + \text{ufp}(b) = \text{ufp}(a + b)$$

となる。また、 $a \neq 0$ かつ $b \neq 0$ の場合は

$$(3.32) \quad \text{ufp}(a) + \text{ufp}(b) = \frac{1 + 2^{e_b - e_a}}{\text{ufp}(m_a + m_b 2^{e_b - e_a})} \text{ufp}(a + b)$$

が成立する。

証明 a と b のどちらか一方が少なくとも 0 の場合は、不等式 (3.30) と不等式 (3.31) の成立は自明であるため、それ以外の場合を証明する。まず

$$(3.33) \quad \text{ufp}(a) + \text{ufp}(b) = 2^{e_a} + 2^{e_b} = 2^{e_a}(1 + 2^{e_b - e_a})$$

であり、

$$(3.34) \quad \text{ufp}(a + b) = \text{ufp}(m_a 2^{e_a} + m_b 2^{e_b}) = 2^{e_a} \text{ufp}(m_a + m_b 2^{e_b - e_a})$$

となるため、式 (3.32) を得る。 $e_a = e_b$ のとき、 $e_b - e_a = 0$ を式 (3.33) と式 (3.34) に代入し、 $2 \leq m_a + m_b < 4$ から $\text{ufp}(m_a + m_b) = 2$ を式 (3.34) に代入すると式 (3.31) を得る。

最後に式 (3.30) の証明を 2 通りの場合に分けて考える。ここでは $|a| \geq |b|$ を仮定しても一般性を失わないために $e_a \geq e_b$ とし、 $e_a = e_b$ のときには式 (3.31) が示されているために $e_a > e_b$ の場合を考える。式 (3.34) から、 $2 \leq m_a + m_b 2^{e_b - e_a}$ のときは $\text{ufp}(a + b) = 2^{e_a + 1}$ となり、 $2 > m_a + m_b 2^{e_b - e_a}$ のときは $\text{ufp}(a + b) = 2^{e_a}$ となる。以上より $\text{ufp}(a + b)$ の最小値は 2^{e_a} である。また、 $\text{ufp}(a) + \text{ufp}(b)$ の最大値は $e_b - e_a = -1$ のときで $1.5 \cdot 2^{e_a}$ となるため、不等式 (3.30) が成立する。 □

上述の補題を用いて、長さが 2 であるベクトルの内積に対する最適な誤差評価を以下のように得た.

定理 3.3 $a, b, c, d \in \mathbb{F}$ に対して,

$$(3.35) \quad |\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| \leq (2.5 - \mathbf{u})\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(|a \cdot b| + |c \cdot d|))$$

が成立する. この不等式において, $ab \neq 0$ かつ $cd \neq 0$ のときに等号が成り立つ例が存在するため, この誤差評価式は最適である. ただし, $\text{fl}(\cdot)$ 内でアンダーフローが発生しないと仮定する.

証明 a, b, c, d のうち 1 つでも 0 があれば, 定理 2.2 から誤差評価式 (3.35) が成立する. よって, a, b, c, d はすべて 0 でないと仮定し,

$$\begin{aligned} \text{fl}(a \cdot b) &= \text{sign}(ab) \cdot m_1 \cdot 2^{e_1}, & \text{fl}(c \cdot d) &= \text{sign}(cd) \cdot m_2 \cdot 2^{e_2}, \\ e_1, e_2 &\in \mathbb{Z}, & 1 \leq m_1, m_2 < 2, & \quad m_1, m_2 \in \mathbb{F} \end{aligned}$$

とする. ここで sign は括弧内の実数の符号を返す関数とする. $ab + cd$ の計算には対称性があるため, $|\text{fl}(a \cdot b)| \geq |\text{fl}(c \cdot d)|$ の場合を証明すれば良い ($e_1 \geq e_2$). アンダーフローが起きない仮定から

$$\text{fl}(a \cdot b) = ab + \delta_1, \quad \text{fl}(c \cdot d) = cd + \delta_2, \quad \text{fl}(a \cdot b + c \cdot d) = \text{fl}(a \cdot b) + \text{fl}(c \cdot d) + \delta_3$$

とすれば

$$\text{fl}(a \cdot b + c \cdot d) = \text{fl}(a \cdot b) + \text{fl}(c \cdot d) + \delta_3 = ab + \delta_1 + cd + \delta_2 + \delta_3$$

となる. δ_1, δ_2 は定理 2.2 より, また δ_3 は定理 2.1 により

$$\begin{aligned} |\delta_1| &= f_1 \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b)), & |\delta_2| &= f_2 \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(c \cdot d)), \\ |\delta_3| &= f_3 \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b) + \text{fl}(c \cdot d)), & 0 \leq f_1, f_2, f_3 &\leq 1 \end{aligned}$$

とおける. ここで,

$$\begin{aligned} |\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| &\leq |\delta_1| + |\delta_2| + |\delta_3| \\ &= f_1 \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b)) + f_2 \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(c \cdot d)) \\ &\quad + f_3 \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b) + \text{fl}(c \cdot d)) \\ (3.36) \quad &\leq f \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(|a \cdot b| + |c \cdot d|)) \end{aligned}$$

となる $f \in \mathbb{R}$ を考える. 以下「 $e_1 = e_2$ 」, 「 $e_1 - e_2 = 1$ かつ $2 \leq m_1 + m_2 2^{e_2 - e_1}$ 」, 「 $e_1 - e_2 = 1$ かつ $2 > m_1 + m_2 2^{e_2 - e_1}$ 」, 「 $e_1 - e_2 \geq 2$ 」の場合について考えていく.

1. $e_1 = e_2$ のとき, 式 (3.36) において $f_1 = f_2 = f_3 = 1$ とする. さらに補題 3.2 の式 (3.31) より

$$\mathbf{ufp}(\text{fl}(|a \cdot b|)) + \mathbf{ufp}(\text{fl}(|c \cdot d|)) = \mathbf{ufp}(\text{fl}(|a \cdot b| + |c \cdot d|))$$

を用いて式 (3.36) を整理すると,

$$|\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| \leq 2\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(|a \cdot b|) + \text{fl}(|c \cdot d|))$$

となる. さらに不等式 (2.1) を用いれば

$$|\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| \leq 2\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(|a \cdot b| + |c \cdot d|))$$

となるため, $f = 2$ を得る.

2. $e_1 - e_2 = 1$ かつ $2 \leq m_1 + m_2 2^{e_2 - e_1}$ の場合, 式 (3.36) において $f_1 = f_2 = f_3 = 1$ とする. さらに補題 3.2 の式 (3.32) より

$$\mathbf{ufp}(\text{fl}(|a \cdot b|)) + \mathbf{ufp}(\text{fl}(|c \cdot d|)) = 0.75\mathbf{ufp}(\text{fl}(|a \cdot b|) + \text{fl}(|c \cdot d|))$$

を用いて式 (3.36) を整理し, 最後に不等式 (2.1) を用いることにより

$$\begin{aligned} |\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| &\leq 1.75\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(|a \cdot b|) + \text{fl}(|c \cdot d|)) \\ &\leq 1.75\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(|a \cdot b| + |c \cdot d|)) \end{aligned}$$

を導け, $f = 1.75$ を得る.

3. $e_1 - e_2 = 1$ かつ $2 > m_1 + m_2 2^{e_2 - e_1}$ の場合を考える. ここではさらに「 $\text{fl}(c \cdot d)$ の仮数部の最下位ビットが 0 の場合」と「 $\text{fl}(c \cdot d)$ の仮数部の最下位ビットが 1 の場合」の 2 通りに場合分けを行う. 以下, $a \in \mathbb{F}$ に対して $a \in 2\mathbf{u} \cdot \mathbf{ufp}(a)\mathbb{Z}$ が成り立つことを利用する (浮動小数点数は仮数部の最下位ビットの位の整数倍で表現されるという意味である).

- (a) $\text{fl}(c \cdot d)$ の最下位ビットが 0 の場合, $\text{fl}(c \cdot d) \in 4\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(c \cdot d))\mathbb{Z} = 2\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b))\mathbb{Z}$ である. また, $\text{fl}(a \cdot b) + \text{fl}(c \cdot d) \in 2\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b))\mathbb{Z}$ であり, 仮定より $|\text{fl}(a \cdot b) + \text{fl}(c \cdot d)| < 2\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b))$ であるから, 定理 2.4 より

$$\text{fl}(a \cdot b) + \text{fl}(c \cdot d) = \text{fl}(a \cdot b + c \cdot d)$$

となり, 誤差が発生しないために $f_3 = 0$ である. よって式 (3.36) において $f_1 = f_2 = 1$ とし, 補題 3.2 の式 (3.32) を用い,

$$\begin{aligned} |\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| &\leq \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(a \cdot b)) + \mathbf{u} \cdot \mathbf{ufp}(\text{fl}(c \cdot d)) \\ &\leq 1.5\mathbf{u} \cdot \mathbf{ufp}(\text{fl}(|a \cdot b|) + \text{fl}(|c \cdot d|)) \end{aligned}$$

を得て, 最後に不等式 (2.1) を用いれば $f = 1.5$ となる.

- (b) 次に $\text{fl}(c \cdot d)$ の最下位ビットが 1 の場合を考えるが, これは偶数丸めが発生しなかったことを表すために $f_2 \neq 1$ である. よって, $\text{fl}(c \cdot d)$ の誤差が $\mathbf{u} \cdot \mathbf{ufp}(c \cdot d)$ となることはなく, $\mathbf{u} \cdot \mathbf{ufp}(c \cdot d)$ 未満の最大の誤差を考える. アンダーフローがない場合, 浮動小数点数の積の誤差は浮動小数点数とな

る^{*1}. また, $\text{fl}(c \cdot d)$ の誤差が大きいときは $2\text{ufp}(c)\text{ufp}(d) = \text{ufp}(cd)$ となる場合である. 浮動小数点数の基本より, $c \in 2\mathbf{u} \cdot \text{ufp}(c)\mathbb{Z}$, $d \in 2\mathbf{u} \cdot \text{ufp}(d)\mathbb{Z}$ であるから, $cd \in 4\mathbf{u}^2 \cdot \text{ufp}(c)\text{ufp}(d)\mathbb{Z}$ となる. $\mathbf{u} \cdot \text{ufp}(cd)$ より小さい最大の誤差は $\mathbf{u} \cdot \text{ufp}(cd) - 4\mathbf{u}^2 \cdot \text{ufp}(c)\text{ufp}(d) = (\mathbf{u} - 2\mathbf{u}^2) \cdot \text{ufp}(cd)$ となるため, $|\delta_2| \leq (1 - 2\mathbf{u})\mathbf{u} \cdot \text{ufp}(cd)$ となり, $f_2 \leq 1 - 2\mathbf{u}$ を得る. さらに, 仮定である $e_1 - e_2 = 1$ かつ $2 > m_1 + m_2 2^{e_2 - e_1}$ から

$$(3.37) \quad \text{ufp}(\text{fl}(a \cdot b) + \text{fl}(c \cdot d)) = \text{ufp}(\text{fl}(a \cdot b)) = 2\text{ufp}(\text{fl}(c \cdot d))$$

を得る. 式 (3.36) において $f_1 = 1$, $f_2 = 1 - 2\mathbf{u}$, $f_3 = 1$ を代入し, 等式 (3.37) から $\text{ufp}(\text{fl}(a \cdot b)) = \text{ufp}(\text{fl}(a \cdot b) + \text{fl}(c \cdot d))$ と $\text{ufp}(\text{fl}(c \cdot d)) = \frac{1}{2}\text{ufp}(\text{fl}(a \cdot b) + \text{fl}(c \cdot d))$ を使用し, さらに不等式 (2.1) を用いることにより

$$\begin{aligned} & |\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| \\ & \leq \mathbf{u} \cdot \text{ufp}(\text{fl}(a \cdot b)) + (1 - 2\mathbf{u})\mathbf{u} \cdot \text{ufp}(\text{fl}(c \cdot d)) + \mathbf{u} \cdot \text{ufp}(\text{fl}(a \cdot b) + \text{fl}(c \cdot d)) \\ & = (2.5 - \mathbf{u})\mathbf{u} \cdot \text{ufp}(\text{fl}(|a \cdot b|) + \text{fl}(|c \cdot d|)) \leq (2.5\mathbf{u} - \mathbf{u}^2)\text{ufp}(\text{fl}(|a \cdot b|) + \text{fl}(|c \cdot d|)) \end{aligned}$$

を得る. よって $f = 2.5 - \mathbf{u}$ となる.

4. $e_1 - e_2 \geq 2$ の場合, 式 (3.36) において $f_1 = f_2 = f_3 = 1$ とする. さらに補題 3.2 の式 (3.32) より

$$\text{ufp}(\text{fl}(|a \cdot b|)) + \text{ufp}(\text{fl}(|c \cdot d|)) \leq 1.25\text{ufp}(\text{fl}(|a \cdot b|) + \text{fl}(|c \cdot d|))$$

を用いて整理し, 最後に不等式 (2.1) を用いることにより $f = 2.25$ を得る.

以上より, $f = 2.5 - \mathbf{u}$ が最大のため, 誤差評価式 (3.35) は成立する. 最後に最適性について例を挙げて説明する.

$$a = 2^2 \cdot 1.25, \quad b = 1 + 12\mathbf{u}, \quad c = 1.5 - 2\mathbf{u}, \quad d = 1.5 - 2\mathbf{u}$$

のときに

$$\begin{aligned} \text{fl}(a \cdot b + c \cdot d) &= \text{fl}((5 + 64\mathbf{u}) + (2.25 - 4\mathbf{u})) = 4(1.5 + 2^{-2} + 2^{-4} + 16\mathbf{u}) \\ ab + cd &= 4(1.5 + 2^{-2} + 2^{-4} + 13.5\mathbf{u} + \mathbf{u}^2) \end{aligned}$$

となり

$$|\text{fl}(a \cdot b + c \cdot d) - (ab + cd)| = 4(2.5\mathbf{u} - \mathbf{u}^2) = (2.5 - \mathbf{u})\mathbf{u} \cdot \text{ufp}(\text{fl}(a \cdot b + c \cdot d))$$

を得る. よって誤差評価式 (3.35) の右辺をこれ以上小さくできない例が存在するために最適な誤差評価式と言える. \square

^{*1} 文献 [1] により, 2 つの浮動小数点数の積は 2 つの浮動小数点数の和に変換できることからわかる.

3.3 n が小さいときの係数の改良

$n \leq 1 - \log_2 \mathbf{u}$ の範囲 (単精度浮動小数点数では $n \leq 25$, 倍精度浮動小数点数では $n \leq 54$) における誤差評価式 (3.5) の改良を議論する. 集合を $\mathbb{F}' = \mathbf{u}_S \mathbb{Z}$ と定める. $0 \leq p \in \mathbb{F}'$ は, ある $l, h \in \mathbb{Z}$, $\log_2 \mathbf{u}_S \leq l \leq h$ を用いて

$$p = \sum_{i=l}^h d_i \cdot 2^i, \quad d_i \in \{0, 1\}$$

と表せる. p に対する d_i の和を

$$\underline{p} = \sum_{i=l}^h d_i, \quad d_i \in \{0, 1\}$$

と下線を引いた表記を用いる^{*1}. これより誤差評価式を改良するために, いくつかの補題を示す.

補題 3.4 $0 \leq a, b \in \mathbb{F}'$, $p = a + b$ に対して

$$(3.38) \quad \underline{p} \leq \underline{a} + \underline{b}$$

が成立する.

証明 $a = 0$ または $b = 0$ の場合は自明に成立するため, それ以外の場合を考える. a, b を

$$a = \sum_{i=l_a}^{h_a} d_i \cdot 2^i, \quad d_i = \begin{cases} 1 & (i = l_a, h_a) \\ 0 \text{ or } 1 & (l_a + 1 \leq i \leq h_a - 1) \\ 0 & (\text{otherwise}) \end{cases}$$

$$b = \sum_{i=l_b}^{h_b} e_i \cdot 2^i, \quad e_i = \begin{cases} 1 & (i = l_b, h_b) \\ 0 \text{ or } 1 & (l_b + 1 \leq i \leq h_b - 1) \\ 0 & (\text{otherwise}) \end{cases}$$

とおくと,

$$p = \sum_{i=\min(l_a, l_b)}^{\max(h_a, h_b)+1} f_i \cdot 2^i, \quad f_i = \begin{cases} 0 \text{ or } 1 & (\min(l_a, l_b) \leq i \leq \max(h_a, h_b) + 1) \\ 0 & (\text{otherwise}) \end{cases}$$

と表現できる. ここで, $a_j, b_j, p_j = a_j + b_j$, $\min(l_a, l_b) \leq j$ に対して

$$a_j = \sum_{i=\min(l_a, l_b)}^j d_i \cdot 2^i, \quad b_j = \sum_{i=\min(l_a, l_b)}^j e_i \cdot 2^i, \quad p_j = \sum_{i=\min(l_a, l_b)}^{j+1} f_{i,j} \cdot 2^i, \quad f_{i,j} \in \{0, 1\}$$

^{*1} 正規化数においては暗黙の 1 (an implicit hidden bit) も含めた仮数部の 1 の個数と同じである. これは暗黙の 1 を保存しない浮動小数点数のビット表現における仮数部の 1 の個数とは異なる.

を定義する．ここで，

$$(3.39) \quad \underline{p}_j \leq \underline{a}_j + \underline{b}_j$$

を j についての帰納法で証明できれば，

$$a_{\max(h_a, h_b)+1} = a, \quad b_{\max(h_a, h_b)+1} = b, \quad p_{\max(h_a, h_b)+1} = p$$

となるため不等式 (3.38) は証明される．

1. $j = \min(l_a, l_b) =: l$ のとき，以下のように場合分けで考える．

(a) $l_a \neq l_b$ のとき $d_l + e_l = 1$, $(f_{l+1,l}, f_{l,l}) = (0, 1)$ となり， $\underline{p}_l = \underline{a}_l + \underline{b}_l = 1$ から不等式 (3.39) は成り立つ．

(b) $l_a = l_b$ のとき $d_l = 1, e_l = 1$, $(f_{l+1,l}, f_{l,l}) = (1, 0)$ となり， $\underline{p}_l = 1 < \underline{a}_l + \underline{b}_l = 2$ から不等式 (3.39) は成り立つ．

2. $j = k$ のとき $\underline{p}_k \leq \underline{a}_k + \underline{b}_k$ が成り立つと仮定する． $j = k+1$ のとき

$$\begin{aligned} \underline{a}_{k+1} + \underline{b}_{k+1} &= d_{k+1} + \underline{a}_k + e_{k+1} + \underline{b}_k \\ \underline{p}_{k+1} &= f_{k+2,k+1} + f_{k+1,k+1} - f_{k+1,k} + \underline{p}_k \end{aligned}$$

となる．ここで $f_{k+1,k}$ について場合分けを行う．

(a) $f_{k+1,k} = 0$ のとき

$$(f_{k+2,k+1}, f_{k+1,k+1}) = \begin{cases} (1, 0) & (d_{k+1} + e_{k+1} = 2) \\ (0, 1) & (d_{k+1} + e_{k+1} = 1) \\ (0, 0) & (d_{k+1} + e_{k+1} = 0) \end{cases}$$

となり， $f_{k+2,k+1} + f_{k+1,k+1} - f_{k+1,k} \leq d_{k+1} + e_{k+1}$ を得る．よって仮定より $\underline{p}_{k+1} \leq \underline{a}_{k+1} + \underline{b}_{k+1}$ となる．

(b) $f_{k+1,k} = 1$ のとき

$$(f_{k+2,k+1}, f_{k+1,k+1}) = \begin{cases} (1, 1) & (d_{k+1} + e_{k+1} = 2) \\ (1, 0) & (d_{k+1} + e_{k+1} = 1) \\ (0, 1) & (d_{k+1} + e_{k+1} = 0) \end{cases}$$

となり， $f_{k+2,k+1} + f_{k+1,k+1} - f_{k+1,k} \leq d_{k+1} + e_{k+1}$ を得る．よって仮定より $\underline{p}_{k+1} \leq \underline{a}_{k+1} + \underline{b}_{k+1}$ となる．

以上より $\min(l_a, l_b) \leq j$ に対して常に成り立つ． □

補題 3.5 $p \in \mathbb{F}'$ に対して

$$\text{fl}(p) \neq p \Leftrightarrow \underline{\text{fl}}(p) < \underline{p}, \quad \text{fl}(p) = p \Leftrightarrow \underline{\text{fl}}(p) = \underline{p}, \quad \underline{\text{fl}}(p) \leq \underline{p}$$

が成立する． $p \in \mathbb{F}'$ に対して $\text{fl}(p)$ は最近偶数の丸めにより p を \mathbb{F} の要素に対応させる表記とする．

証明 $p \geq 0$ を仮定する.

1. $\text{fl}(p) = p \Rightarrow \underline{\text{fl}(p)} = \underline{p}$ は明らかに成り立つ.
2. $\text{fl}(p) \neq p \Rightarrow \underline{\text{fl}(p)} < \underline{p}$ を示す. p を \mathbb{F} の要素に原点方向に丸めた結果 (切り捨て) を a とし, b は $p = a + b$ を満たすものとする (図 4 参照). ここでは, $a > 2\mathbf{u}_N$ となり, a は非正規化数や 0 にならないことに留意する.

$$\begin{aligned} p &= \pm 2^e \times \boxed{1} \boxed{d_1} \boxed{d_2} \dots \boxed{d_{51}} \boxed{d_{52}} \boxed{d_{53}} \boxed{d_{54}} \dots \\ a &= \pm 2^e \times \boxed{1} \boxed{d_1} \boxed{d_2} \dots \boxed{d_{51}} \boxed{d_{52}} \boxed{0} \boxed{0} \dots \end{aligned}$$

Fig. 4. The relation between $p \in \mathbb{F}'$ and $a \in \mathbb{F}$ ($d_i \in \{0, 1\}$, $\mathbf{u} = 2^{-53}$)

このとき $\underline{p} = \underline{a} + \underline{b}$ であり, $\mathbf{u} \cdot \text{ufp}(p) \geq \text{ufp}(b) > 0$ が成立する. また, 仮定 $\text{fl}(p) \neq p$ より $\underline{b} \geq 1$ となる. ここで, 以下 4 つの場合について考える.

- (a) $\mathbf{u} \cdot \text{ufp}(p) > \text{ufp}(b) > 0$ の場合を考える. これは p の仮数部の下位ビットがそのまま切り捨てられて $\text{fl}(p)$ が得られた場合に相当する. $\text{fl}(p) = a$ となるため, $\underline{p} = \underline{a} + \underline{b} > \underline{a} = \underline{\text{fl}(p)}$ となる.
- (b) $\mathbf{u} \cdot \text{ufp}(p) = \text{ufp}(b)$, $\underline{b} \geq 2$ のときを考える. これは p を $\text{fl}(p)$ に丸める際に切り上げが起きたことを意味する. $\text{fl}(p) = a + 2\mathbf{u} \cdot \text{ufp}(a)$ となるため補題 3.4 より $\underline{\text{fl}(p)} \leq \underline{a} + 1$ となる. よって $\underline{p} = \underline{a} + \underline{b} > \underline{a} + 1 \geq \underline{\text{fl}(p)}$ となる.
- (c) $\mathbf{u} \cdot \text{ufp}(p) = \text{ufp}(b)$, $\underline{b} = 1$, a の仮数部最終ビットが 0 のときを考える. これは $\text{fl}(p)$ が偶数丸め (切り捨て) の結果であることを意味する. $\text{fl}(p) = a$ となるため (a) と同様に $\underline{p} > \underline{\text{fl}(p)}$ となる.
- (d) $\mathbf{u} \cdot \text{ufp}(p) = \text{ufp}(b)$, $\underline{b} = 1$, a の仮数部最終ビットが 1 のときを考える. これは $\text{fl}(p)$ が偶数丸め (切り上げ) の結果であることを意味する. $\text{fl}(p) = a + 2\mathbf{u} \cdot \text{ufp}(a)$ となり, a の仮数部最終ビットに 1 を加えるため繰り上がりが必ず発生する. ここで $w = -\log_2 \mathbf{u}$ とし,

$$p = 2^e \sum_{i=1}^{w+1} d_i 2^{1-i}, \quad d_i = \begin{cases} 0 \text{ or } 1 & (2 \leq i \leq w-1) \\ 1 & (i = 1, w, w+1) \end{cases}, \quad e \in \mathbb{Z}$$

とおく. $\underline{\text{fl}(p)}$ は \underline{p} に対してまず $i = w, w+1$ の 2 つの 1 を失い, 繰り上げのための 1 を 1 つ追加するため $\underline{\text{fl}(p)} < \underline{a} + 1$ となり, $\underline{p} = \underline{a} + \underline{b} > \underline{a} \geq \underline{\text{fl}(p)}$ となる.

以上より $\text{fl}(p) \neq p \Rightarrow \underline{\text{fl}(p)} < \underline{p}$ が成り立つ.

3. $\text{fl}(p) = p \Leftarrow \underline{\text{fl}(p)} = \underline{p}$ の対偶は 2 から証明される.
4. $\text{fl}(p) \neq p \Leftarrow \underline{\text{fl}(p)} < \underline{p}$ は, 1 の対偶から証明される.
5. $\underline{\text{fl}(p)} \leq \underline{p}$ は 1 と 2 より成り立つ.

p に対する $\text{fl}(p)$ と, $-p$ に対する $\text{fl}(-p)$ では同様の証明ができるため, $p < 0$ についても補題は成立する. \square

成分が 2 のべき乗数または 0 であるベクトルの総和について、次の補題が成立する。

補題 3.6 $p \in \mathbb{F}^n$, $\text{ufp}(p_i) = p_i$, $1 \leq i \leq n$, $\alpha := \text{float}(\sum_{i=1}^n p_i)$ に対して,

$$\sum_{i=1}^n p_i = \text{float}(\sum_{i=1}^n p_i) + \Delta, \quad |\Delta| \leq (n - \underline{\alpha})\mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n p_i))$$

が成立する。ただし、上式にある 2 つの $\text{float}(\cdot)$ 内の計算順序は同じとする。

証明 2 章で説明をした 2 分木（ただし葉は p_i ）を用いた帰納法で証明する。まず、2 分木の葉に関しては計算が存在しないために自明に成立する。 $j = 1, 2$ に対して、 $\alpha_j := \text{float}(\sum_{i \in I_j} p_i)$ としたとき

$$\begin{aligned} \sum_{i \in I_j} p_i &= \text{float}(\sum_{i \in I_j} p_i) + \Delta_j = \alpha_j + \Delta_j, \\ |\Delta_j| &\leq (n_j - \underline{\alpha_j})\mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i \in I_j} p_i)) = (n_j - \underline{\alpha_j})\mathbf{u} \cdot \text{ufp}(\alpha_j) \end{aligned}$$

を仮定する。このとき、 $I_3 = I_1 \cup I_2$ に関する計算に関して、定理 2.1 より

$$\alpha_3 := \text{fl}(\alpha_1 + \alpha_2) = \alpha_1 + \alpha_2 + \delta, \quad |\delta| \leq \mathbf{u} \cdot \text{ufp}(\alpha_1 + \alpha_2) \leq \mathbf{u} \cdot \text{ufp}(\alpha_3)$$

を得る。この α_3 に対して

$$\sum_{i \in I_3} p_i = \alpha_3 + \Delta_3, \quad |\Delta_3| \leq (n_3 - \underline{\alpha_3})\mathbf{u} \cdot \text{ufp}(\alpha_3), \quad n_3 = n_1 + n_2$$

の成立を示す。 δ は定理 2.1 から求まる和の誤差上限であり、 $\Delta_3 = \Delta_1 + \Delta_2 - \delta$ の関係が成立している。左の子ノードまでの計算結果と右の子ノードまでの計算結果の和を

$$\alpha_1 + \alpha_2 =: b$$

とすると、補題 3.5 と補題 3.4 より

$$(3.40) \quad \underline{\alpha_3} \leq \underline{b} \leq \underline{\alpha_1} + \underline{\alpha_2}$$

となる。ここで不等式 (3.40) に関して次の場合を考える。

1. $\underline{\alpha_3} = \underline{b} \leq \underline{\alpha_1} + \underline{\alpha_2}$ のとき、補題 3.5 より $\underline{\alpha_3} = \underline{b} \Rightarrow \alpha_3 = b$ となるため、

$$\alpha_3 = \alpha_1 + \alpha_2$$

となるため $\delta = 0$ である。よって

$$\begin{aligned} \Delta_3 &= \Delta_1 + \Delta_2 \\ &\leq (n_1 - \underline{\alpha_1})\mathbf{u} \cdot \text{ufp}(\alpha_1) + (n_2 - \underline{\alpha_2})\mathbf{u} \cdot \text{ufp}(\alpha_2) \\ &\leq (n_1 + n_2 - (\underline{\alpha_1} + \underline{\alpha_2}))\mathbf{u} \cdot \text{ufp}(\alpha_3) \leq (n_3 - \underline{\alpha_3})\mathbf{u} \cdot \text{ufp}(\alpha_3) \end{aligned}$$

となる.

2. $\underline{\alpha}_3 < \underline{b} \leq \underline{\alpha}_1 + \underline{\alpha}_2$ のとき $\underline{\alpha}_3 \leq \underline{\alpha}_1 + \underline{\alpha}_2 - 1$ となるため,

$$\begin{aligned}\Delta_3 &= \Delta_1 + \Delta_2 - \delta \\ &\leq (n_1 - \underline{\alpha}_1)\mathbf{u} \cdot \text{ufp}(\alpha_1) + (n_2 - \underline{\alpha}_2)\mathbf{u} \cdot \text{ufp}(\alpha_2) + \mathbf{u} \cdot \text{ufp}(\alpha_3) \\ &\leq (n_1 + n_2 - (\underline{\alpha}_1 + \underline{\alpha}_2 - 1))\mathbf{u} \cdot \text{ufp}(\alpha_3) \leq (n_3 - \underline{\alpha}_3)\mathbf{u} \cdot \text{ufp}(\alpha_3)\end{aligned}$$

を得る.

以上より補題は証明された. □

補題 3.7 $1 \leq m < 2$, $m \in \mathbb{F}'$ に対して

$$(3.41) \quad m \leq \sum_{i=1}^m 2^{1-i}$$

が成立する. さらに $p \in \mathbb{F}^n$, $p_i = \text{ufp}(p_i)$ が与えられ, $\sum_{i=1}^n p_i \neq 0$, $\sum_{i=1}^n p_i = m \cdot 2^e$, $e \in \mathbb{Z}$ と置く. このとき

$$(3.42) \quad \sum_{i=1}^m 2^{1-i} \leq \sum_{i=1}^n 2^{1-i}$$

が成立する.

証明 不等式 (3.41) に関しては自明である. また, 不等式 (3.42) に関しては, 補題 3.4 を再帰的に使用すれば $\underline{m} \leq n$ を得ることから証明される. □

補題 3.8 $p \in \mathbb{F}^n$, $p_i = \text{ufp}(p_i)$ が与えられ, $1 \leq m_1, m_2 < 2$, $m_1 \in \mathbb{F}'$, $m_2 \in \mathbb{F}$, $e_1, e_2 \in \mathbb{Z}$ とする. $\sum_{i=1}^n p_i \neq 0$ かつ $\text{float}(\sum_{i=1}^n p_i) \neq 0$ を仮定し,

$$\sum_{i=1}^n p_i = m_1 \cdot 2^{e_1}, \quad \text{float}(\sum_{i=1}^n p_i) = m_2 \cdot 2^{e_2}$$

と置く. このとき $n \leq 1 - \log_2 \mathbf{u}$ ならば $e_1 \leq e_2$ が成立する.

証明 $e_1 > e_2$ が成り立つと仮定する. このとき,

$$m_1 2^{e_1} = \sum_{i=1}^n p_i > \text{float}(\sum_{i=1}^n p_i) = m_2 2^{e_2}$$

となる. ある実数 $\Delta > 0$ と $\alpha \in \mathbb{F}$ を

$$m_1 2^{e_1} = m_2 2^{e_2} + \Delta, \quad \alpha := \text{float}(\sum_{i=1}^n p_i)$$

と置く．補題 3.6 から

$$\Delta \leq (n - \underline{\alpha})\mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n p_i)) = (n - \underline{\alpha})\mathbf{u}2^{e_2}$$

となる．補題 3.7 の不等式 (3.41) より $m_2 \leq \sum_{i=1}^{\underline{\alpha}} 2^{1-i}$ であるから

$$\begin{aligned} m_1 2^{e_1} &\leq m_2 2^{e_2} + (n - \underline{\alpha})\mathbf{u}2^{e_2} \leq 2^{e_2} \sum_{i=1}^{\underline{\alpha}} 2^{1-i} + (n - \underline{\alpha})\mathbf{u}2^{e_2} = 2^{e_2}(2 - 2^{-\underline{\alpha}+1} + (n - \underline{\alpha})\mathbf{u}) \\ &= 2^{e_2}(2 + n\mathbf{u} - (2^{-\underline{\alpha}+1} + \underline{\alpha}\mathbf{u})) \end{aligned}$$

を得る．ここで $1 \leq \underline{\alpha} \leq -\log_2 \mathbf{u}$ であり， $\underline{\alpha} = -\log_2 \mathbf{u}$ のときに $2^{-\underline{\alpha}+1} + \underline{\alpha}\mathbf{u}$ は最小になる．さらに $n \leq 1 - \log_2 \mathbf{u}$ ならば

$$m_1 2^{e_1} \leq 2^{e_2}(2 + (n - 2 + \log_2 \mathbf{u})\mathbf{u}) < 2 \cdot 2^{e_2}$$

となり， $e_1 \leq e_2$ が得られ，仮定した $e_1 > e_2$ に矛盾する． □

補題 3.9 $p \in \mathbb{F}^n$ に対して， $n \leq 1 - \log_2 \mathbf{u}$ ならば

$$(3.43) \quad \sum_{i=1}^n \text{ufp}(p_i) \leq (2 - 2^{1-n})\text{ufp}(\text{float}(\sum_{i=1}^n |p_i|))$$

が成立する．

証明 すべての i において $p_i = 0$ ならば不等式 (3.43) は自明に成立する．よって， $\sum_{i=1}^n \text{ufp}(p_i) \neq 0$ かつ $\text{float}(\sum_{i=1}^n \text{ufp}(p_i)) \neq 0$ の場合を考え，下記の表記を導入する．

$$\begin{aligned} \sum_{i=1}^n \text{ufp}(p_i) &= m_1 \cdot 2^{e_1}, \quad \text{float}(\sum_{i=1}^n \text{ufp}(p_i)) = m_2 \cdot 2^{e_2}, \\ e_1, e_2 &\in \mathbb{Z}, \quad 1 \leq m_1, m_2 < 2, \quad m_1 \in \mathbb{F}', m_2 \in \mathbb{F}. \end{aligned}$$

$a \in \mathbb{R}$ に対して不等式 (1.2) より $|a| \geq \text{ufp}(a)$ であるから

$$\text{ufp}(\text{float}(\sum_{i=1}^n |p_i|)) \geq \text{ufp}(\text{float}(\sum_{i=1}^n \text{ufp}(p_i))) = \text{ufp}(m_2 \cdot 2^{e_2}) = 2^{e_2}$$

となる．ここで， $\sum_{i=1}^n \text{ufp}(p_i)$ と $\text{ufp}(\text{float}(\sum_{i=1}^n |p_i|))$ の比を考える．補題 3.8 より $e_1 \leq e_2$ で

あるから

$$(3.44) \quad \frac{\sum_{i=1}^n \text{ufp}(p_i)}{\text{ufp}(\text{float}(\sum_{i=1}^n |p_i|))} \leq \frac{m_1 \cdot 2^{e_1}}{2^{e_2}} \leq m_1$$

を得る．また，補題 3.7 より $m_1 \leq \sum_{i=1}^n 2^{1-i} = 2 - 2^{1-n}$ であるため，この m_1 の上限を式 (3.44) に代入し，分母を払うことにより不等式 (3.43) が得られる． \square

定理 3.10 $x, y \in \mathbb{F}^n$ に対して， $n \leq 1 - \log_2 \mathbf{u}$ のとき

$$|\text{float}(x^T y) - x^T y| \leq (n + 1 - 2^{1-n}) \mathbf{u} \cdot \text{ufp}(\text{float}(|x|^T |y|))$$

が成り立つ．ただし， $\text{float}(\cdot)$ 内でアンダーフローが発生しないとする．さらに不等式の左辺と右辺にある $\text{float}(\cdot)$ の計算順は同じとする．

証明 $p_i = x_i y_i$ ， $\tilde{p}_i = \text{fl}(x_i \cdot y_i)$ とする．内積における誤差の上限を計算するため，不等式 (3.11) のように総和の計算の誤差と 2 項の積に関する誤差に分ける．

$$(3.45) \quad |\text{float}(x^T y) - x^T y| \leq |\text{float}(\sum_{i=1}^n \tilde{p}_i) - \sum_{i=1}^n \tilde{p}_i| + \sum_{i=1}^n |\tilde{p}_i - p_i|$$

ここで，定理 2.2 と補題 3.9 より

$$(3.46) \quad \sum_{i=1}^n |\tilde{p}_i - p_i| \leq \mathbf{u} \sum_{i=1}^n \text{ufp}(\tilde{p}_i) \leq (2 - 2^{1-n}) \mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|))$$

を得る．よって，誤差評価式 (1.5) と不等式 (3.46) をそれぞれ不等式 (3.45) の右辺に使用すれば

$$\begin{aligned} |\text{float}(x^T y) - x^T y| &\leq (n - 1) \mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|)) + (2 - 2^{1-n}) \mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|)) \\ &= (n + 1 - 2^{1-n}) \mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|)) \end{aligned}$$

となる． \square

誤差に関する係数の改良により，定理 3.10 の誤差評価式はポイント (b) の観点では悪化した， $n \leq 1 - \log_2 \mathbf{u}$ のとき

$$|\text{float}(x^T y) - x^T y| \leq (n + 1) \mathbf{u} \cdot \text{ufp}(\text{float}(|x|^T |y|))$$

とポイント (b) を達成する簡単な係数で評価してもよい．

3.4 一般の n に対する係数の改良

ここでは, \mathbf{ufp} を用いた内積の誤差評価について, ポイント (a) を満たす誤差評価式を導出する.

補題 3.11 $p \in \mathbb{F}^n$ に対して,

$$(3.47) \quad \sum_{i=1}^n \mathbf{ufp}(p_i) \leq (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})\mathbf{ufp}(\text{float}(\sum_{i=1}^n |p_i|))$$

が成り立つ.

証明 すべての i に対して $p_i = 0$ ならば不等式 (3.47) は自明に成立する. 以下, 少なくともある i において $p_i \neq 0$ とし, 下記の表記を導入する.

$$\begin{aligned} \sum_{i=1}^n \mathbf{ufp}(p_i) &= m_1 2^{e_1}, \quad \text{float}(\sum_{i=1}^n \mathbf{ufp}(p_i)) = m_2 2^{e_2}, \\ e_1, e_2 &\in \mathbb{Z}, \quad 1 \leq m_1, m_2 < 2, \quad m_1 \in \mathbb{F}', \quad m_2 \in \mathbb{F} \end{aligned}$$

まず

$$(3.48) \quad m_1 2^{e_1} \leq 2^{e_2} (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})$$

となることを以下の 2 通りに場合分けをして示す.

1. $m_1 2^{e_1} \leq m_2 2^{e_2}$ のとき $e_1 \leq e_2$ となる. ここで $f(n) = (n + \log_2 2^{-E} - 2)2^{-E} + 2^{1-n}$ と置く (ただし E は自然数). $n < E$ のとき $f(n)$ は単調減少, $f(E) = f(E+1) = 0$, $E+1 < n$ のとき単調増加であるため, 自然数 n に対しては $f(n) \geq 0$ が証明される. 以上と補題 3.7 により

$$m_1 2^{e_1} \leq 2^{e_1} (2 - 2^{1-n}) \leq 2^{e_2} (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})$$

を得る. よって式 (3.48) は成り立つ.

2. $m_1 2^{e_1} > m_2 2^{e_2}$ のときを考える. $\alpha := m_2 2^{e_2}$ とすると, 補題 3.6 より

$$m_1 2^{e_1} = m_2 2^{e_2} + \Delta, \quad \Delta \leq (n - \underline{\alpha})\mathbf{u} \cdot \mathbf{ufp}(\text{float}(\sum_{i=1}^n \mathbf{ufp}(p_i))) = (n - \underline{\alpha})\mathbf{u} 2^{e_2}$$

が成り立つ. さらに補題 3.7 を用いて

$$\begin{aligned} m_1 2^{e_1} &\leq m_2 2^{e_2} + (n - \underline{\alpha})\mathbf{u} 2^{e_2} \leq 2^{e_2} \sum_{i=1}^{\underline{\alpha}} 2^{1-i} + (n - \underline{\alpha})\mathbf{u} 2^{e_2} \\ &= 2^{e_2} (2 - 2^{-\underline{\alpha}+1} + (n - \underline{\alpha})\mathbf{u}) = 2^{e_2} (2 - \mathbf{u} 2^{-\underline{\alpha}+1-\log_2 \mathbf{u}} + (n - \underline{\alpha})\mathbf{u}) \\ &= 2^{e_2} (2 + n\mathbf{u} - \mathbf{u} (2^{-\underline{\alpha}+1-\log_2 \mathbf{u}} + \underline{\alpha})) \end{aligned}$$

を得る．このとき， α が浮動小数点数であることから $1 \leq \underline{\alpha} \leq -\log_2 \mathbf{u}$ であり， $\underline{\alpha} = -\log_2 \mathbf{u}$ のとき $2^{-\underline{\alpha}+1-\log_2 \mathbf{u}} + \underline{\alpha}$ は最小となるため

$$m_1 2^{e_1} \leq 2^{e_2} (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})$$

となる．よって式 (3.48) は成り立つ．

次に式 (3.48) より

$$\frac{\sum_{i=1}^n \text{ufp}(p_i)}{\text{ufp}(\text{float}(\sum_{i=1}^n \text{ufp}(p_i)))} \leq \frac{2^{e_2} (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})}{2^{e_2}} = 2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u}$$

となり，上式の分母を払い，不等式 (1.2) を用いて

$$\begin{aligned} \sum_{i=1}^n \text{ufp}(p_i) &\leq (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u}) \cdot \text{ufp}(\text{float}(\sum_{i=1}^n \text{ufp}(p_i))) \\ &\leq (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u}) \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |p_i|)) \end{aligned}$$

を得る． □

定理 3.12 $x, y \in \mathbb{F}^n$ に対して

$$|\text{float}(x^T y) - x^T y| \leq (n + 1 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})\mathbf{u} \cdot \text{ufp}(\text{float}(|x|^T |y|))$$

が成立する．ただし，浮動小数点演算の際にアンダーフローが発生しないことを仮定する．また，不等式の左辺と右辺で $\text{float}(\cdot)$ 内の演算順序は同じとする．

証明 $p_i = x_i y_i$, $\tilde{p}_i = \text{fl}(x_i \cdot y_i)$ とする．内積における誤差の上限を計算するため，式 (3.11) のように総和の計算の誤差と積に関する誤差に分ける．

$$(3.49) \quad |\text{float}(x^T y) - x^T y| \leq |\text{float}(\sum_{i=1}^n \tilde{p}_i) - \sum_{i=1}^n \tilde{p}_i| + \sum_{i=1}^n |\tilde{p}_i - p_i|$$

ここで定理 2.2 と補題 3.11 より

$$(3.50) \quad \sum_{i=1}^n |\tilde{p}_i - p_i| \leq \sum_{i=1}^n \mathbf{u} \cdot \text{ufp}(\tilde{p}_i) \leq (2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})\mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|))$$

を得る．よって，誤差評価式 (1.5) と不等式 (3.50) をそれぞれ不等式 (3.49) の右辺にあてはめれば

$$|\text{float}(x^T y) - x^T y| \leq (n - 1)\mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|))$$

$$\begin{aligned}
& +(2 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})\mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|)) \\
& = (n + 1 + (n + \log_2 \mathbf{u} - 2)\mathbf{u})\mathbf{u} \cdot \text{ufp}(\text{float}(\sum_{i=1}^n |\tilde{p}_i|))
\end{aligned}$$

と上限をとれる. □

定理 3.12 の結果は任意の n に対して成立するうえに, 誤差評価式の係数は $(n + 2)\mathbf{u} < 1$ の範囲で誤差評価式 (3.5) の係数よりも小さい.

3.5 内積におけるアンダーフローの取り扱い

誤差評価式 (3.6) と (3.7) では最後に \mathbf{u}_N を足すことでアンダーフローへの対策を施している. これは $n < \mathbf{u}^{-1}$ のときは過大評価であり, また右辺の最小値は \mathbf{u}_N よりも小さくはない. ただし誤差評価式 (3.4) のように非正規化数である $\frac{n}{2}\mathbf{u}_S$ を足しては, 計算のパフォーマンスを下げることがある. 本節ではポイント (e) に関して誤差評価式を改良する.

定理 3.13 $x, y \in \mathbb{F}^n$, $2(n + 1)\mathbf{u} \leq 1$ に対して

$$|\text{float}(x^T y) - x^T y| \leq \text{fl}((n + 2)\mathbf{u} \cdot (\text{ufp}(\text{float}(|x|^T |y|)) + \mathbf{u}_N))$$

が成立する. ここでは, 不等式の左辺と右辺にある $\text{float}(\cdot)$ の計算順は同じとする.

証明 まず $\text{fl}(|x_i \cdot y_i|) \geq \mathbf{u}_N$ ならば $d_i = 1$ であり, $\text{fl}(|x_i \cdot y_i|) < \mathbf{u}_N$ ならば $d_i = 0$ とし, $d = \sum_{i=1}^n d_i$ とする. 表記の簡略化のために $A = \text{ufp}(\text{float}(|x|^T |y|))$ と置く. $d = n$ であれば, 定理 3.1 の誤差評価式 (3.8) と n に関する仮定から得る $\text{fl}((n + 2)\mathbf{u}) = (n + 2)\mathbf{u}$ から

$$|\text{float}(x^T y) - x^T y| \leq \text{fl}((n + 2)\mathbf{u} \cdot A) \leq \text{fl}((n + 2)\mathbf{u} \cdot (A + \mathbf{u}_N))$$

と証明できるため, 以後は $0 \leq d \leq n - 1$ について考える. 以降, A について 2 通りに場合を分けて考える.

1. $A \geq \mathbf{u}^{-1}\mathbf{u}_N$ のとき, $(n + 2)\mathbf{u}$ と 2 のべき乗数 A との積においてアンダーフローが発生しないため

$$(3.51) \quad (n + 2)\mathbf{u}A = \text{fl}(((n + 2)\mathbf{u}) \cdot A) \leq \text{fl}(((n + 2)\mathbf{u}) \cdot (A + \mathbf{u}_N))$$

となる. 定理 3.1 の誤差評価式 (3.8) より

$$\begin{aligned}
|\text{float}(x^T y) - x^T y| & \leq (n + 1 + n\mathbf{u} - \mathbf{u})\mathbf{u}A + \frac{n - d}{2}\mathbf{u}_S \\
& = (n + 2 + n\mathbf{u} - \mathbf{u} - 1)\mathbf{u}A + \frac{n - d}{2}\mathbf{u}_S
\end{aligned}$$

を得る．ここで，仮定による n の制限から得る $n\mathbf{u} - \mathbf{u} - 1 < 0$ と，仮定である $A \geq \mathbf{u}^{-1}\mathbf{u}_N$ から得られる $\mathbf{u}A \geq \mathbf{u}_N$ より

$$(n\mathbf{u} - \mathbf{u} - 1)\mathbf{u}A + \frac{n-d}{2}\mathbf{u}_S \leq (n\mathbf{u} - \mathbf{u} - 1)\mathbf{u}_N + \frac{n}{2}\mathbf{u}_S < 0$$

となり *1,

$$|\text{float}(x^T y) - x^T y| \leq (n+2)\mathbf{u}A$$

が得られ，この式と不等式 (3.51) をつなげればよい．

2. $A < \mathbf{u}^{-1}\mathbf{u}_N$ のときを考える． A が 2 のべき乗数であるから定理 2.4 より

$$(3.52) \quad \text{fl}(A + \mathbf{u}_N) = A + \mathbf{u}_N$$

を導出できる．これは定理 2.4 において $A \leq \mathbf{u}_N$ のとき， $\mathbf{u}c = \mathbf{u}_S$ として，それ以外のときには $\mathbf{u}c = \mathbf{u}_N$ と置くことにより証明される．等式 (2.3) より $\mathbf{u} \cdot \mathbf{u}_N = \frac{1}{2}\mathbf{u}_S$ であることから，

$$(3.53) \quad -(n+1+n\mathbf{u}-\mathbf{u})\mathbf{u} \cdot \mathbf{u}_N + \frac{n}{2}\mathbf{u}_S \leq -\frac{1}{2}\mathbf{u}_S$$

が得られる．ここで，定理 3.1 の誤差評価式 (3.8) から始めて

$$\begin{aligned} |\text{float}(x^T y) - x^T y| &\leq (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}A + \frac{n-d}{2}\mathbf{u}_S \\ &\leq (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}A + \frac{n}{2}\mathbf{u}_S \\ &= (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}(A + \mathbf{u}_N - \mathbf{u}_N) + \frac{n}{2}\mathbf{u}_S \\ (3.54) \quad &= (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}(A + \mathbf{u}_N) - (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u} \cdot \mathbf{u}_N + \frac{n}{2}\mathbf{u}_S \end{aligned}$$

$$(3.55) \quad \leq (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}(A + \mathbf{u}_N) - \frac{1}{2}\mathbf{u}_S$$

$$(3.56) \quad < (n+1+n\mathbf{u}-\mathbf{u})\mathbf{u}(A + \mathbf{u}_N)$$

が導出される．(3.54) から (3.55) の変形は不等式 (3.53) による．ここで，定理 2.2 より $|\delta| < \mathbf{u}$, $|\eta| \leq \frac{1}{2}\mathbf{u}_S$, $\delta\eta = 0$ ，さらに等式 (3.52) を用いて

$$\begin{aligned} \text{fl}((n+2)\mathbf{u} \cdot (A + \mathbf{u}_N)) &= (1+\delta)\text{fl}((n+2)\mathbf{u})\text{fl}(A + \mathbf{u}_N) + \eta \\ (3.57) \quad &= (1+\delta)(n+2)\mathbf{u}(A + \mathbf{u}_N) + \eta \end{aligned}$$

が成立する． $\delta\eta = 0$ より，ここから「 $\eta = 0$ 」と「 $\delta = 0$ 」の場合に分ける．

*1 最後の不等式は，等式 (2.3) より $\mathbf{u}_N = \frac{1}{2}\mathbf{u}^{-1} \cdot \mathbf{u}_S$ であることから，すぐに成立が確認できる．

- (a) 式 (3.57) において $\eta = 0$ を考える. n の制限から $1 - n\mathbf{u} - 2\mathbf{u} > n\mathbf{u} - \mathbf{u}$ が成立するため,

$$\begin{aligned}\text{fl}((n+2)\mathbf{u} \cdot (A + \mathbf{u}_N)) &= (1 + \delta)(n+2)\mathbf{u}(A + \mathbf{u}_N) \geq (1 - \mathbf{u})(n+2)\mathbf{u}(A + \mathbf{u}_N) \\ &= (n+2 - n\mathbf{u} - 2\mathbf{u})\mathbf{u}(A + \mathbf{u}_N) \\ &\geq (n+1 + n\mathbf{u} - \mathbf{u})\mathbf{u}(A + \mathbf{u}_N)\end{aligned}$$

となり, これを式 (3.56) へと不等式をつなげればよい.

- (b) 式 (3.57) において $\delta = 0$ の場合は

$$\begin{aligned}\text{fl}((n+2)\mathbf{u} \cdot (A + \mathbf{u}_N)) &= (n+2)\mathbf{u}(A + \mathbf{u}_N) + \eta \geq (n+2)\mathbf{u}(A + \mathbf{u}_N) - \frac{1}{2}\mathbf{u}_S \\ &\geq (n+1 + n\mathbf{u} - \mathbf{u})\mathbf{u}(A + \mathbf{u}_N) - \frac{1}{2}\mathbf{u}_S\end{aligned}$$

となり, この不等式と式 (3.55) をつなげればよい.

□

定理 3.13 の誤差上限は $n = 1$ かつ $\text{float}(|x|^T|y|) = 0$ のとき最小値 $2\mathbf{u}_S$ となり, 誤差評価式 (3.6) と (3.7) では最小値が \mathbf{u}_N である点が異なる. ここで \mathbf{u}_N を足すことにより誤差上限の変化を考える. 誤差評価式 (3.6) と (3.7) では $(n+2)\mathbf{u}$ をかけて A より小さくなった値に \mathbf{u}_N を足し, 定理 3.13 の誤差上限は A そのものに \mathbf{u}_N を足すことから, $\text{fl}(A + \mathbf{u}_N) = A$ になる場合は多くなるため, 定理 3.13 のほうが誤差上限が良い場合が多いと考える.

3.6 FMA を用いた場合の内積に対する誤差評価

最後に Fused Multiply-Add を用いた $x, y \in \mathbb{F}^n$ に対する内積の最適な誤差評価を与える. $\text{fl}_{\text{Fd}}(x, y)$ は

```
t = fl(x1 · y1)
for i = 2 : n
    t = flF(xi, yi, t)
end
```

と計算する表記とする. 一般に内積計算の丸め誤差解析は, 任意の計算順序も考慮に入れた式である $\text{float}(\cdot)$ を考えてきた. ただし, FMA を使用する場合は任意の計算順序を考えず, 逐次的な計算方法についてのみ結果を与える. $x, y \in \mathbb{F}^4$ に対する内積 $x^T y$ に対して, 前から順に計算する

$$\text{fl}_F(x_4, y_4, \text{fl}_F(x_3, y_3, \text{fl}_F(x_2, y_2, \text{fl}(x_1 \cdot y_1))))$$

においては FMA は 3 回使用可能である. また, 2 項どうしの計算結果を最後に足し合わせる

$$\text{fl}(\text{fl}_F(x_2, y_2, \text{fl}(x_1 \cdot y_1))) + \text{fl}_F(x_4, y_4, \text{fl}(x_3 \cdot y_3))$$

では FMA を 2 回しか使用できない．このため，FMA を使用する回数が計算順により異なるため，本節では計算順を固定する．FMA を用いた内積計算に対して，次の定理が成立する．

定理 3.14 $x, y \in \mathbb{F}^n$ に対して，FMA による $\text{fl}_F(x_i, y_i, \text{fl}_{\text{Fd}}(x_{1:i-1}, y_{1:i-1})) \geq \mathbf{u}_N$ のとき $d_i = 1$ とし，それ以外のときに $d_i = 0$ とする．ただし， d_1 については $\text{fl}(x_1 y_1) \geq \mathbf{u}_N$ のとき $d_1 = 1$ ，そうでないときに $d_1 = 0$ とする． $d = \sum_{i=1}^n d_i$ として

$$\begin{aligned} |\text{fl}_{\text{Fd}}(x, y) - (x^T y)| &\leq d\mathbf{u} \cdot \text{ufp}(\text{fl}_{\text{Fd}}(|x|, |y|)) + \frac{n-d}{2}\mathbf{u}_S \\ &\leq n\mathbf{u} \cdot \text{ufp}(\text{fl}_{\text{Fd}}(|x|, |y|)) + \frac{n}{2}\mathbf{u}_S \end{aligned}$$

が成り立つ．アンダーフローが発生しないとき，この誤差評価式は最適である．

証明 ベクトル $x_{1:k} = (x_1, x_2, \dots, x_k)^T \in \mathbb{F}^k$ という表記を導入し（ y についても同様とする），帰納法により証明する．

$n = 1$ のとき，定理 2.2 より自明に成立する． $n = k$ のとき，

$$|\text{fl}_{\text{Fd}}(x_{1:k}, y_{1:k}) - x^T y| = \sum_{i=1}^k d_i \cdot \mathbf{u} \cdot \text{ufp}(\text{fl}_{\text{Fd}}(|x_{1:k}|, |y_{1:k}|)) + \sum_{i=1}^k \frac{1-d_i}{2}\mathbf{u}_S$$

が成り立つと仮定する．

$n = k+1$ のとき，定理 2.3 を用いて

$$\begin{aligned} &|\text{fl}_{\text{Fd}}(x_{1:k+1}, y_{1:k+1}) - \sum_{i=1}^{k+1} x_i y_i| \\ &= |\text{fl}_F(x_{k+1}, y_{k+1}, \text{fl}_{\text{Fd}}(x_{1:k}, y_{1:k})) - (\sum_{i=1}^k x_i y_i + x_{k+1} y_{k+1})| \\ &= |\text{fl}_{\text{Fd}}(x_{1:k}, y_{1:k}) + x_{k+1} y_{k+1} + \delta + \eta - (\sum_{i=1}^k x_i y_i + x_{k+1} y_{k+1})| \\ &\quad \left(|\delta| \leq d_{k+1} \cdot \mathbf{u} \cdot \text{ufp}(\text{fl}_{\text{Fd}}(|x_{1:k}|, |y_{1:k}|) + |x_{k+1}| |y_{k+1}|), \quad |\eta| \leq \frac{1-d_{k+1}}{2}\mathbf{u}_S \right) \\ &\leq |\text{fl}_{\text{Fd}}(x_{1:k}, y_{1:k}) - \sum_{i=1}^k x_i y_i| + |\delta| + |\eta| \\ &\leq \sum_{i=1}^k d_i \cdot \mathbf{u} \cdot \text{ufp}(\text{fl}_{\text{Fd}}(|x_{1:k}|, |y_{1:k}|)) + \sum_{i=1}^k \frac{1-d_i}{2}\mathbf{u}_S \\ &\quad + d_{k+1} \mathbf{u} \cdot \text{ufp}(\text{fl}_{\text{Fd}}(|x_{1:k+1}|, |y_{1:k+1}|)) + \frac{1-d_{k+1}}{2}\mathbf{u}_S \end{aligned}$$

$$\leq \sum_{i=1}^{k+1} d_i \cdot \mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}_{\mathbf{Fd}}(|x_{1:k+1}|, |y_{1:k+1}|)) + \sum_{i=1}^{k+1} \frac{1-d_i}{2} \mathbf{u}_S$$

となり, $n = k+1$ のときも成り立つ. 以上より, すべての自然数に対して成立する. ここで $0 \leq d \leq n$ であるから

$$|\mathbf{fl}_{\mathbf{Fd}}(x, y) - (x^T y)| \leq n\mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}_{\mathbf{Fd}}(|x|, |y|)) + \frac{n}{2} \mathbf{u}_S$$

が成立する.

もし, アンダーフローがすべてに発生しない場合には

$$(3.58) \quad |\mathbf{fl}_{\mathbf{Fd}}(x, y) - (x^T y)| \leq n\mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}_{\mathbf{Fd}}(|x|, |y|))$$

となる. この不等式において $x = (1 + 2^{-2}, 1, \dots, 1)^T$, $y = (1 + 2^2\mathbf{u}, \mathbf{u}, \dots, \mathbf{u})^T \in \mathbb{F}^n$ の場合を考える. ベクトルの内積は

$$\begin{aligned} x^T y &= (1 + 2^{-2}, 1, \dots, 1)(1 + 2^2\mathbf{u}, \mathbf{u}, \dots, \mathbf{u})^T \\ &= (1 + 2^{-2})(1 + 2^2\mathbf{u}) + 1 \cdot \mathbf{u} + \dots + 1 \cdot \mathbf{u} \\ &= (1 + 2^{-2} + 2^2\mathbf{u} + \mathbf{u} + \dots + \mathbf{u}) = 1 + 2^{-2} + 2^2\mathbf{u} + n\mathbf{u} \end{aligned}$$

となる. また,

$$\mathbf{fl}_{\mathbf{F}}(1, \mathbf{u}, \mathbf{fl}((1 + 2^{-2})(1 + 2^2\mathbf{u}))) = \mathbf{fl}_{\mathbf{F}}(1, \mathbf{u}, 1 + 2^{-2} + 2^2\mathbf{u}) = 1 + 2^{-2} + 2^2\mathbf{u}$$

であるから

$$\mathbf{fl}_{\mathbf{Fd}}(x, y) = 1 + 2^{-2} + 2^2\mathbf{u}$$

を得る. また,

$$\mathbf{ufp}(\mathbf{fl}_{\mathbf{Fd}}(|x|, |y|)) = \mathbf{ufp}(1 + 2^{-2} + 2^2\mathbf{u}) = 1$$

となる. よって

$$\begin{aligned} |\mathbf{fl}_{\mathbf{Fd}}(x, y) - x^T y| &= |(1 + 2^{-2} + 2^2\mathbf{u}) - (1 + 2^{-2} + 2^2\mathbf{u} + n\mathbf{u})| \\ &= n\mathbf{u} = n\mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}_{\mathbf{Fd}}(|x|, |y|)) \end{aligned}$$

となるため, $|x|^T |y| \neq 0$ に対して不等式 (3.58) の等号が成り立つ例がすべての n に存在するため, 不等式は最適である. \square

定理 3.14 の結果については, ポイント (a), (b), (d), (e) を達成できている. ポイント (c) を達成するならば, $n\mathbf{u} \leq 1$ のときに

$$|\mathbf{fl}_{\mathbf{Fd}}(x, y) - (x^T y)| \leq \mathbf{fl}(n\mathbf{u} \cdot \mathbf{ufp}(\mathbf{fl}_{\mathbf{Fd}}(|x|, |y|)))$$

と右辺全体を浮動小数点演算にて評価できる.

FMA を用いた内積のアンダーフローにも対応する誤差評価について, 定理 3.13 のように \mathbf{u}_N を加える位置を工夫できたらよいが, それが不可能なことを以下に示す.

注意 1 $x, y \in \mathbb{F}^n$ に対して, $n\mathbf{u} \leq 1$ のとき

$$(3.59) \quad |\text{fl}_{\text{Fd}}(x, y) - x^T y| \leq \text{fl}(n\mathbf{u} \cdot (\text{ufp}(\text{fl}_{\text{Fd}}(|x|, |y|)) + \mathbf{u}_N))$$

は少なくとも $n = 4k - 3$, $k \in \mathbb{N}$ のときに成立しない.

証明 $x = (0.5, \dots, 0.5)^T$, $y = (\mathbf{u}_S, \dots, \mathbf{u}_S)^T$ を考える. このとき, $x^T y = \frac{n}{2} \mathbf{u}_S$ であり, また $\text{fl}_{\text{Fd}}(x, y) = 0$ となる. ただし, $n = 4k - 3$, $k \in \mathbb{N}$ のとき, $\text{fl}(n\mathbf{u} \cdot (\text{ufp}(\text{fl}_{\text{Fd}}(|x|, |y|)) + \mathbf{u}_N)) = \frac{n-1}{2} \mathbf{u}_S$ であることから不等式 (3.59) は成立しない. \square

4. おわりに

本稿では, 浮動小数点演算における内積の丸め誤差の評価について, 先行研究を改良または拡張した誤差評価式を導出できた. 具体的には, 長さが 2 のベクトルの内積における最適な誤差評価式の導出, 長さが小さいときの内積の誤差評価の改良, 長さに制約のない誤差評価式の導出である. さらに, 丸め誤差解析におけるアンダーフローに関する対策法や, Fused Multiply-Add を用いた場合の内積の丸め誤差解析について提案した.

謝辞 終始丁寧に査読をいただき, 有益なコメントをいただきました 2 名の査読者の方に深く感謝をいたします. 本研究は, 科学研究費補助金・若手研究 B, 誤差解析の意味で高精度な数値線形計算の基盤の創出 (研究課題番号: 25730076) の支援を受けた.

参考文献

- [1] Dekker, T.J., A floating-point technique for extending the available precision, *Numerische Mathematik*, **18** (1971), 224–242.
- [2] Higham, N.J., *Accuracy and Stability of Numerical Algorithms*, second edition, SIAM Publications, Philadelphia, 2002.
- [3] IEEE Standard for Floating-Point Arithmetic, Std 754–2008, 2008.
- [4] Jeannerod, C.-P., Rump, S.M., Improved error bounds for inner products in floating-point arithmetic, *SIAM. J. Matrix Anal. & Appl. (SIMAX)*, **34** (2013), 338–344.
- [5] Jeannerod, C.-P., Rump, S.M., On relative errors of floating-point operations: optimal bounds and applications, *Research Report*, 2014, pp.15, hal-00934443.
- [6] Rump, S.M., Error estimation of floating-point summation and dot product, *BIT Numer. Math.*, **52** (2012), 201–220.
- [7] Rump, S.M., Computable backward error bounds for basic algorithms in linear algebra. *Nonlinear Theory and Its Applications, IEICE*, **6** (2015), 360–363.

- [8] Winograd, S., A new algorithm for inner product, IEEE Trans. Comput., C-18 (1968), 693–694.

樋口 裕幸 (非会員) 〒337-8570 さいたま市見沼区深作 307

2014 年芝浦工業大学システム理工学部数理科学科卒業。現在，芝浦工業大学大学院理工学研究科システム理工学専攻修士課程。応用数学，特に浮動小数点演算の誤差解析の理論に興味を持つ。

尾崎 克久 (正会員) 〒337-8570 さいたま市見沼区深作 307

2007 年早稲田大学大学院理工学研究科博士課程修了，博士 (工学)。現在，芝浦工業大学システム理工学部数理科学科准教授。精度保証付き数値計算，誤差解析，高精度計算の研究に従事。