

УДК 004.622:004.222.3:004.222.5

**С.В. Иваница**

Донецкий национальный технический университет (ДонНТУ) МОН Украины  
Украина, 83001, Донецк, ул. Артема, 58

## Оценка погрешности представления вещественных чисел в постбинарных форматах с плавающей запятой

**S.V. Ivanitsa**

Donetsk national technical university (DonNTU) MES of Ukraine  
Ukraine, 83001, c. Donetsk, Artema st., 58

## *Error Estimation of Representation of Real Numbers in the Postbinary Floating-Point Formats*

**С.В. Іваниця**

Донецький національний технічний університет (ДонНТУ) МОН України  
Україна, 83001, м. Донецьк, вул. Артема, 58

## Оцінка похибки подання дійсних чисел у постбінарних форматах з рухомою комою

В статье рассмотрена точность представления чисел с плавающей запятой и показаны формулы для расчета абсолютной и относительной погрешностей для бинарных и постбинарных форматов. Получены формулы для расчета погрешностей всех форматов и выявлена закономерность между значениями погрешности постбинарного и бинарного форматов одного класса точности. Рассмотрены стандартные способы округления чисел с плавающей запятой и предложен новый способ округления чисел в постбинарных форматах.

**Ключевые слова:** постбинарный формат числа, погрешность, округление чисел, интервал.

In the article, precision of floating point numbers are considered, the formulas for calculating of the absolute and relative errors for binary and postbinary formats are shown. The formulas for calculating the errors of all sizes and pattern are obtained, rules for values of error of postbinary and binary formats one class of accuracy are found. Standard methods of rounding floating point numbers are considered. New method of number rounding in postbinary formats is proposed.

**Key words:** postbinary number format, accuracy, number rounding, interval.

Розглянуто точність подання чисел з рухомою комою і показані формули для розрахунку абсолютної і відносної похибок для бінарних і постбінарних форматів. Отримано формули для розрахунку похибок всіх форматів і виявлена закономірність між значеннями похибки постбінарного і бінарного форматів одного класу точності. Розглянуто стандартні способи округлення чисел з плаваючою комою і запропоновано новий спосіб округлення чисел в постбінарних форматах.

**Ключові слова:** постбінарний формат числа, похибка, округлення чисел, інтервал.

## Введение

Функциональная структура использования систем искусственного интеллекта (СИИ) состоит из трех комплексов вычислительных средств [1], основным среди которых является определяющий систему программных и аппаратных средств интеллектуальный интерфейс. Интеллектуальный интерфейс обеспечивает использование

компьютерных средств для решения задач, которые возникают в среде профессиональной деятельности конечного пользователя. При этом СИИ имеет возможность решать задачи различной сложности и оперировать различными типами данных, в том числе и вещественными.

Произвольное вещественное число представляется бесконечной систематической (например, десятичной или двоичной) дробью. На практике в научных и инженерных вычислениях вещественные числа приходится представлять в компьютере конечными дробями, чаще всего числами с плавающей запятой (числами, представленными в формате IEEE 754). Следовательно, числа с плавающей запятой представляют конечное множество, на которое отображается бесконечное множество вещественных чисел. Поэтому исходное число может быть представлено в формате IEEE 754 не точно. Это один из ряда недостатков чисел формата IEEE 754–2008, которые подробно рассмотрены в [2–4].

В качестве увеличения точности представления чисел, а также повышения надежности вычислений, для чисел с плавающей точкой были предложены постбинарные форматы чисел от одинарной до квадратичности [2, с. 202]. Используемый в постбинарных форматах способ кодирования данных основан на принципах кодо-логического базиса [5]: в качестве кодовой системы выступает тетракод  $T$ , а в качестве единицы хранения одного разряда – тетрит  $t$ ,

$$T = \{t\}, \quad t \in \{0, 1, A, M\}, \quad (1)$$

кодирующий одно из четырех значений: двоичный ноль (0), двоичную единицу (1), значение неопределенности (A), значение множественности (M). Причем, при кодировании числовых значений, тетрит  $t = A$  может принимать любое (случайное) значение 0 или 1, а тетрит  $t = M$  – и 0 и 1 одновременно (т.е. представлять два числовых набора).

Такое «гибкое» кодирование количественных значений позволяет с высокой степенью точности представлять числа в форматах с плавающей запятой.

**Целью данной статьи** является описание алгоритмов реализации как существующих стандартных видов округления, так и нового – постбинарного округления, а также получение формул для расчета погрешностей представления чисел для стандартных двоичных (binary) и постбинарных форматов pbinary (от англ. «postbinary» – фактически «постбинарный») [6].

## Точность представления вещественных чисел в формате с плавающей запятой

При представлении числа в виде полей порядка, мантииссы и знака, на вещественной оси можно отложить конечный набор значений, в общем случае не превосходящий

$$P_{\Omega} = 2^{\Omega} = 2^{s+l+m}, \quad (2)$$

где  $P_{\Omega}$  – количество чисел в формате с плавающей запятой, представленное  $\Omega$ -разрядным двоичным значением;  $s, l, m$  – разрядности знака, порядка и мантииссы соответственно.

Имея в арсенале конечное количество кодируемых в формате с плавающей запятой значений (**базовых точек**) вещественной оси, вся дальнейшая процедура представления любого вещественного числа сводится к отображению его одной (чаще всего ближайшей) базовой точкой. Процедура подбора такой точки для исходного вещественного числа называется **округлением числа**, а расстояние между действительной позицией числа на вещественной оси и базовой точкой его отображения – **абсолютной погрешностью**  $\Delta$  представления числа в формате с плавающей запятой.

Естественно, двоичные эквиваленты некоторых десятичных дробей совпадают с базовыми точками (при условии, что само число входит в диапазон представления данного формата с плавающей запятой). В этом случае такое число представляется в формате с плавающей точкой без округления (и, соответственно, без потери точности). К ним относятся числа, не имеющие дробную часть (множество целых чисел), и числа, которые можно представить в виде конечной двоичной дроби, дробная часть которых кратна  $2^{-z}$ , где  $z \in \mathbf{Z}$  («кратность отрицательной степени двойки»).

Остальные числа представлены в формате с плавающей запятой приближенно, т.е. с некоторой ошибкой точности.

На рис. 1 приведен график ошибки точности (погрешности) представления числа в формате IEEE 754–2008 при увеличении порядка на 1 (т.е. от  $k$  до  $k+1$ ). При этом базовые точки  $n_i$  – фактически разрядная сетка поля мантииссы формата с плавающей запятой (на рис. 1 – светлые точки).

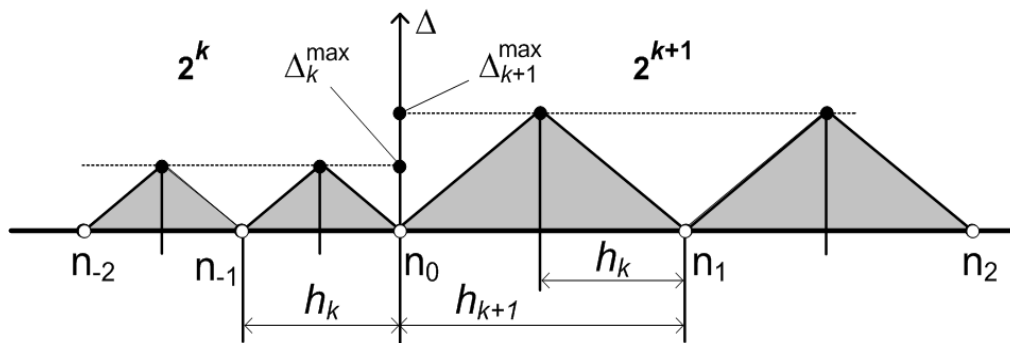


Рисунок 1 – График ошибки точности представления чисел в формате IEEE 754–2008

По форме графика видно, что максимальная абсолютная ошибка точности приходится на число, равноудаленное от двух подряд лежащих на вещественной оси базовых точек, а в самих базовых точках ошибка равна нулю (что говорит о точном представлении числа). Следовательно, чем больше расстояние между позицией исходного числа и ближайшей базовой точкой, тем с меньшей точностью оно может быть представлено этой точкой. Расстояние  $h_k$  между соседними базовыми точками, т.е. позициями чисел, представленных полями форматов с единым десятичным значением порядка  $k$  ( $k \neq 0$ ) и с различающимися в один бит мантииссами, можно определить по формуле:

$$h_k = 2^{E_k - \text{off} - m}, \quad (3)$$

где  $\text{off}$  – десятичное смещение порядка  $E$ ;  $m$  – число разрядов поля мантииссы;  $E_k$  – десятичное число, полученное из поля порядка двоичного числа с плавающей запятой (значение смещенного порядка).

Расстояние  $h_k$  – это фактически шаг чисел одного порядка, который удваивается с увеличением порядка числа с плавающей запятой на единицу:

$$h_{k+1} = 2 \cdot h_k. \quad (4)$$

Иными словами, чем дальше от нуля, тем шире шаг чисел в формате IEEE 754–2008 по числовой оси.

Следует отметить, что формулы (3) и (4) справедливы для расчета шага нормализованных чисел [7]. Для денормализованных чисел (при  $E_k = 0$  и ненулевой мантииссе) примем значение  $k = 0$ , показывая что такие числа близки к нулю. Шаг  $h_0$  для множества денормализованных чисел можно определить следующим образом:

$$h_0 = 2^{1 - \text{off} - m}. \quad (5)$$

Очевидно, что  $h_0$  – минимальный шаг, равный по своему значению минимальному денормализованному числу.

Абсолютная максимальная ошибка для нормализованных ( $\Delta_k^{\max}$ ) и денормализованных ( $\Delta_0^{\max}$ ) чисел в формате IEEE 754–2008 равна в пределе половине шага:

$$\Delta_k^{\max} = \frac{1}{2} h_k = \frac{1}{2} (2^{E_k - \text{off} - m}) = 2^{E_k - \text{off} - m - 1}; \quad (6)$$

$$\Delta_0^{\max} = \frac{1}{2} h_0 = \frac{1}{2} (2^{1 - \text{off} - m}) = 2^{-(\text{off} + m)}. \quad (7)$$

Учитывая (4) и (6), для нормализованных чисел можно получить зависимость:

$$\Delta_{k+1}^{\max} = 2 \cdot \Delta_k^{\max}, \quad (8)$$

то есть значение абсолютной максимальной ошибки точности представления удваивается с увеличением порядка числа с плавающей запятой на единицу.

Относительная максимальная ошибка нормализованных ( $\delta_k^{\max}$ ) и денормализованных ( $\delta_0^{\max}$ ) чисел в формате IEEE 754–2008 равна

$$\delta_k^{\max} = \frac{\Delta_k^{\max}}{|F_k|} \cdot 100\%; \quad (9)$$

$$\delta_0^{\max} = \frac{\Delta_0^{\max}}{|F_0|} \cdot 100\%, \quad (10)$$

где  $F_k$  и  $F_0$  – десятичные числа, полученные из форматов представления нормализованных и денормализованных чисел с плавающей запятой соответственно:

$$|F_k| = 2^{E_k - \text{off}} \cdot \left(1 + \frac{M_n}{2^m}\right); \quad (11)$$

$$|F_0| = 2^{1 - \text{off}} \cdot \left(\frac{M_n}{2^m}\right), \quad (12)$$

где  $M_n$  – десятичное число, полученное из  $m$ -разрядного поля мантииссы двоичного числа с плавающей запятой ( $0 < n < 2^m$ ).

Окончательно получаем:

$$\delta_k^{\max} = \frac{1}{2^{m+1} + 2 \cdot M_n} \cdot 100\%; \quad (13)$$

$$\delta_0^{\max} = \frac{1}{2 \cdot M_n} \cdot 100\%, \quad (14)$$

Для наглядности, формулы нахождения значений максимальной абсолютной и относительной погрешностей представления нормализованных и денормализованных чисел для бинарных и постбинарных форматов сведены в табл. 1. Отличия в полученных числовых значениях погрешностей для  $\Omega$ -разрядных бинарных (binary $\Omega$ ) и аналогичных им постбинарных форматов (pbinary $\Omega$ ), связаны, прежде всего, с организацией полей постбинарного формата, разрядность мантииссы  $m$  которого меньше на число разрядов  $id$  поля идентификатора формата (табл. 1 и [2, с. 203, рис. 4.16]).

Таблица 1 – Формулы максимальной абсолютной и относительной ошибок представления нормализованных и денормализованных чисел для бинарных (binary) и постбинарных (pbinary) форматов

Формат	$off$	$m$	$id$	$\Delta_0^{\max}$	$\Delta_k^{\max}$	$\delta_k^{\max}, \%$
binary16*	15	10	–	$2^{-25}$	$2^{E_k - 26}$	$(2^{11} + 2M_n)^{-1} \cdot 100\%$
pbinary16	15	9	1	$2^{-24}$	$2^{E_k - 25}$	$(2^{10} + 2M_n)^{-1} \cdot 100\%$

Продолж. табл. 1

binary32	127	23	–	$2^{-150}$	$2^{E_k-151}$	$(2^{24} + 2M_n)^{-1} \cdot 100\%$
<b>pbinary32</b>	127	21	2	$2^{-148}$	$2^{E_k-149}$	$(2^{22} + 2M_n)^{-1} \cdot 100\%$
binary64	1023	52	–	$2^{-1075}$	$2^{E_k-1076}$	$(2^{53} + 2M_n)^{-1} \cdot 100\%$
<b>pbinary64</b>	1023	48	4	$2^{-1071}$	$2^{E_k-1072}$	$(2^{49} + 2M_n)^{-1} \cdot 100\%$
binary128	16383	112	–	$2^{-16495}$	$2^{E_k-16496}$	$(2^{113} + 2M_n)^{-1} \cdot 100\%$
<b>pbinary128</b>	16383	104	8	$2^{-16487}$	$2^{E_k-16488}$	$(2^{105} + 2M_n)^{-1} \cdot 100\%$
<b>pbinary256</b>	524287	219	16	$2^{-524506}$	$2^{E_k-524507}$	$(2^{220} + 2M_n)^{-1} \cdot 100\%$

\* – не входит в стандарт IEEE 754–2008, однако определен как формат чисел половинной точности.

Таким образом, величины абсолютной ошибки  $\delta_k$  числа, представленного в формате binary $\Omega$ , и абсолютной ошибки  $\delta'_k$  числа, представленного в аналогичном ему формате pbinary $\Omega$  (т.е. форматы с равными порядками  $E_k$  ( $k \geq 0$ ) и разрядностью  $\Omega$ ) связаны следующим соотношением:

$$\delta'_k = 2^{id} \cdot \delta_k, \quad (15)$$

где  $id$  – разрядность идентификатора формата pbinary $\Omega$ , являющаяся фактически разницей между разрядностями мантисс форматов binary $\Omega$  и pbinary $\Omega$ .

Из формулы (15) следует, что за счет «укороченной» мантиссы формат pbinary $\Omega$  имеет абсолютную ошибку представления числа, большую в  $2^{id}$  раза, чем у формата binary $\Omega$ .

Этот недостаток постбинарного формата с плавающей запятой компенсируется, во-первых, поддержкой динамической разрядности, при которой число в процессе вычисления может выражаться в форматах с возрастающей точностью, что приведет к уменьшению ошибки точности, и, во-вторых, использованием тетракода для дополнения стандартных способов округления новым, т.н. **постбинарным округлением**.

## Способы округления чисел формата IEEE 754–2008

Стандарт IEEE 754–2008 предусматривает четыре способа округления чисел:

1. **К нулю** (рис. 2 а).

При округлении к нулю нужно просто отбросить незначимые разряды числа, поэтому этот способ самый легкий в аппаратной реализации.

2. **К ближайшему числу** (к ближайшей базовой точке) (рис. 2 б).

При округлении к ближайшему числу нужно **к мантиссе прибавить значение старшего незначимого разряда числа**. Данный способ округления в математическом сопроцессоре используется по умолчанию.

3. **К положительной бесконечности** ( $+\infty$ ) (рис. 2 в).

Округление к  $+\infty$  применяется при кодировании интервальных чисел [8].

При округлении к  $+\infty$  нужно к полю мантиссы прибавить инверсное значение поля знака  $\#S$ , поскольку:

– если число положительное ( $S = 0$ ) – нужно к мантиссе **всегда прибавлять 1**;

– если число отрицательное ( $S = 1$ ) – нужно просто **отбросить незначимые разряды**.

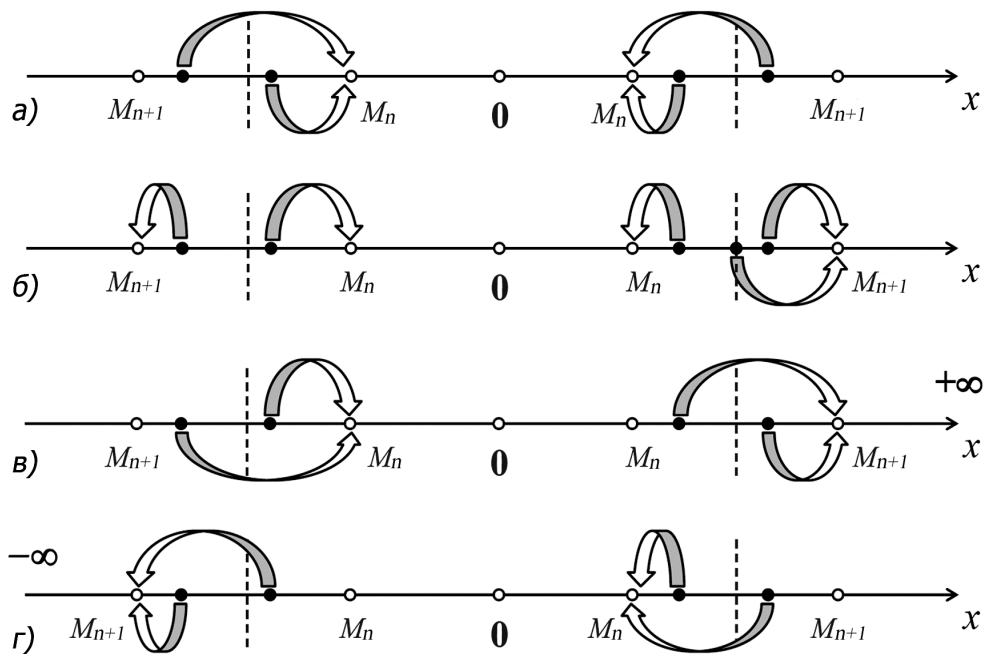


Рисунок 2 – Принцип округления чисел стандарта IEEE 754–2008 (а – к нулю, б – к ближайшему числу, в – к положительной бесконечности, г – к отрицательной бесконечности)

#### 4. К отрицательной бесконечности ( $-\infty$ ) (рис. 2 г).

Округление к  $-\infty$  применяется при кодировании интервальных чисел.

При округлении к  $-\infty$  нужно к полю мантиссы прибавить значение поля знака  $S$ , поскольку:

– если число положительное ( $S = 0$ ) – нужно просто **отбросить незначащие разряды**;

– если число отрицательное ( $S = 1$ ) – нужно к мантиссе **всегда прибавлять 1**.

На рис. 2 светлые точки являются базовыми, а темные – позициями исходных чисел. Пунктирными линиями обозначена середина между соседними базовыми точками.

Таким образом, для реализации стандартных видов округления достаточно проанализировать **старший незначащий разряд двоичной дроби исходного числа**, т.е. первый разряд двоичного числа, не попавший в поле мантиссы, а также учитывать знак числа.

Очевидно, что рассмотренные формулы абсолютной и относительной максимальной ошибки (6), (7) и (13), (14) справедливы только при округлении к ближайшему числу.

## Постбинарное округление чисел с плавающей запятой

Введение тетракода как системы кодирования постбинарных форматов позволяет рассмотреть наряду со стандартными способами и новый способ округления. На рис. 4 приведен способ постбинарного округления. Причем числа, чьи позиции находятся на светлой части вещественной оси (между пунктирными линиями в участках II и III), представляются в постбинарных форматах в виде **интервального числа** с границами соседних базовых точек  $[M_n, M_{n+1}]$ . Такая возможность достигается появлением разрядов множественности (M) и неопределенности (A) в младших разрядах мантиссы [2, с. 155-173].

Таким образом, **постбинарное округление** – округление, оперирующее значением  $\frac{1}{4}$  шага между базовыми точками.

При постбинарном округлении рассматриваются два старших незначащих разряда:

- 00 – число находится в участке I вещественной оси – округление к **меньшей базовой точке** (т.е. к ближайшему числу): отбрасывание незначащих разрядов;
- 01 или 10 – число находится в участках II и III вещественной оси – формирование **нормированного тетракода** [2, с. 158-160] из поля мантииссы, т.е. фактически прибавление к мантииссе **значения множественности**;
- 11 – число находится в участке IV вещественной оси – округление к **большей базовой точке** (т.е. к ближайшему числу): прибавление 1 к полю мантииссы.

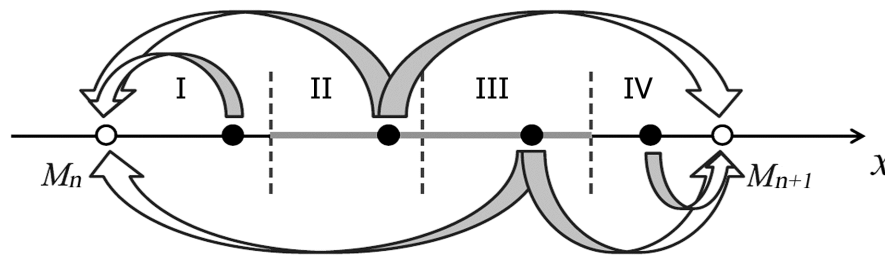


Рисунок 4 — Принцип постбинарного округления чисел

Рассмотрим пример представления числа  $F = 0,9871625$  в формате `rbinary32` с использованием постбинарного округления. Получаем нормализованное двоичное число:

$$F_2 = 1,111110010110110101011\underline{10011111}_2 \times 2^{-1} ..$$

Здесь подчеркнуты незначащие разряды для формата `rbinary32`, среди которых выделены два подлежащих рассмотрению старших разряда.

Таким образом, при постбинарном округлении получаем поле формата

$$\text{rbinary}(F_T) = \underbrace{0}_{\text{знак}} \underbrace{01111110}_{\text{порядок}} \underbrace{111110010110110101}_{\text{мантисса}} \underbrace{MAA}_{\text{идентификатор}} \underbrace{00}_{\text{идентификатор}}$$

Полученный тетракод  $F_T$  сводится к двум двоичным 32-разрядным значениям, являющимися границами интервального числа  $F$  (методы сведения тетракода к интервалу подробно рассмотрены в [2, с. 155-173]):

$$F_T \rightarrow F = [3F7CB6AC, 3F7CB6B0]_h \approx [0,98716235, 0,98716259]_{10},$$

гарантированно содержащее исходное число  $F$ .

Ширина  $d$  полученного интервала равна расстоянию между соседними базовыми точками, то есть, используя (3), получаем:

$$d = h_k = 2^{126-127-21} = 2^{-22} = 2,384185791015625 \times 10^{-7} \approx 2,4 \times 10^{-7}.$$

Это же значение можно получить и как разность границ интервала:

$$d = 0,98716259 - 0,98716235 = 2,4 \times 10^{-7}.$$

Таким образом, при использовании постбинарного округления стало возможным:

1) для чисел, лежащих в областях I и IV вещественной оси, уменьшить ошибку точности представления числа, причем абсолютная максимальная ошибка  $\Delta_k^{\max}$  для этих чисел в постбинарном формате оказалась равной в пределе четверти шага чисел (рис. 5), что в два раза меньше аналогичной абсолютной ошибки  $\Delta_k^{\max}$  для стандарта IEEE 754–2008:

$$\Delta_k^{\max} = \frac{1}{2} \Delta_k^{\max} = \frac{1}{4} h_k = 2^{E_k - \text{off} - m - 2}. \quad (16)$$

2) числа, лежащие в областях II и III вещественной оси, привести к интервалу, ширина которого равна шагу чисел постбинарного формата.

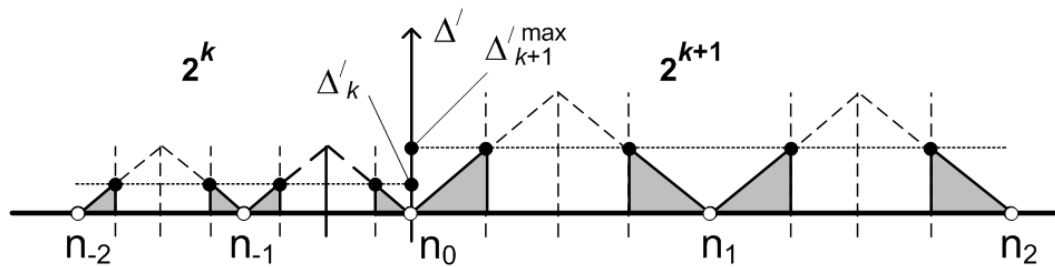


Рисунок 5 – График ошибки точности представления чисел в постбинарном формате

## Выводы

Округление чисел с плавающей запятой – очень важная проблема существующей плавающей арифметики, поскольку числа приходится округлять как при вводе исходных значений, так и после каждой арифметической операции. Переход к постбинарным форматам частично исправляет данную ситуацию, так как в постбинарном представлении числа используется расширенный кодо-логический базис, позволяющий использовать «гибкий» аппарат кодирования чисел, что в свою очередь обеспечивает:

- представление одним значением некоего диапазона чисел на вещественной оси (например, в результате постбинарного округления);
- фиксирование и хранение незначущих (недоопределенных) разрядов мантииссы, приводящее к грамотному (с математической точки зрения) увеличению точности представления чисел;
- организацию плавающей постбинарной арифметики, привносящей новые алгоритмы нормализации чисел, обработки исключительных ситуаций и базовых арифметических операций.

## Литература

1. Курс лекций по теме «Системы искусственного интеллекта» / [сост. Даурова А.А.]. – Владикавказ : Северо-Кавказский горно-металлургический институт, 2008. — [Электронный ресурс]. – Режим доступа : <http://www.skgmi-gtu.ru/aoi/Method/>.
2. Аноприенко А.Я. Постбинарный компьютер и интервальные вычисления в контексте кодо-логической эволюции / А.Я. Аноприенко, С.В. Иваница – Донецк : ДонНТУ, УНИТЕХ, 2011. – 248 с. [Электронный ресурс]. – Режим доступа : <http://ea.donntu.edu.ua/handle/123456789/7544>.
3. Яшкардин В. IEEE 754 – стандарт двоичной арифметики с плавающей точкой / В. Яшкардин. – [Электронный ресурс]. – Режим доступа : <http://softelectro.ru/ieee754.html>.
4. Юровицкий В.М. IEEE754-тика угрожает человечеству. – Москва : МФТИ, РГСУ [Электронный ресурс]. – Режим доступа : <http://www.yur.ru/science/computer/IEEE754.htm>.
5. Аноприенко А.Я. Обобщенный кодо-логический базис в вычислительном моделировании и представлении знаний: эволюция идеи и перспективы развития / А.Я. Аноприенко // Научные труды Донецкого национального технического университета. Серия «Информатика, кибернетика и вычислительная техника» (ИКВТ-2005). – Донецк : ДонНТУ, 2005. – Выпуск 93. – С. 289-316. [Электронный ресурс]. — Режим доступа : <http://ea.donntu.edu.ua/handle/123456789/2780>.
6. Аноприенко А.Я. Гибкая разрядность и постбинарные форматы представления вещественных чисел / А.Я. Аноприенко, С.В. Иваница // Вісник Інженерної академії України. – 2012. – № 1. – С. 92–98.
7. Steve Hollasch. IEEE Standard 754 Floating Point Numbers / Steve Hollasch. – 2004. – [Электронный ресурс]. – Режим доступа : <http://steve.hollasch.net/cgindex/coding/ieeefloat.html>.
8. Аноприенко А.Я. Интервальные вычисления и перспективы их развития в контексте кодо-логической эволюции / А.Я. Аноприенко, С.В. Иваница // Научные труды Донецкого национального технического университета. Серия «Проблемы моделирования и автоматизации проектирования динамических систем» (МАП-2010). – Донецк : ДонНТУ, 2010. – Выпуск 8 (168). – С. 150-160.



## Literatura

1. Kurs lekcij po teme "Sistemy iskusstvennogo intellekta" (sost. Daurova A.A.). Vladikavkaz, Severo-Kavkazskij gorno-metallurgicheskij institut, 2008. <http://www.skgmi-gtu.ru/aoi/Method/>.
2. Anoprienko A.Ja. Postbinarnyj komp'juting i interval'nye vychislenija v kontekste kodo-logicheskoj jevoljucii. Doneck, DonNTU. UNITEH, 2011. 248 s. <http://ea.donntu.edu.ua/handle/123456789/7544>.
3. Jashkardin V. IEEE 754 – standart dvoichnoj arifmetiki s plavajushhej tochkoj. <http://softelectro.ru/ieee754.html>.
4. Jurovickij V.M. IEEE754-tika ugrozhaet chelovechestvu – MFTI, RGSU, Moskva. <http://www.yur.ru/science/computer/IEEE754.htm>.
5. Anoprienko A.Ja. Nauchnye trudy Doneckogo nacional'nogo tehničeskogo universiteta. Serija "Informatika, kibernetika i vychislitel'naja tehnika" (IKVT-2005) vypusk 93. Doneck: DonNTU. 2005. S. 289-316. <http://ea.donntu.edu.ua/handle/123456789/2780>.
6. Anoprienko A.Ja. Visnyk Inženernoi akademii Ukrainy. 2012. № 1. S. 92-98.
7. Steve Hollasch IEEE Standard 754 Floating Point Numbers - 2004. <http://steve.hollasch.net/cgiindex/coding/ieeefloat.html>.
8. Anoprienko A.Ja. Nauchnye trudy Doneckogo nacional'nogo tehničeskogo universiteta. Serija "Problemy modelirovanija i avtomatizacii proektirovanija dinamičeskikh sistem" (MAP-2010). Vypusk 8 (168). Doneck. DonNTU, 2010. S. 150-160.

### РЕЗУМЕ

*S.V. Ivanitsa*

### *Error Estimation of Representation of Real Numbers in the Postbinary Floating-Point Formats*

As the numbers increase for the accuracy of repose, as well as improve the reliability of the calculations, have been proposed postbinary numbers from single to quad precision have been proposed for floating-point formats [2, p. 202].

Coding method used in posbinary formats is based on the principles of code-logical basis [5]: tetracode serves as a code system, and tetrigit serves a storage unit discharge. This "flexibility" can encode quantitative values with a high degree of precision to represent numbers in floating point format.

The introduction of a coding system tetracode postbinary format allows us to consider the standard methods and a new way of rounding, i.e. postbinary rounding, which operates the value of  $\frac{1}{4}$  pitch between base points.

When using postbinary rounding it is possible:

1) Decrease of the error of accuracy of representation, and the absolute maximum error in postbinary format proved to be a step in the limit of a quarter numbers, which is two times less than the same absolute error for the standard IEEE 754–2008.

2) Bringing of the value of the field postbinary format to the interval [8], whose width is equal to the step of postbinary numbers.

*Статья поступила в редакцию 05.06.2012.*