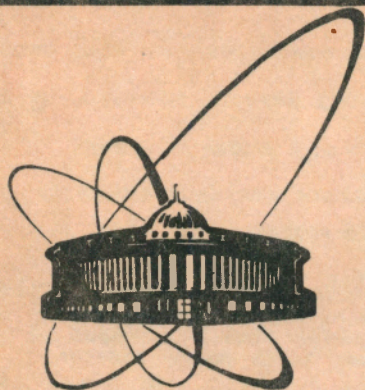


91-187



**СООБЩЕНИЯ
ОБЪЕДИНЕННОГО
ИНСТИТУТА
ЯДЕРНЫХ
ИССЛЕДОВАНИЙ
ДУБНА**

P10-91-187

Е. В. Белякова

**МЕТОД РЕГУЛИРОВАНИЯ ТОЧНОСТИ ДАННЫХ
ДЛЯ УВЕЛИЧЕНИЯ СКОРОСТИ ВЫЧИСЛЕНИЙ**

1991

1. При обработке числовых данных на ЭВМ возникают погрешности, связанные с чисто машинными особенностями.

Работающая в настоящее время в ОИЯИ вычислительная машина VAX-8350 допускает точность 33 десятичных знака после запятой. Однако применение двойной или тем более четверной точности существенно влияет на скорость вычислений.

В настоящей работе предлагается метод, который позволит получать данные с любой промежуточной точностью вплоть до максимальной.

2. На VAX-8350 существуют 4 формата для представления действительных чисел:

- F-format (float) имеет точность 7 десятичных знаков;
- D-format (double precision) - 16 знаков;
- G-format (grant) - 16 знаков;
- H-format (huge) - 33 знака.

3. Рассмотрим задачу о суммировании N чисел. a_1, a_2, \dots, a_N хранятся в памяти в определенном формате, и $\sum_{i=1}^N a_i = S$.

Число S известно и тоже хранится в памяти. Для иллюстрации можно взять задачу суммирования вероятностей, здесь $S = 1$.

Предполагается, что числа a_i после округления до n -го знака после запятой используются в дальнейших вычислениях. Возникает вопрос, как следует произвести округление, чтобы сумма S' округленных чисел a'_i имела бы минимальное отклонение от S ?

Пример: $a = 6,768239$. Пусть $n=3$, тогда $a' = 6,768$.

Введем погрешность округления ζ суммы S : $\sum_{i=1}^N a'_i = S + \zeta$.

Если $\zeta > 0$, то округление выполняется с избытком. Если $\zeta < 0$ - то с недостатком. Случай $\zeta = 0$ практически невозможен.

Введем функцию $\eta_n(x)$: $0 \leq \eta_n(x) \leq 10^{-n}$ для любого $n \geq 1$.

$\eta_n(x)$ можно выразить элементарными функциями:

$\eta_n(x) = x - 10^{-n} * [10^n * x]$, где операция $[]$ означает взятие целой части числа.

Заметим, что функция $\eta_n(x)$ сохраняет знак аргумента.

Например, $\eta_2(6,01785) = 0,00785$.

Положим: $\eta_i = \eta_n(a_i)$ $S_\eta = \sum_{i=1}^N \eta_i$ и для n фиксированного возьмем $c = 10^{-n}/2$.

Если $\eta_i < c$, то округление до n -го знака после запятой выполняется с недостатком, иначе - с избытком.

Пусть $I_{(a>b)}$ - функция-индикатор.

$$I_{(a>b)} = \begin{cases} 0, & a \leq b \\ 1, & a > b. \end{cases}$$

Представим сумму S_η в следующем виде:

$$S_\eta = \sum_{i=1}^N \eta_i * I_{(\eta_i > c)} + \sum_{j=1}^N \eta_j * I_{(\eta_j < c)} = S_\eta^{up} + S_\eta^{down}$$

Пусть $k = \sum_{i=1}^N I_{(\eta_i > c)}$ и $K = \sum_{j=1}^N I_{(\eta_j < c)}$.

Следовательно, справедливы оценки:

$$k * c \leq S_\eta^{up} \leq 10^{-n} * k$$

$$0 \leq S_\eta^{down} < k * c.$$

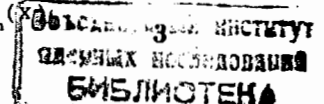
4. Вспомним, что при округлении сумма $\sum_{i=1}^N a_i$ меняется на

величину ζ :

$$\zeta = \sum_{i=1}^N (10^{-n} - \eta_i * I_{(\eta_i > c)}) - \sum_{j=1}^N \eta_j * I_{(\eta_j < c)} =$$

$$= 10^{-n} * N - S_\eta^{up} - S_\eta^{down} = 10^{-n} * N - S_\eta.$$

Итак, $\zeta = 10^{-n} * N - S_\eta$, то есть ζ зависит от S_η и от выбора параметра n функции η_n .



$$\text{Обозначим } L_{\zeta} = \frac{|\zeta - \eta_n(\zeta)|}{10^{-n}} = 10^n * (|\zeta - \eta_n(\zeta)|),$$

L_{ζ} - целочисленная функция.

Число поправок после первого округления, выполняемого "в лоб", можно определить так:

$$L = \begin{cases} L_{\zeta} + 1, & \text{если } \zeta > 0 \text{ и } \eta_n(\zeta) \geq c \text{ и } \eta_n(S) < c, \\ L_{\zeta}, & \text{иначе.} \end{cases}$$

Определим окончательную погрешность $\delta = |S - S''|$, $\delta \leq c = 10^{-n}/2$.

5. Примеры использования алгоритма в задачах суммирования.

Пример 1.

$N=6, n=2$.

a_i	a'_i	a''_i
$a_1=1,1129$	$a'_1=1,11$	$a''_1=1,11$
$a_2=1,1568$	$a'_2=1,16$	$a''_2=1,16$
$a_3=1,8860$	$a'_3=1,89$	$a''_3=1,89$
$a_4=1,0009$	$a'_4=1,00$	$a''_4=1,00$
$a_5=1,0177$	$a'_5=1,02$	$a''_5=1,02$
$a_6=1,1551$	$a'_6=1,16$	$a''_6=1,15$
$S=7,3285$	$S'=7,34$	$S''=7,33$

$$\zeta = 0,0115 > 0, \quad c = 10^{-2} = 0,01$$

$$\eta_2(S) = 0,0085 \geq c, \quad \text{но}$$

$$\eta_2(\zeta) = 0,0015 < c,$$

следовательно,

$$L = L_{\zeta} = \frac{|0,0115 - 0,0015|}{0,01} = 1,$$

то есть достаточно отнять величину 0,01 от любого числа среди a'_i ,

чтобы выполнялось: $\delta = |S - S''| = |-0,0015| < c = 0,005$.

В конкретном примере изменено $a''_6 = 1,15$.

Пример 2.

$N=10, n=2$.

Пусть $S=7,9814$, а $S' = \sum_{i=1}^{10} a'_i = 7,95$.

$\zeta = -0,0314$ - округление выполнено с недостатком.

$$L_{\zeta} = \frac{|-0,0314 - 0,0014|}{0,01} = 3$$

$$L = L_{\zeta} = 3.$$

Следовательно, к первым трем a'_i прибавляем по 0,01. Итак, $S'' = 0,03 + S' = 7,98$ и $\delta = |S - S''| = 0,0014 < 0,05$.

Пример 3.

$N=20, n=2$. Пусть $S=7,2441$, а $S'=7,28$. Тогда $\zeta = 0,0359$.

Находим $L = \frac{0,03}{0,01} + 1 = 4$. Вносим поправки и получаем: $S'' = 7,24$.

6. Заключение

Рассмотренный алгоритм оптимального округления чисел с заведомо заданной погрешностью очень удобен и прост благодаря введению функции $\eta_n(x)$.

Таким образом, используя $\eta_n(x)$ - функцию выделения дробной части числа x , начиная с $(n+1)$ -го десятичного знака, можно работать исключительно с выбранной и удобной точностью, и колоссально увеличить скорость вычислений на компьютерах, так как несравненно быстрее работать с числами, которые имеют ограниченное количество значащих цифр после запятой.

Процедуру с функцией $\eta_n(x)$ целесообразно вставлять после групп промежуточных операций.

Область применения метода достаточно широка - от частных небольших задач до обработки данных с физических установок, в особенности для оптимальной упаковки информации о событиях.

В перспективе это означает экономию памяти ЭВМ и магнитных носителей.

Автор благодарит Кривожижина В.Г. и Кухтина В.В. за ценные замечания о работе.

Рукопись поступила в издательский отдел
23 апреля 1991 года.

Белякова Е.В.

P10-91-187

Метод регулирования точности данных
для увеличения скорости вычислений

Предлагается метод, который позволяет получить данные с любой промежуточной точностью /вплоть до машинной максимальной/ и увеличить тем самым скорость вычислений.

Работа выполнена в Лаборатории сверхвысоких энергий ОИЯИ.

Сообщение Объединенного института ядерных исследований. Дубна 1991

Перевод автора

Belyakova E.V.

P10-91-187

The Method of Accuracy Data Regularization
for the Speed Calculation Increase

The proposed method allows to get the data with any transient accuracy up to max permissible at a concrete computer. So it's possible to increase the speed of calculations.

The investigation has been performed at the Particle Physics Laboratory, JINR.

Communication of the Joint Institute for Nuclear Research. Dubna 1991