

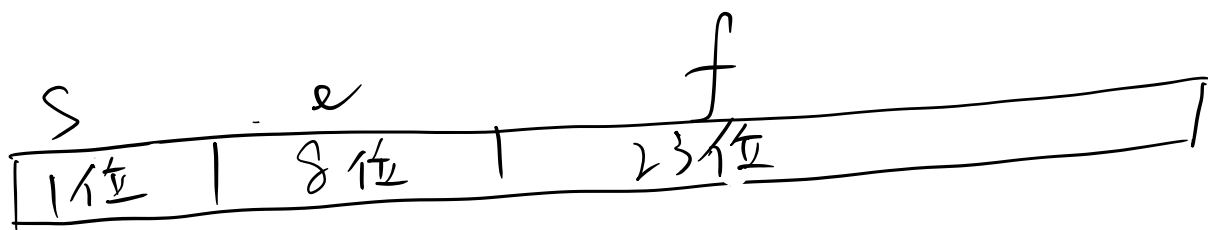
对于小数点后确定只有保留几位的情况，
比如超市 $xxx.xx$ 元，用定点数表示
法中的 BCD 码即可，定点意味着小
数点位置固定，比如始终把十进制
数最后两位视做小数部分，点在
倒数第三位。

BCD 码用 4 位二进制表示一个十进
制数，在准确表达的同时，也牺牲了
可用范围（16 个数只用 10 个），且应

用场易限制较大。

对于小数点后保留位数不确定,或者表示的数范围很大时,定点数法就不适用了。

浮点数的十进制位置不确定,由科学计数法中10的指数来控制,以下实例为证(32位浮点数)



以上表示:

非全数 ~ ~ ~

$$\text{number} = (-1)^s \times 2^e \times 1.f$$

当 $e=0$ 且 $f=0$ 时 $\text{number} \xrightarrow{\text{设定}} 0$

$e=0$ 且 $f \neq 0$ 时 $\text{number} \xrightarrow{\text{设定}} 0.f$

$e=255$ 时 范围无穷大或负无穷大 (看 s)

S 必用一位, 剩下 31 位中, e 用得越多, 可表示范围越大, 有效位越少; 反之 e 越小, 可表示范围越小, 有效位越多。综合考虑, IEEE 标准定为 e 用 8 位, f 用 23 位。

$2^{23} = 8388608$, 十进制 7位有效数

$2^{126} (2^8 \rightarrow -127 \sim 126) = 4 \times 10^{37}$, 约可表示范围

十进制浮点数 \longleftrightarrow 二进制

① $2 \rightarrow 10$

6	5	4	3	2	1	0	-1	-2	-3	-4	-5	
1	0	0	1	0	1	0	.	1	0	0	0	1

$$= 2^6 + 2^3 + 2^1 + 2^{-1} + 2^{-5}$$

$$= 64 + 8 + 2 + 0.5 + 0.03125$$

$$= 74.53125$$

② $10 \rightarrow 2$

1 1 1 0 0 0 1 1

$$85.256 = 1010101.010000110\dots$$

$$\begin{array}{r} 2 \overline{)85} \dots 1 \\ 2 \overline{)42} \dots 0 \\ 2 \overline{)21} \dots 1 \\ 2 \overline{)10} \dots 0 \\ 2 \overline{)5} \dots 1 \\ 2 \overline{)2} \dots 0 \\ 2 \overline{)1} \dots 1 \\ 0 \end{array} \quad \uparrow$$

$$\begin{array}{r} 0.256 \times 2 \dots 0 \\ 0.512 \times 2 \dots 1 \\ 0.024 \times 2 \dots 0 \\ 0.048 \times 2 \dots 0 \\ 0.096 \times 2 \dots 0 \\ 0.192 \times 2 \dots 0 \\ 0.384 \times 2 \dots 0 \\ 0.768 \times 2 \dots 1 \\ 0.536 \times 2 \dots 1 \\ 0.072 \times 2 \dots 0 \end{array} \quad \downarrow$$

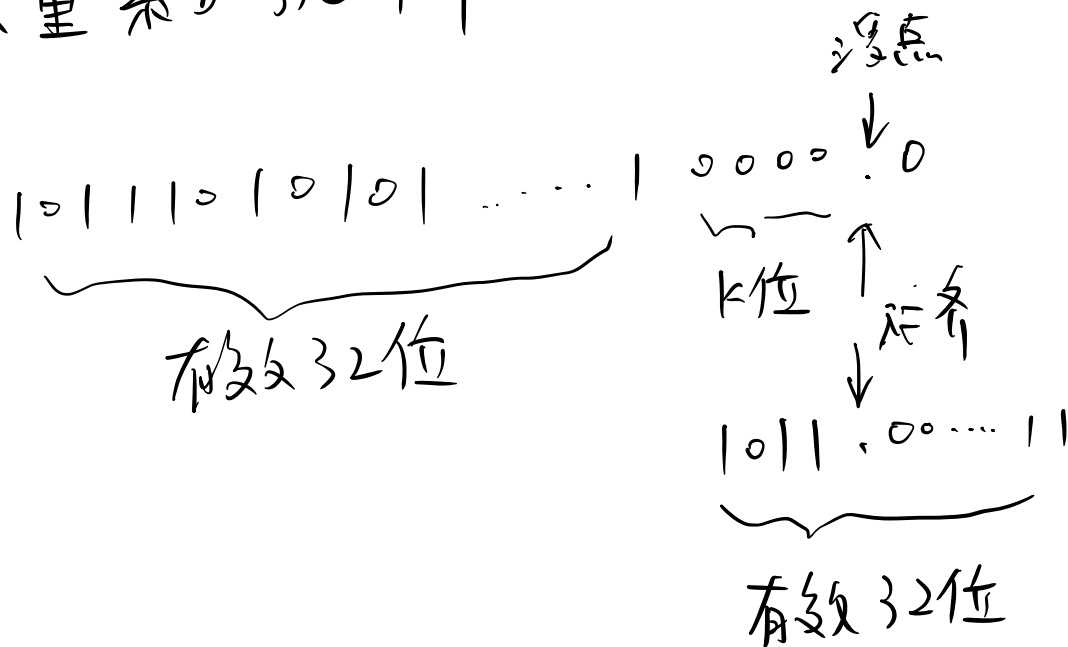
精度损失.

浮点有效位23位时, 进行加法时要
 注意: 某两数量级相差太大, 对

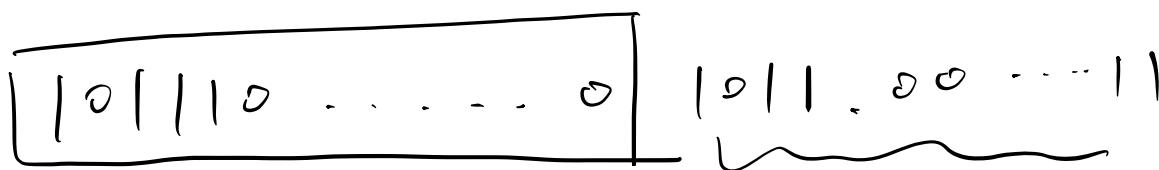
小数点前11, 12, ...

齐右小的一个精度会部份甚至完全丢失,

在大量累加统计中问题很大.



+



只能保留 13 位有效二进制数 ↓

完全丢失

解决方案：

两个已判断小者会损失的精度，提前将损失部份摘出，放在另一个内存地址中，只加上不会损失部份，每次都将摘出部份加在一起，最多统计时再将摘出和加回总和中。