

计算机“二—十”进制数据转换算法*

钟志新

(电子计算机系)

摘要 “二—十”转换是计算机中经常用到的,将二进制浮点数 $0. B_1 B_2 \dots B_n * 2^B$ 转换成十进制数 $0. D_1 D_2 \dots D_m * 10^D$ 的关键是首先求出 D 。本文通过全面研制二进制浮点数的性质,推导出了一个关于精确求 D 的重要结论,并由此构造出一个实际应用时总误差最小的高效求 D 公式。

关键词 数据, 转换, 浮点数, 二进制, 十进制, 变指尾数点, 下确界, 最小值

分类号 TP301.6

数据输出是用户获取程序运行结果的重要手段之一。从系统软件角度看,数据输出就是将数据的内部形式转换成为外部形式。在诸多转换中,使用频度最高、实现难度最大的是将二进制浮点数转换成十进制数,即完成形如: $0. B_1 B_2 \dots B_n * 2^B \Rightarrow 0. D_1 D_2 \dots D_m * 10^D$ 的数据转换。通常,这种转换(记为BTD)需经历如下三个步骤:

(1) 由给定的规格化二进制浮点数 $0. B_1 B_2 \dots B_n * 2^B$ 计算出十进制数的指数 D 。

其中, $0. B_1 B_2 \dots B_n = \sum_{i=1}^n B_i * 2^{-i}$ 称为尾数, $B_i = 0$ 或 1 , 规格化指 $B_1 \neq 0$, n 是硬件能精确表示的位数。显然, $0. B_1 B_2 \dots B_n \in [0.5 \quad 1.0)$, B 是二进制数的阶码。

(2) 接着计算十进制数的规格化小数 $0. D_1 D_2 \dots D_m = 0. B_1 B_2 \dots B_n * 2^B * 10^{-D}$ 。

其中, $0. D_1 D_2 \dots D_m = \sum_{i=1}^m D_i * 10^{-i}$ 称为小数, $D_i = 0 \sim 9$, 规格化指 $D_1 \neq 0$, m 是用户要求的输出位数,显然, $0. D_1 D_2 \dots D_m \in [0.1 \quad 1.0)$ 。

(3) 将 $0. D_1 D_2 \dots D_m$ 和 D 按 ASCII 码以适当形式输出。

可见,第一步求出的 D 是贯穿 BTD 的全过程,特别是对第二步计算小数至关重要,因为小数 $0. D_1 D_2 \dots D_m$ 是通过 $K (\geq 1)$ 次调用“双精度乘法子程序”按算式 $(\dots ((0. B_1 B_2 \dots B_n * 2^B * 10^{D_1}) * 10^{D_2}) * \dots) * 10^{D_K}$ 求出来的,其中, $D_1 + D_2 + \dots + D_K = -D$, 10^{D_i} ($i=1, 2, \dots, K$) 查表取得。如果第一步求出的 D 有误,那么第二步求出的小数就不能满足规格化要求,故需修正,即再调用一次“双精乘”。这显然会增加 BTD 的运行时间,所以,寻找快速可靠的求 D 公式,避免或减少修正,是 BTD 算法的核心问题。

传统 BTD 算法中的求 D 公式是:

* 1992年4月27日收稿

$$D = \lfloor B * \log_{10} 2 \rfloor + 1 \quad (1)$$

式(1)中没有引入尾数,是不是它对求 D 过程没有影响呢?考查 $B=4$ 的情况,由式(1)可求得 $D=2$ 。但分析表明,对于 $B=4$,存在一个特殊的尾数点 $t_B=0.625$,若给定的尾数 $\geq t_B$,式(1)总是对的;若给定的尾数 $< t_B$,式(1)总是错的(计算值比实际值大1)。可见, t_B 在求 D 过程中有“门限”的作用。

考查规格化二进制浮点数 $0.B_1B_2\cdots B_n * 2^B$,其值域为 $[2^{B-1} \quad 2^B)$ 。若令 $2^{B-1} = 0.D'_1D'_2\cdots D'_m * 10^{D'}$, $2^B = 0.D_1D_2\cdots D_m * 10^D$,且 D'_1 和 D_1 均不为0,显然, D' 与 D 可能相等,也可能不相等(差1)。由此引出了二个定义:

定义1 对任意给定的规格化二进制浮点数,若阶码 B 使得 2^{B-1} 和 2^B 用规格化十进制数表示时,有指数 $D'=D$,则称 B 是稳定阶;若有指数 $D'+1=D$,则称 B 是不稳定阶。

根据定义, $B=4$ 是一个不稳定阶。

定义2 设 B 是不稳定阶,称 t_B 是关于 B 的变指尾数点。如果尾数 $\geq t_B$ 时,有 $0.B_1B_2\cdots B_n * 2^B = 0.D_1D_2\cdots D_m * 10^D$;如果尾数 $< t_B$ 时,有 $0.B_1B_2\cdots B_n * 2^B = 0.D'_1D'_2\cdots D'_m * 10^{D-1}$ 。

根据定义,对所有不稳定阶 B ,因尾数的不同,十进制数的指数 D 有两个可能的取值,但式(1)忽略尾数在求 D 过程中的影响,只能求出其中之一,因而正确性受到损失。

我们关心的是:

(1)对于一个给定的计算机系统,在二进制浮点数阶码的取值范围内究竟有多少个是不稳定阶?这个数可以帮助我们了解式(1)出错的严重程度。

(2)对于一个给定的二进制浮点数,如何断定其阶码是不稳定阶?

(3)对所有不稳定阶,如何求出相应的变指尾数点?

下面我们将逐一回答这些问题,并推导出一个理论上的精确求 D 公式。

1 精确求 D 公式

对于一个给定的计算机系统,从二进制浮点数阶码 B 的取值范围,可以导出十进制数指数 D 的取值范围。下面的定理将告诉我们: D 有多少个可能的取值, B 的取值范围中就会有多个不稳定阶。

引理 对任意给定的指数 D ,存在唯一一组阶码 B_D 和相应的规格化小数 $0.D_1D_2\cdots D_m$,使得 $2^{B_D} = 0.D_1D_2\cdots D_m * 10^D$ 。

证明 记 $\text{Dom}(B_D) = \left[\frac{D-1}{\log_{10} 2}, \frac{D}{\log_{10} 2} \right)$,它有如下性质:

(1) $\text{Dom}(B_D)$ 与 D 是一一对应的;

(2) $\text{Dom}(B_D)$ 中的元素(即阶码值)个数至少3个,至多4个;

(3)任意两个不同的 Dom ,交集为空,因而任意阶码 B 只能属于某一个 Dom 。上述性质不难验证,故证明略去。

任取 $B \in \text{Dom}(B_D)$,因为

$$\frac{D-1}{\log_{10} 2} \leq B < \frac{D}{\log_{10} 2}$$

故

$$D - 1 \leq B * \log_{10} 2 < D$$

设 $2^B = X$, 取对数得 $B * \log_{10} 2 = \log_{10} X$

从而有: $D - 1 \leq \log_{10} X < D$, 即有: $10^{D-1} \leq X < 10^D$

故可令 $X = 0. D_1 D_2 \cdots D_m * 10^D$, 且 $D_1 \neq 0$

即有: $2^B = 0. D_1 D_2 \cdots D_m * 10^D$, 且 $D_1 \neq 0$

(2)

由 B 的取法和 $\text{Dom}(B_D)$ 的性质知引理得证。

定理 1 对于任意给定的指数 D , 存在唯一一个不稳定阶 B 。

证明 令 $B = \min\{B' \mid B' \in \text{Dom}(B_D)\}$

$$B = \begin{cases} 0 & \text{若 } D=1 \\ \lfloor \frac{D-1}{\log_{10} 2} \rfloor + 1 & \text{若 } D \neq 1 \end{cases}$$

对于 $D=1$ 的情况, $B=0$ 是不稳定阶容易检证。

假设 $D \neq 1$, 由引理知有:

$$2^B = 2^{\lfloor \frac{D-1}{\log_{10} 2} \rfloor + 1} = 0. D_1 D_2 \cdots D_m * 10^D, \text{ 且 } D_1 \neq 0$$

下面证有 $2^{B-1} = 0. D'_1 D'_2 \cdots D'_m * 10^{D-1}$, 且 $D'_1 \neq 0$, 我们从考查 $0. D_1 D_2 \cdots D_m$ 的大小着手:

$$\text{因 } \lfloor \frac{D-1}{\log_{10} 2} \rfloor < \frac{D-1}{\log_{10} 2}$$

$$\text{故 } 0. D_1 D_2 \cdots D_m * 10^D < 2^{\frac{D-1}{\log_{10} 2} + 1} = 2^{\frac{D-1 + \log_{10} 2}{\log_{10} 2}}$$

$$\text{故有 } D + \log_{10} 0. D_1 D_2 \cdots D_m < D - 1 + \log_{10} 2$$

$$\text{即 } \log_{10} 0. D_1 D_2 \cdots D_m < \log_{10} (2 * 10^{-1})$$

$$\text{故有 } 0. 1 \leq 0. D_1 D_2 \cdots D_m < 0. 2$$

(3)

$$\text{由此可得 } 2^{B-1} = \frac{2^B}{2} = \frac{0. D_1 D_2 \cdots D_m}{2} * 10^D = 0. 0 D'_1 D'_2 \cdots D'_m * 10^D = 0. D'_1 D'_2 \cdots D'_m$$

$* 10^{D-1}$, 且 $D'_1 \neq 0$

所以, $B \lfloor \frac{D-1}{\log_{10} 2} \rfloor + 1$ ($D \neq 1$) 是不稳定阶。

由 B 的取法知, B 是唯一的, 定理 1 得证。

根据定理 1 的结论和 $\text{Dom}(B_D)$ 的大小知, 在 B 的取值范围内, 有近 1/3 的阶码可能导致式 (1) 计算出错。因此, 寻找更为可靠的求 D 公式是必要的。下面的定理为这个努力提供了线索。

定理 2 对任意给定的阶码 B , B 是不稳定阶当且仅当 $\lfloor B * \log_{10} 2 \rfloor = \lfloor (B-1) * \log_{10} 2 \rfloor + 1$; B 是稳定阶当且仅当 $\lfloor B * \log_{10} 2 \rfloor = \lfloor (B-1) * \log_{10} 2 \rfloor$ 。

证明 设 B 是不稳定阶, 由定义有: $2^B = 0. D_1 D_2 \cdots D_m * 10^D$, $2^{B-1} = 0. D'_1 D'_2 \cdots D'_m * 10^{D-1}$, 且 D_1, D'_1 不为 0 取对数再取下确界得:

$$\lfloor B * \log_{10} 2 \rfloor = D - 1, \lfloor (B - 1) * \log_{10} 2 \rfloor = D - 2 \quad (4)$$

显然有: $\lfloor B * \log_{10} 2 \rfloor = \lfloor (B-1) * \log_{10} 2 \rfloor + 1$

必要条件得证。

下面证充分条件: 设 $\lfloor B * \log_{10} 2 \rfloor = \lfloor (B-1) * \log_{10} 2 \rfloor + 1$, 由引理中的式 (2)

中, 可以假定:

$$2^B = 0. D_1 D_2 \cdots D_m * 10^D \quad (5)$$

从而有: $\lfloor B * \log_{10} 2 \rfloor = D - 1$

由充分性假设, 有: $\lfloor (B - 1) * \log_{10} 2 \rfloor = D - 2$

即 $(B - 1) * \log_{10} 2 = D - 1 + \Delta$, 其中 $-1 \leq \Delta < 0$

令: $\Delta = \log_{10} 0. D'_1 D'_2 \cdots D'_m$, 且 $D'_1 \neq 0$

$$\text{可得 } 2^{B-1} = 0. D'_1 D'_2 \cdots D'_m * 10^{D-1} \quad (6)$$

由式 (5) 和 (6) 知, B 是不稳定阶, 充分条件得证。

B 是稳定阶的充要条件证明类似。(证略)

从定理 2 证明过程中导出的式 (4) 可以看到: $\lfloor B * \log_{10} 2 \rfloor + 1 = D$, $\lfloor (B - 1) * \log_{10} 2 \rfloor + 1 = D - 1$, 它们有相同的计算模式, 若用 B 参与运算可求得 D , 若用 $B - 1$ 参与运算可求得 $D - 1$, 这使我们设想: 能否引入一种控制, 使得式 (1) 在求 D 时, 可根据所给二进制浮点数尾数值的不同, 正确选择 B 或 $B - 1$ 参与运算, 从而获得正确的指数 D 。下面的定理找到了这种控制。

定理 3 若 B 是不稳定阶, 则关于 B 的变指尾数点 $t_B = 10^{\lfloor B * \log_{10} 2 \rfloor} / 2^B$ 。

证明 因 B 是不稳定阶, 易知有: $2^{B-1} < 10^{D-1} \leq 2^B$

即有: $2^{-1} < 10^{D-1} / 2^B \leq 1$

由定理 2 与式 (4) 知: $t_B = 10^{D-1} / 2^B$

因为对所有小于 t_B 的尾数有:

$$0. B_1 B_2 \cdots B_n * 2^B < t_B * 2^B = 10^{D-1}$$

且 $0. B_1 B_2 \cdots B_n * 2^B = 0. B_1 B_2 \cdots B_n * 0. D_1 D_2 \cdots D_m * 10^D \geq 0.5 * 0.1 * 10^D = 0.5 * 10^{D-1}$

故有 $0. B_1 B_2 \cdots B_n * 2^B = 0. D'_1 D'_2 \cdots D'_m * 10^{D-1}$, 且 $D'_1 (\geq 5) \neq 0$

对所有大于等于 t_B 的尾数有:

$$0. B_1 B_2 \cdots B_n 2^B \geq t_B * 2^B = 10^{D-1} = 0.1 * 10^D$$

且由定理 1 式 (3) 知, $0. B_1 B_2 \cdots B_n * 2^B < 2^B < 0.2 * 10^D$

故有 $0. B_1 B_2 \cdots B_n * 2^B = 0. D_1 D_2 \cdots D_m * 10^D$, 且 $D_1 (=1) \neq 0$

所以, t_B 是关于 B 的变指尾数点。定理 3 得证。

若事先求出全部的 t_B 并存于一张表, 则可构造如下求 D 公式:

$$D = \lfloor \text{INT}(B. B_1 B_2 \cdots B_n - t_B) * \log_{10} 2 \rfloor + 1 \quad (8)$$

其中, $B. B_1 B_2 \cdots B_n$ 是分离二进制浮点数的阶码与尾数后组合而成的定点数; INT 是取整操作。

显然, 对于不稳定阶 B , 式 (8) 能够根据尾数 $0. B_1 B_2 \cdots B_n$ 与 t_B 的大小关系正确选择 B 或 $B - 1$ 参与运算, 得到正确结果。对于稳定阶, 式 (8) 未加判断地也做了同样的选择。由定理 2 知, 此举不影响正确性。因此, 若不考虑硬件精度位数限制带来的误差, 式 (8) 是一个精确求 D 公式。

下一节我们将从工程的角度评价式 (8), 进而给出一个可实际应用的最小总误差求 D 公式。

2 最小总误差求 D 公式

将式 (8) 应用于工程实际时, 有如下三个问题: (1) 受机器字长限制, 二进制浮点数的尾数和 t_B 值存在误差。因此, 式 (8) 的计算结果仍含有误差; (2) t_B 的个数由定理 1 知是一个不小的数, 将 t_B 存于表中会占用一大片主存单元; (3) 为得到 t_B 需查表。这会严重地影响算法效率, 因此, 要找到一个空间比式 (8) 小, 正确性比式 (1) 高的实用求 D 公式。

方法之一是取一常数 t 来代替表中的全部 t_B , 即令:

$$D = \lfloor \text{INT}(B_1 B_2 \dots B_n - t) * \log_{10} 2 \rfloor + 1 \quad (9)$$

显然它没有空间占用, 但 t 取何值才能使式 (9) 具有最小误差呢?

用 t 代替各个 t_B , 出错的尾数区间是 $[t_B t]$ 或 $[t t_B]$ 。若将所有的 t_B 按大小排序并赋一下标, 可得: $t_{B1} < t_{B2} < \dots < t_{BN}$, 其中 N 为变指尾数点的最大个数, $t_{B_N} = 1.0, t_{B_1} > 0$ 。

5, 则 $f(t) = \sum_{i=1}^N |t - t_{B_i}|$, $t \in [t_{B_1} t_{B_N}]$ 表示用 t 代替各个 t_B 时, 出错区间的总和。下面求解使 $f(t)$ 最小的 t 值。

该问题更一般的描述是:

已知: $0.5 = t_0 < t_1 < \dots < t_m < t_{m+1} < \dots < t_N = 1.0$, N 为正整数, 函数 $f(t) = \sum_{i=0}^N |t - t_i|$, $t \in [t_0 t_N]$, 试求使 $f(t)$ 在 $[t_0 t_N]$ 上取小值的点 t' 。

解 函数 $f(t)$ 在 $[t_0 t_N]$ 上是连续的。(易证, 略)

由连续函数在闭区间上的性质知, 在 $[t_0 t_N]$ 上至少存在一点 t' 使 $f(t')$ 是函数 $f(t)$ 在该区间上的最小值。

将闭区间 $[t_0 t_N]$ 分割成 N 个子区间 $[t_m t_{m+1}]$, $m = 0, 1, \dots, N-1$ 。设 $t \in [t_m t_{m+1}]$, 则有

$$\begin{aligned} f(t) &= \sum_{i=0}^N |t - t_i| = \sum_{i=0}^m (t - t_i) + \sum_{i=m+1}^N (t_i - t) \\ &= (2m - N + 1)t + C, C = \sum_{i=m+1}^N t_i - \sum_{i=0}^m t_i \end{aligned}$$

可见, 函数 $f(t)$ 在各子区间上表现为线性函数, 故其极小值必在端点处取得, 所以使函数 $f(t)$ 在 $[t_0 t_N]$ 上取最小值的点必在 t_i 中找到。

因为对于所有的 $m < \lceil \frac{N-1}{2} \rceil$, 函数 $f(t)$ 在 $[t_m t_{m+1}]$ 上均是递减的, 对于所有的 $m > \lceil \frac{N-1}{2} \rceil$, $f(t)$ 在 $[t_m t_{m+1}]$ 上均是递增的; 对于 $m = \lceil \frac{N-1}{2} \rceil$, 若 N 为偶数, $f(t)$ 在 $[t_m t_{m+1}]$ 上是递增的, 若 N 为奇数, $f(t)$ 在 $[t_m t_{m+1}]$ 上恒取 C , 所以 $t' = t_{\lceil \frac{N-1}{2} \rceil}$ 必使函数 $f(t)$ 在 $[t_0 t_N]$ 上取得最小值。(解毕)

式 (1) 是式 (9) 的一个特例, 即取 $t' = 0.5$ 。由函数 $f(t)$ 在区间 $[t_0 t_{\lceil \frac{N-1}{2} \rceil}]$ 上的递减性知: $f(0.5) > f(t_{\lceil \frac{N-1}{2} \rceil})$, 就是说式 (9) 的误差区间总和比式 (1) 要小, 因此正确性更高。

3 结束语

式(9)中取 $t = t_{\lceil \frac{N-1}{2} \rceil}$ 是基于公式可能遇到系统允许的全部二进制浮点数而推导出来的。如果一个计算机系统是应用于某个专门领域的,公式遇到的浮点数中, t_B 有聚类性,例如偏向基本区间 $[0.5 \quad 1.0)$ 的某个部分,那么 t 的取值可以在一个更小的范围考虑,方法与结论形式均同上,但公式更具针对性,效率也将更高。

本文的研究曾得到李晓梅教授、蔡放讲师的帮助,在此表示感谢。

参 考 文 献

- 1 国防科技大学六系. YH10 运行系统设计说明. 1983
- 2 国防科技大学六系. YFT77 设计说明书. 1988

The Research on BTD Algorithm of Data Conversion in Computer

Zhong Zhixin

(Department of Computer science)

Abstract

The binary-to-decimal (BTD) conversion is often used in computer. It is the key problem to calculate D first for converting the binary floating point number (BFPN) $0. B_1 B_2 \cdots B_n \cdot 2^b$ into the decimal $0. D_1 D_2 \cdots D_m \cdot 10^p$. In this paper, we comprehensively studied the nature of BFPN and inferred an important conclusion for calculating D precisely. Based on the conclusion we constructed a high efficient formula which can not only calculate D but also keep the total error minimum in application.

Key words data, conversion, floating point number, binary system, decimal system.