

情况 1: 规格化的值

这是最普遍的情况。当 exp 的位模式既不全为 0 (数值 0), 也不全为 1 (单精度数值为 255, 双精度数值为 2047) 时, 都属于这类情况。在这种情况下, 阶码字段被解释为以偏置 (biased) 形式表示的有符号整数。也就是说, 阶码的值是 $E = e - \text{Bias}$, 其中 e 是无符号数, 其位表示为 $e_{k-1} \cdots e_1 e_0$, 而 Bias 是一个等于 $2^{k-1} - 1$ (单精度是 127, 双精度是 1023) 的偏置值。由此产生指数的取值范围, 对于单精度是 $-126 \sim +127$, 而对于双精度是 $-1022 \sim +1023$ 。

小数字段 frac 被解释为描述小数 f , 其中 $0 \leq f < 1$, 其二进制表示为 $0.f_{n-1} \cdots f_1 f_0$, 也就是二进制小数点在最高有效位的左边。尾数定义为 $M = 1 + f$ 。有时, 这种方式也叫做隐含的以 1 开头的 (implied leading 1) 表示, 因为我们可以把 M 看成一个二进制表达式为 $1.f_{n-1} f_{n-2} \cdots f_0$ 的数字。既然我们总是能够调整阶码 E , 使得尾数 M 在范围 $1 \leq M < 2$ 之中 (假设没有溢出), 那么这种表示方法是一种轻松获得一个额外精度位的技巧。既然第一位总是等于 1, 那么我们就不需要显式地表示它。

情况 2: 非规格化的值

当阶码域为全 0 时, 所表示的数是非规格化形式。在这种情况下, 阶码值是 $E = 1 - \text{Bias}$, 而尾数的值是 $M = f$, 也就是小数字段的值, 不包含隐含的开头的 1。

旁注 对于非规格化值为什么要这样设置偏置值

使阶码值为 $1 - \text{Bias}$ 而不是简单的 $-\text{Bias}$ 似乎是违反直觉的。我们将很快看到, 这种方式提供了一种从非规格化值平滑转换到规格化值的方法。

非规格化数有两个用途。首先, 它们提供了一种表示数值 0 的方法, 因为使用规格化数, 我们必须总是使 $M \geq 1$, 因此我们就不能表示 0。实际上, $+0.0$ 的浮点表示的位模式为全 0: 符号位是 0, 阶码字段全为 0 (表明是一个非规格化值), 而小数域也全为 0, 这就得到 $M = f = 0$ 。令人奇怪的是, 当符号位为 1, 而其他域全为 0 时, 我们得到值 -0.0 。根据 IEEE 的浮点格式, 值 $+0.0$ 和 -0.0 在某些方面被认为是不同的, 而在其他方面是相同的。

非规格化数的另外一个功能是表示那些非常接近于 0.0 的数。它们提供了一种属性, 称为逐渐溢出 (gradual underflow), 其中, 可能的数值分布均匀地接近于 0.0。

情况 3: 特殊值

最后一类数值是指阶码全为 1 的时候出现的。当小数域全为 0 时, 得到的值表示无穷, 当 $s=0$ 时是 $+\infty$, 或者当 $s=1$ 时是 $-\infty$ 。当我们把两个非常大的数相乘, 或者除以零时, 无穷能够表示溢出的结果。当小数域为非零时, 结果值被称为 “NaN”, 即 “不是一个数 (Not a Number)” 的缩写。一些运算的结果不能是实数或无穷, 就会返回这样的 NaN 值, 比如当计算 $\sqrt{-1}$ 或 $\infty - \infty$ 时。在某些应用中, 表示未初始化的数据时, 它们也很有用处。

2.4.3 数字示例

图 2-34 展示了一组数值, 它们可以用假定的 6 位格式来表示, 有 $k=3$ 的阶码位和 $n=2$ 的尾数位。偏置量是 $2^{3-1} - 1 = 3$ 。图中的 a 部分显示了所有可表示的值 (除了 NaN)。两个无穷值在两个末端。最大数量值的规格化数是 ± 14 。非规格化数聚集在 0 的附近。图的 b 部分中, 我们只展示了介于 -1.0 和 $+1.0$ 之间的数值, 这样就能够看得更加清楚了。两个零是特殊的非规格化数。可以观察到, 那些可表示的数并不是均匀分布的——越靠近原点处它们越稠密。