# 15 Apéndice C: Aritmética de Punto Flotante

La aritmética de punto flotante es considerada un tema esotérico para muchas personas. Esto es sorprendente porque el punto flotante es omnipresente en los sistemas informáticos. Casi todos los lenguajes de programación tienen un tipo de datos de punto flotante, i.e. los números no enteros como 1.2 ó 1e+45.

Los problemas de precisión del punto flotante son una preocupación crítica en la informática y la computación numérica, donde la representación y manipulación de números reales puede conducir a resultados inesperados y erróneos. Estos obstáculos surgen debido alas limitaciones inherentes de la aritmética de punto flotante, que aproxima números reales con una precisión finita.

Esto puede causar problemas importantes en diversas aplicaciones, desde simulaciones científicas hasta cálculos financieros, donde incluso las imprecisiones menores pueden propagarse y amplificarse, dando lugar a errores sustanciales. Comprender estos obstáculos es esencial para que los desarrolladores, ingenieros y científicos implementen algoritmos numéricos sólidos y confiables, garantizando que las matemáticas realizadas por las computadoras sigan siendo confiables y precisas.

Comprensión de los Errores de Redondeo en Aritmética de Punto Flotante El meollo de la cuestión es la representación de números de punto flotante. Las computadoras utilizan una cantidad finita de bits para almacenar estos números, y generalmente cumplen con el estándar IEEE 754<sup>137</sup>. Este estándar define cómo se almacenan los números en formato binario, con un número fijo de bits asignados para el signo, el exponente y la mantisa. Si bien esto permite representar una amplia gama de valores, también impone limitaciones a la precisión. No todos los números decimales se pueden representar exactamente en forma binaria, lo que genera pequeñas discrepancias conocidas como errores de redondeo.

Una de las primeras cosas que uno encuentra con sorpresa cuando hace cálculos en una computadora es que por ejemplo, si usamos números muy grandes y seguimos incrementando su valor eventualmente el resultado será

<sup>&</sup>lt;sup>137</sup>El estándar IEEE 754-2008 define varios tamaños de números de punto flotante: media precisión (binary16), precisión simple (binary32), precisión doble (binary64), precisión cuádruple (binary128), etc., cada uno con su propia especificación.

negativo ... ¿qué pasó? Esto se llama desbordamiento aritmético al intentar crear un valor numérico que está fuera del rango que puede representarse con un número dado de dígitos, ya sea mayor que el máximo o menor que el mínimo valor representable. Algo similar pasa al restar al menor número representable en la máquina, el resultado será positivo y se denomina subdesbordamiento.

Por otro lado, tenemos errores de redondeo, estos ocurren porque el sistema binario no puede representar con precisión ciertas fracciones. Por ejemplo, el número decimal 0.1 no se puede representar exactamente en binario, lo que da como resultado una aproximación, por ejemplo al sumar 0.1 + 0.2 en una computadora usando por ejemplo el lenguaje Python obtenemos:

```
\begin{array}{l} \text{print}(0.1+0.2) \\ 0.30000000000000000004 \end{array}
```

que no es exactamente lo que esperábamos<sup>138</sup>. Cuando se utilizan tales aproximaciones en los cálculos, los errores pueden acumularse. Este fenómeno es particularmente problemático en procesos iterativos, donde el mismo cálculo se realiza repetidamente y los errores se acumulan con el tiempo.

Estos errores de precisión se vuelven particularmente problemáticos cuando se realizan controles de igualdad. En un mundo ideal, comparar la igualdad de dos números de punto flotante sería sencillo. Sin embargo, debido a los pequeños errores introducidos durante las operaciones aritméticas, dos números que deberían ser iguales pueden no ser exactamente iguales en su representación binaria.

Por ejemplo, el resultado de sumar 0.1 y 0.2 puede no ser exactamente igual a 0.3 debido a la forma en que estos números se representan en binario, por ejemplo:

```
 \begin{array}{c} \text{if } 0.1 + 0.2 == 0.3; \\ \text{print("si")} \\ \text{else:} \\ \text{print("no")} \end{array}
```

```
\begin{array}{c} \text{print}(0.1\ *\ 3) \\ 0.30000000000000000004 \end{array}
```

 $<sup>^{138}\</sup>mathrm{Lo}$ mismo obtenemos si usamos la multiplicación, por ejemplo:

la respuesta será "no". Esta discrepancia puede provocar que fallen las comprobaciones de igualdad, lo que provocará un comportamiento inesperado en los programas.

Para mitigar estos problemas, los desarrolladores suelen utilizar una técnica conocida como "comparación épsilon". En lugar de verificar la igualdad exacta, verifican si la diferencia absoluta entre dos números de punto flotante es menor que un valor pequeño predefinido, conocido como épsilon. Este enfoque reconoce la imprecisión inherente de la aritmética de punto flotante y proporciona una forma más sólida de comparar números. Sin embargo, elegir un valor épsilon apropiado puede resultar complicado, ya que depende del contexto específico y del rango de valores involucrados.

El problema se ve exacerbado por la precisión finita de los números de punto flotante. Al realizar operaciones aritméticas, el resultado a menudo debe redondearse para que se ajuste a los bits disponibles. Este redondeo puede introducir más imprecisiones. Por ejemplo, sumar dos números de punto flotante de magnitudes muy diferentes puede provocar una pérdida de precisión, ya que el número más pequeño puede ignorarse de hecho. Esto se conoce como cancelación catastrófica y puede afectar gravemente a la precisión de los cálculos.

Por ejemplo, ¿qué podría tener de interesante la humilde fórmula cuadrática?. Después de todo, es una fórmula. Simplemente le pones números. Bueno, hay un detalle interesante. Cuando el coeficiente lineal b es grande en relación con los otros coeficientes, la fórmula cuadrática puede dar resultados incorrectos cuando se implementa en aritmética de punto flotante. Eso es cierto, pero veamos qué sucede cuando tenemos a = c = 1 y b = 10e8 en

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \text{ y } x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$$
 (15.1)

regresa

$$x_1 = -7.450580596923828e - 09 \text{ y } x_2 = -1000000000.0$$

La primera raíz está equivocada en aproximadamente un 25%, aunque la segunda es correcta.

¿Qué pasó? La ecuación cuadrática violó la regla cardinal del análisis numérico: evitar restar números casi iguales. Cuanto más similares sean dos números, más precisión puedes perder al restarlos. En este caso  $(b^2 - 4ac)$  es casi igual a b. Si evaluamos, obtenemos 1.49e - 8 cuando sería la respuesta correcta 2.0e - 8.

antoniocarrillo@ciencias.unam.mx 715 Antonio Carrillo Ledesma, Et alii

Si usamos la otra formula<sup>139</sup>

$$x_1 = \frac{-2c}{b + \sqrt{b^2 - 4ac}}$$
 y  $x_2 = \frac{-2c}{b - \sqrt{b^2 - 4ac}}$  (15.2)

obtenemos

$$x_1 = -1e - 08 \text{ y } x_2 = -134217728.0$$

Entonces, ¿cuál formula cuadrática es mejor? Da la respuesta correcta para la primera raíz, exacta dentro de la precisión de la máquina. Pero ahora la segunda raíz está equivocada en un 34%. ¿Por qué la segunda raíz es incorrecta? La misma razón que antes: ¡restamos dos números casi iguales!

La versión familiar de la fórmula cuadrática calcula correctamente la raíz más grande y la otra versión calcula correctamente la raíz más pequeña. Ninguna versión es mejor en general. No estaríamos ni mejor ni peor usando siempre la nueva fórmula cuadrática que la anterior. Cada uno es mejor cuando evita restar números casi iguales. La solución es utilizar ambas fórmulas cuadráticas, utilizando la apropiada para la raíz que estás intentando calcular.

Otro error común es la suposición de que la aritmética de punto flotante es asociativa y distributiva, como lo es en matemáticas puras. En realidad, el orden de las operaciones puede afectar significativamente el resultado debido a errores de redondeo. Por ejemplo, la expresión (a+b) + c puede producir un resultado diferente que a + (b + c) cuando se trata de números de punto flotante. Esta no asociatividad puede provocar errores sutiles, especialmente en algoritmos complejos que se basan en cálculos precisos.

Además, la aritmética de punto flotante puede introducir problemas en los algoritmos que requieren comparaciones exactas, como los que se utilizan en la clasificación o la búsqueda. Cuando los números de punto flotante se utilizan como claves en estructuras de datos como tablas hash o árboles de búsqueda binarios, los errores de precisión pueden provocar un comportamiento incorrecto, como no encontrar un elemento existente o insertar incorrectamente un duplicado. Para solucionar este problema, es posible que los desarrolladores necesiten implementar funciones de comparación personalizadas que tengan en cuenta la imprecisión del punto flotante.

Pasando a implicaciones prácticas, estos errores de redondeo pueden tener consecuencias de gran alcance. En las simulaciones científicas, pequeñas im-

 $<sup>\</sup>overline{\phantom{a}}^{139}$ Para obtener la segunda fórmula, multiplique el numerador y denominador de  $x_1$  por  $-b - \sqrt{b^2 - 4ac}$  (y de manera similar para  $x_2$ ).

precisiones pueden dar lugar a predicciones incorrectas, lo que podría socavar la validez de la investigación. En aplicaciones financieras, los errores de redondeo pueden dar lugar a importantes discrepancias monetarias, afectando todo, desde el análisis del mercado de valores hasta las transacciones bancarias. Incluso en aplicaciones cotidianas como la representación de gráficos, los errores de redondeo pueden causar artefactos visuales que restan valor a la experiencia del usuario.

Además, comprender las limitaciones de la aritmética de punto flotante puede ayudar a diseñar sistemas más robustos. Al anticipar dónde es probable que se produzcan errores de redondeo, los desarrolladores pueden implementar controles y equilibrios para detectar y corregir imprecisiones. Por ejemplo, la aritmética de intervalos puede proporcionar límites a los posibles valores de un cálculo, ofreciendo una manera de cuantificar la incertidumbre introducida por los errores de redondeo.

Estrategias para Mitigar los Problemas de Precisión del Punto Flotante en el Desarrollo de Software Los problemas de precisión del punto flotante son un desafío común en el desarrollo de Software y a menudo conducen a resultados inesperados y errores sutiles. Estos problemas surgen debido a las limitaciones inherentes a la representación de números reales en formato binario, lo que puede provocar errores de redondeo y pérdida de precisión. Para mitigar estos obstáculos, los desarrolladores deben emplear una variedad de estrategias que garanticen la precisión numérica y la confiabilidad en sus aplicaciones.

Una estrategia eficaz es utilizar tipos de datos de mayor precisión cuando sea necesario. Si bien los tipos de punto flotante estándar como "flotante" y "doble" son suficientes para muchas aplicaciones, es posible que no proporcionen la precisión requerida para cálculos más sensibles. En tales casos, el uso de tipos de precisión extendidos, como "long double" en C++ o bibliotecas de precisión arbitraria como GMP (Biblioteca aritmética de precisión múltiple GNU), puede reducir significativamente el riesgo de pérdida de precisión.

Sin embargo, es importante equilibrar la necesidad de precisión con consideraciones de rendimiento, ya que los tipos de mayor precisión pueden resultar costosos desde el punto de vista computacional. Otro enfoque consiste en implementar algoritmos numéricos que sean inherentemente más estables. Algunos algoritmos son más susceptibles a errores de redondeo que otros y elegir el algoritmo correcto puede marcar una diferencia significativa.

Por ejemplo, al resolver sistemas de ecuaciones lineales, utilizar métodos como la descomposición LU o la descomposición QR puede ser más estable que la eliminación gaussiana. Además, se pueden emplear técnicas de refinamiento iterativo para mejorar la precisión de la solución corrigiendo iterativamente los errores introducidos por la aritmética de punto flotante.

Los desarrolladores deben tener en cuenta el orden de las operaciones en sus cálculos. Las propiedades asociativas y distributivas de la aritmética no siempre se cumplen en la aritmética de punto flotante debido a errores de redondeo. Por lo tanto, reorganizar el orden de las operaciones a veces puede conducir a resultados más precisos. Por ejemplo, al sumar una gran cantidad de números de punto flotante, sumarlos en orden de magnitud ascendente puede minimizar la acumulación de errores de redondeo. Esta técnica, conocida como suma de Kahan, ayuda a preservar la precisión al compensar los pequeños errores que ocurren durante el proceso de suma.

Además de estas estrategias, es fundamental realizar pruebas y validaciones exhaustivas del Software numérico. Las pruebas unitarias deben diseñarse para cubrir una amplia gama de valores de entrada, incluidos casos extremos que probablemente expongan problemas de precisión. Comparar los resultados de los cálculos de punto flotante con soluciones analíticas conocidas o utilizar aritmética de mayor precisión como referencia puede ayudar a identificar discrepancias.

Además, se puede realizar un análisis de sensibilidad para evaluar cómo pequeños cambios en los valores de entrada afectan la salida, proporcionando información sobre la estabilidad y confiabilidad de los algoritmos numéricos.

Por último, la documentación y la comunicación desempeñan un papel fundamental a la hora de mitigar los problemas de precisión del punto flotante. Los desarrolladores deben documentar las limitaciones y suposiciones de sus algoritmos numéricos, así como cualquier fuente potencial de error. Esta información es invaluable para otros desarrolladores que necesiten mantener o ampliar el Software. Además, una comunicación clara con las partes interesadas sobre la precisión esperada del Software puede ayudar a gestionar las expectativas y evitar mal entendidos.

# 15.1 Aritmética de Punto Flotante IEEE

La aritmética de punto flotante es considerada un tema esotérico para muchas personas. Esto es sorprendente porque el punto flotante es omnipresente en los sistemas informáticos. Casi todos los lenguajes de programación tienen un tipo de datos de punto flotante, i.e. los números no enteros como 1.2 ó 1e+45.

Un número de punto flotante de 64 bits (double en lenguaje C) tiene aproximadamente 16 dígitos decimales de precisión con un rango del orden de  $1.7 \times 10^{-308}$  a  $1.7 \times 10^{308}$  de acuerdo con el estándar 754 de IEEE -El estándar IEEE 754-2008 define varios tamaños de números de punto flotante: media precisión (binary16), precisión simple (binary32), precisión doble (binary64), precisión cuádruple (binary128), etc., cada uno con su propia especificación-, la implementación típica de punto flotante<sup>140</sup>.

Una de las primeras cosas que uno encuentra con sorpresa cuando hace cálculos en una computadora es que por ejemplo, si usamos números muy grandes y seguimos incrementando su valor eventualmente el resultado será negativo ... ¿qué pasó? Esto se llama desbordamiento aritmético al intentar crear un valor numérico que está fuera del rango que puede representarse con un número dado de dígitos, ya sea mayor que el máximo o menor que el mínimo valor representable. Algo similar pasa al restar al menor número representable en la máquina, el resultado será positivo y se denomina subdesbordamiento.

Las computadoras, desde las PC hasta las supercomputadoras, tienen aceleradores de punto flotante; la mayoría de los compiladores deberán compilar algoritmos de punto flotante de vez en cuando y prácticamente todos los sistemas operativos deben responder a excepciones de punto flotante como el desbordamiento.

Desde ya hace mucho tiempo, los procesadores Intel x86 y todos los procesadores de las siguientes generaciones de todas las marcas admiten un formato de precisión extendido de 80 bits con un significado de 64 bits, que es compatible con el especificado en el estándar IEEE. Cuando un compilador usa este formato con registros de 80 bits para acumular sumas y productos internos, está trabajando efectivamente con un redondeo unitario de  $2^{-64}$  en vez de  $2^{-53}$  para precisión doble, dando límites de error más pequeños en un factor de hasta  $2^{11} = 2048$ .

¿Dieciséis lugares decimales es mucho? Casi ninguna cantidad medida se conoce con tanta precisión. Por ejemplo: en la constante en la ley de gravedad de Newton sólo se conoce con seis cifras significativas. En la carga

 $<sup>^{140}</sup>$  Además, existe el número de 80 bits ( $long\ double$  en lenguaje C) que tiene 18 dígitos decimales de precisión con un rango del orden de  $3.4\times10^{-4096}$  a  $1.1\times10^{4096}$ . Y no podemos olvidar, el número de 32 bits (float en lenguaje C) que tiene 14 dígitos decimales de precisión con un rango del orden de  $3.4\times10^{-38}$  a  $3.4\times10^{38}$ .

de un electrón se conoce con 11 cifras significativas, mucha más precisión que la constante gravitacional de Newton, pero aún menos que un número de punto flotante <sup>141</sup>.

Entonces, ¿cuándo no son suficientes 16 dígitos de precisión? Un área problemática es la resta. Las otras operaciones elementales (suma, multiplicación, división) son muy precisas. Siempre que no se presenten desbordamientos y subdesbordamientos, estas operaciones suelen producir resultados que son correctos hasta el último bit. Pero la resta puede ser desde exacta hasta completamente inexacta. Si dos números concuerdan con n cifras, puede perder hasta n cifras de precisión en su resta. Este problema puede aparecer inesperadamente en medio de otros cálculos.

uuQué pasa con el desbordamiento o con el subdesbordamiento? uuCuándo se necesitan números mayores que uu10<sup>308</sup>? No tan a menudo, pero en los cálculos de probabilidad, por ejemplo, se usan todo el tiempo a menos que se haga un uso inteligente.

Es común en probabilidad calcular un número de tamaño mediano que es el producto de un número astronómicamente grande y un número infinitesimalmente pequeño. El resultado final encaja perfectamente en una computadora, pero es posible que los números intermedios no se deban a un desbordamiento o un subdesbordamiento. Por ejemplo, el número máximo de punto flotante en la mayoría de las computadoras está entre 170 factorial y 171 factorial. Estos grandes factoriales aparecen frecuentemente en aplicaciones, a menudo en proporciones con otros grandes factoriales.

- 3.1415 para diseñar los mejores motores.
- 3.1415926535 para obtener la circunferencia de la Tierra dentro de una fracción de pulgada.
- 3.141592653589793 para los cálculos de navegación interplanetaria de la NASA-JPL.
- 3.1415926535897932384626433832795028842 para medir el radio del universo con una precisión igual al tamaño de un átomo de hidrógeno.

En el 2022 se calcularon los primero 100 billones de decimales del número  $\pi$ , para lograr este hito necesitó 157 días, 23 horas, 31 minutos y 7,651 segundos de cálculos, 515 Terabytes de almacenamiento y desplegar un abanico de tecnologías de computación en Compute Engine, un servicio de computación de Google Cloud.

 $<sup>^{141}</sup>$ ; Cuántos dígitos de  $\pi$  necesitamos?

Anatomía de un Número de Punto Flotante Un número de punto flotante de 64 bits codifica un número de la forma  $\pm p \times 2^e$ . El primer bit codifica el signo, 0 para números positivos y 1 para números negativos. Los siguientes 11 bits codifican el exponente e, y los últimos 52 bits codifican la precisión p.

El exponente se almacena con un sesgo de 1023. Es decir, los exponentes positivos y negativos se almacenan todos en un solo número positivo almacenando e+1023 en lugar de almacenarlos directamente. Once bits pueden representar números enteros desde 0 hasta 2047. Restando el sesgo, esto corresponde a valores de e de -1023 a +1024. Definimos  $e_{min}=-1022$  y  $e_{max}=+1023$ . Los valores  $e_{min}-1$  y  $e_{max}+1$  están reservados para usos especiales.

Los números de punto flotante se almacenan normalmente en forma normalizada. En base 10, un número está en notación científica normalizada si el significando es  $\geq 1$  y < 10. Por ejemplo,  $3.14 \times 10^2$  está en forma norma-lizada, pero  $0.314 \times 10^3$  y  $31.4 \times 10^2$  no lo están.

En general, un número en base  $\beta$  está en forma normalizada si tiene la forma  $p \times \beta^e$  donde  $1 \le p < \beta$ . Esto dice que para expresar un número binario, es decir,  $\beta = 2$ , el primer bit del significado de un número normalizado es siempre 1. Dado que este bit nunca cambia, no es necesario almacenarlo. Por lo tanto, podemos expresar 53 bits de precisión en 52 bits de almacenamiento. En lugar de almacenar el significado directamente, almacenamos f, la parte fraccionaria, donde el significado es de la forma 1.f.

El esquema anterior no explica cómo almacenar 0. Es imposible especificar valores de f y e de modo que  $1.f \times 2^e = 0$ . El formato de punto flotante hace una excepción a las reglas establecidas anteriormente. Cuando  $e = e_{min} - 1$  y f = 0, los bits se interpretan como 0. Cuando  $e = e_{min} - 1$  y  $f \neq 0$ , el resultado es un número desnormalizado. Los bits se interpretan como  $0.f \times 2^{e_{min}}$ . En resumen, el exponente especial reservado debajo de  $e_{min}$  se usa para representar 0 y números de punto flotante desnormalizados.

El exponente especial reservado arriba de  $e_{max}$  se usa para representar  $\infty$  y NaN. Si  $e=e_{max}+1$  y f=0, los bits se interpretan como  $\infty$ . Pero si  $e=e_{max}+1$  y  $f\neq 0$ , los bits se interpretan como un NaN o "no es un número" (Not a Number).

Dado que el exponente más grande es 1023 y el significativo más grande es 1.f donde f tiene 52 unidades, el número de punto flotante (en C y C++, esta

constante se define como  $DBL\_MAX$  definido en < float.h>, en Python<sup>142</sup> en  $sys.float\_info)$  más grande es  $2^{1023}(2-2^{-52})=2^{1024}-2^{971}\approx 2^{1024}\approx 1.8\times 10^{308}$ . Los números mayores que  $2^{1024}$  i.e.  $(2-2^{52})$  producen un desbordamiento conocido como Overflow.

Dado que el exponente más pequeño es -1022, el número normalizado positivo más pequeño es  $1.0 \times 2^{-1022} \approx 2.2 \times 10^{-308}$  (en C y C++, esta constante se define como  $DBL\_MIN$  definido en <float.h>, en Python en  $sys.float\_info$ ). Sin embargo, no es el número positivo más pequeño representable como un número de punto flotante, solo el número de punto flotante normalizado más pequeño. Los números más pequeños se pueden expresar en forma desnormalizada, aunque con una pérdida de significado. El número positivo desnormalizado más pequeño ocurre con f que tiene 51 números 0 seguidos de un solo 1. Esto corresponde a  $2^{-52}*2^{-1022}=2^{-1074}\approx 4.9 \times 10^{-324}$ .

Los números que aparecen en los cálculos y tienen magnitud menor que  $2^{-1023}$  i.e.  $(1+2^{-52})$  producen un desbordamiento de la capacidad mínima o subdesbordamiento también conocido como Underflow y, por lo general se igualan a cero.

C y C++ da el nombre de  $DBL\_EPSILON$  al número positivo más pequeño  $\epsilon$  tal que  $1+\epsilon\approx 1$ , también llamado la precisión de la máquina. Dado que el significativo tiene 52 bits, está claro que  $DBL\_EPSILON=2^{-52}\approx 2.2\times 10^{-16}$ . Por eso decimos que un número de punto flotante tiene entre 15 y 16 cifras significativas (decimales).

Formatos de Punto Flotante Se han propuesto varias representaciones diferentes de números reales, pero por mucho la más utilizada es la representación de punto flotante. Las representaciones de punto flotante tienen una base (que siempre se asume que es par) y una precisión p. Si  $\beta = 10$  y p = 3, entonces el número 0.1 se representa como  $1.00 \times 10^{-1}$ . Si  $\beta = 2$  y

```
import sys
print(sys.float_info)

sys.float_info(max=1.7976931348623157e+308, max_exp=1024,
max_10_exp=308, min=2.2250738585072014e-308, min_exp=-1021,
min_10_exp=-307, dig=15, mant_dig=53, epsilon=2.220446049250313e-
16, radix=2, rounds=1)
```

p=24, entonces el número decimal 0.1 no se puede representar exactamente, pero es aproximadamente 1.10011001100110011001101101  $\times$  2<sup>-4</sup>.

En general, un número de punto flotante se representará como  $\pm d.dd...d \times \beta^e$ , donde d.dd...d se llama mantisa y tiene p dígitos. De forma más precisa  $\pm d_0.d_1d_2...d_{p-1} \times \beta^e$  representa el número

$$\pm (d_0 + d_1 \beta^{-1} + \dots + d_{p-1} \beta^{-(p-1)}) \beta^e, (0 \le d_i < \beta).$$
 (15.3)

El término número de punto flotante se utilizará para referirse a un número real que se puede representar exactamente en el formato en discusión. Otros dos parámetros asociados con las representaciones de punto flotante son los exponentes más grandes y más pequeños permitidos,  $e_{max}$  y  $e_{min}$ . Dado que hay  $\beta^p$  posibles significados, y  $e_{max} - e_{min} + 1$  posibles exponentes, un número de punto flotante se puede codificar en

$$[\log_2(e_{max} - e_{min} + 1)] + [\log_2(\beta^p)] + 1$$

bits, donde el +1 final es para el bit de signo. La codificación precisa no es importante por ahora.

Hay dos razones por las que un número real podría no ser exactamente representable como un número de punto flotante. La situación más común se ilustra con el número decimal 0.1. Aunque tiene una representación decimal finita, en binario tiene una representación repetida infinita. Por lo tanto, cuando  $\beta=2$ , el número 0.1 se encuentra estrictamente entre dos números de punto flotante y ninguno de ellos lo puede representar exactamente.

Una situación menos común es que un número real esté fuera de rango, es decir, su valor absoluto sea mayor que  $\beta \times \beta^{e_{max}}$  o menor que  $1.0 \times \beta^{e_{min}}$ . La mayor parte de esta sección analiza cuestiones debidas a la primera razón.

Las representaciones de punto flotante no son necesariamente únicas. Por ejemplo, tanto  $0.01 \times 10^1$  como  $1.00 \times 10^{-1}$  representan 0.1. Si el dígito inicial es distinto de cero  $(d_0 \neq 0$  en la Ec. (15.3)), se dice que la representación está normalizada. El número de punto flotante  $1.00 \times 10^{-1}$  está normalizado, mientras que  $0.01 \times 10^1$  no lo está.

¿Cuándo es Exacta la Conversión de Base de Punto Flotante de Ida y Vuelta? Suponga que almacenamos un número de punto flotante en la memoria, lo imprimimos en base 10 legible por humanos y lo regresamos a memoria. ¿Cuándo se puede recuperar exactamente el número original?

Suponga que comenzamos con la base  $\beta$  con p lugares de precisión y convertimos a la base  $\gamma$  con q lugares de precisión, redondeando al más cercano, luego volvemos a convertir a la base original  $\beta$ . El teorema de Matula dice que si no hay enteros positivos i y j tales que

$$\beta^i = \gamma^j$$

entonces una condición necesaria y suficiente para que la conversión de ida y vuelta sea exacta (suponiendo que no haya desbordamiento o subdesbordamiento) es que

$$\gamma^{q-1} > \beta^p.$$

En el caso de números de punto flotante (por ejemplo doble en C) tenemos  $\beta=2$  y p=53. (ver Anatomía de un Número de Punto Flotante). Estamos imprimiendo a base  $\gamma=10$ . Ninguna potencia positiva de 10 también es una potencia de 2, por lo que se mantiene la condición de Matula en las dos bases.

Si imprimimos q = 17 decimales, entonces

$$10^{16} > 2^{53}$$

por lo que la conversión de ida y vuelta será exacta si ambas conversiones se redondean al más cercano. Si q es menor, algunas conversiones de ida y vuelta no serán exactas.

También puede verificar que para un número de punto flotante de precisión simple (p = precisión de 24 bits) necesita q = 9 dígitos decimales, y para un número de precisión cuádruple (precisión de p = 113 bits) necesita q = 36 dígitos decimales<sup>143</sup>.

Mirando hacia atrás en el teorema de Matula, claramente necesitamos

$$\gamma^q > \beta^p$$
.

¿Por qué? Porque el lado derecho es el número de fracciones de base  $\beta$  y el lado izquierdo es el número de fracciones de base  $\gamma$ . No puede tener un mapa

antoniocarrillo@ciencias.unam.mx 724 Antonio Carrillo Ledesma, Et alii

 $<sup>^{143}\</sup>mathrm{El}$ número de bits asignados para la parte fraccionaria de un número de punto flotante es 1 menos que la precisión: la cifra inicial es siempre 1, por lo que los formatos IEEE ahorran un bit al no almacenar el bit inicial, dejándolo implícito. Entonces, por ejemplo, un doble en C tiene una precisión de 53 bits, pero 52 bits de los 64 bits en un doble se asignan para almacenar la fracción.

uno a uno de un espacio más grande a un espacio más pequeño. Entonces, la desigualdad anterior es necesaria, pero no suficiente. Sin embargo, es casi suficiente. Solo necesitamos una base  $\gamma$  con más cifras significativas, es decir, Matula nos dice

$$\gamma^{q-1} > \beta^p$$

es suficiente. En términos de base 2 y base 10, necesitamos al menos 16 decimales para representar 53 bits. Lo sorprendente es que un decimal más es suficiente para garantizar que las conversiones de ida y vuelta sean exactas. No es obvio a priori que cualquier número finito de decimales adicionales sea siempre suficiente, pero de hecho solo uno más es suficiente.

A continuación, se muestra un ejemplo para mostrar que el decimal adicional es necesario. Suponga que p=5. Hay más números de 2 dígitos que números de 5 bits, pero si solo usamos dos dígitos, la conversión de base de ida y vuelta no siempre será exacta. Por ejemplo, el número 17/16 escrito en binario es  $1.0001_{dos}$  y tiene cinco bits significativos. El equivalente decimal es  $1.0625_{diez}$ , que redondeado a dos dígitos significativos es  $1.1_{diez}$ . Pero el número binario más cercano a  $1.1_{diez}$  con 5 bits significativos es  $1.0010_{dos} = 1.125_{diez}$ . En resumen, redondeando al más cercano da

$$1.0001_{dos} - > 1.1_{diez} - > 1.0010_{dos}$$

y así no volvemos al punto de partida.

Error de Representación se refiere al hecho de que la mayoría de las fracciones decimales no pueden representarse exactamente como fracciones binarias (en base 2). Esta es la razón principal de por qué muchos lenguajes de programación (Python, Perl, C, C++, Java, Fortran, y tantos otros) frecuentemente no mostrarán el número decimal exacto que esperas.

¿Por qué es eso? 1/10 no es representable exactamente como una fracción binaria. Casi todas las máquinas de hoy en día usan aritmética de punto flotante: IEEE-754, y casi todas las plataformas mapean los flotantes al «doble precisión» de IEEE-754. Estos «dobles» tienen 53 bits de precisión, por lo tanto en la entrada la computadora intenta convertir 0.1 a la fracción más cercana que puede de la forma  $J/(2^N)$  donde J es un entero que contiene exactamente 53 bits. Reescribiendo

$$1/10\approx J/2^N$$

antoniocarrillo@ciencias.unam.mx 725 Antonio Carrillo Ledesma, Et alii

como

$$J \approx 2^N/10$$

y recordando que J tiene exactamente 53 bits (es  $\geq 2^{52}$  pero  $< 2^{53}$ ), el mejor valor para N es 56. O sea, 56 es el único valor para N que deja J con exactamente 53 bits. El mejor valor posible para J es entonces el cociente redondeado, si usamos Python para los cálculos tenemos

Ya que el resto es más que la mitad de 10, la mejor aproximación se obtiene redondeándolo

$$>>> q+1$$
 $7205759403792794$ 

Por lo tanto la mejor aproximación a 1/10 en doble precisión 754 es

$$7205759403792794/2**56$$

el dividir tanto el numerador como el denominador reduce la fracción a

$$3602879701896397/2 **55$$

Notemos que como lo redondeamos, esto es un poquito más grande que 1/10; si no lo hubiéramos redondeado, el cociente hubiese sido un poquito menor que 1/10. ¡Pero no hay caso en que sea exactamente 1/10!

Entonces la computadora nunca «ve» 1/10: lo que ve es la fracción exacta de arriba, la mejor aproximación al flotante doble de 754 que puede obtener

Si multiplicamos esa fracción por  $10^{55}$ , podemos ver el valor hasta los 55 dígitos decimales

```
>>> 3602879701896397*10**55 // 2**55 1000000000000000055511151231257827021181583404541015625
```

antoniocarrillo@ciencias.unam.mx 726 Antonio Carrillo Ledesma, Et alii

lo que significa que el valor exacto almacenado en la computadora es igual al valor decimal

0.10000000000000000055511151231257827021181583404541015625.

en lugar de mostrar el valor decimal completo, muchos lenguajes, redondean el resultado a 17 dígitos significativos.

Error de Redondeo Comprimir infinitos números reales en un número finito de bits requiere una representación aproximada. Aunque hay un número infinito de enteros, en la mayoría de los programas el resultado de los cálculos de números enteros se puede almacenar en 32 bits o 64 bits. Por el contrario, dado cualquier número fijo de bits, la mayoría de los cálculos con números reales producirán cantidades que no se pueden representar exactamente usando tantos bits. Por lo tanto, el resultado de un cálculo de punto flotante a menudo debe redondearse para volver a ajustarse a su representación finita. Este error de redondeo es el rasgo característico del cálculo de punto flotante.

Dado que la mayoría de los cálculos de punto flotante tienen errores de redondeo de todos modos, ¿importa si las operaciones aritméticas básicas introducen un poco más de error de redondeo de lo necesario? Esa pregunta es un tema principal a lo largo de esta sección. La sección Dígitos de Guarda analiza los dígitos de protección, un medio para reducir el error al restar dos números cercanos. IBM consideró que los dígitos de guarda eran lo suficientemente importantes que en 1968 añadió un dígito de guarda al formato de doble precisión en la arquitectura System / 360 (la precisión simple ya tenía un dígito de guarda) y modernizó todas las máquinas existentes en el campo.

El estándar IEEE va más allá de sólo requerir el uso de un dígito de protección. Proporciona un algoritmo para la suma, resta, multiplicación, división y raíz cuadrada y requiere que las implementaciones produzcan el mismo resultado que ese algoritmo. Por lo tanto, cuando un programa se mueve de una máquina a otra, los resultados de las operaciones básicas serán los mismos en todos los bits si ambas máquinas admiten el estándar IEEE. Esto simplifica enormemente la portabilidad de programas. Otros usos de esta especificación precisa se dan en operaciones exactamente redondeadas.

Error Relativo y Ulps Dado que el error de redondeo es inherente al cálculo de punto flotante, es importante tener una forma de medir este error. Considere el formato de punto flotante con  $\beta = 10$  y p = 3, que se utilizará

antoniocarrillo@ciencias.unam.mx 727 Antonio Carrillo Ledesma, Et alii

en esta sección. Si el resultado de un cálculo de punto flotante es  $3.12 \times 10^{-2}$ , y la respuesta cuando se calcula con precisión infinita es 0.0314, está claro que tiene un error de 2 unidades en el último lugar. De manera similar, si el número real 0.0314159 se representa como  $3.14 \times 10^{-2}$ , entonces tiene un error de 0.159 unidades en el último lugar.

En general, si el número de punto flotante  $dd..d \times \beta^e$  se usa para representar z, entonces tiene un error de  $|d.d...d - (z/\beta^e)|^{\beta^{p-1}}$  unidades en el último lugar. El término ulps se utilizará como abreviatura de (units in the last place) "unidades en último lugar". Si el resultado de un cálculo es el número de punto flotante más cercano al resultado correcto, aún podría tener un error de hasta 0.5 ulp.

Otra forma de medir la diferencia entre un número de punto flotante y el número real al que se aproxima es el error relativo, que es simplemente la diferencia entre los dos números divididos por el número real. Por ejemplo, el error relativo cometido al aproximar 3.14159 por  $3.14 \times 10^0$  es  $0.00159/3.141590 \approx 0005$ .

Para calcular el error relativo que corresponde a .5 ulp, observe que cuando un número real es aproximado por el número de punto flotante

más cercano posible  $d.dd...dd \times \beta^e$ , el error puede ser tan grande como  $0.00...00\beta' \times \beta^e$ , donde  $\beta'$  es el dígito  $\beta/2$ , hay p unidades en el significado del número de punto flotante y p unidades de 0 en el significado del error. Este error es  $((\beta/2)\beta^{-p}) \times \beta^e$ . Dado que los números de la forma  $d.dd...dd \times \beta^e$  tienen todos el mismo error absoluto, pero tienen valores que oscilan entre  $\beta^e$  y  $\beta \times \beta^e$ , el error relativo varía entre  $((\beta/2)\beta^{-p}) \times \beta^e/\beta^e$  y  $((\beta/2)\beta^{-p}) \times \beta^e/\beta^{e+1}$ . Eso es,

$$\frac{1}{2}\beta^{-p} \le \frac{1}{2}ulp \le \frac{\beta}{2}\beta^{-p} \tag{15.4}$$

En particular, el error relativo correspondiente a 0.5 ulp puede variar en un factor de  $\beta$ . Este factor se llama bamboleo. Estableciendo  $\epsilon = (\beta/2)\beta^{-p}$  en el mayor de los límites en Ec. (15.4) anterior, podemos decir que cuando un número real se redondea al número de punto flotante más cercano, el error relativo siempre está limitado por  $\epsilon$ , que se conoce como épsilon de la máquina.

En el ejemplo anterior, el error relativo fue  $0.00159/3.14159 \approx 0005$ . Para evitar números tan pequeños, el error relativo normalmente se escribe como

antoniocarrillo@ciencias.unam.mx 728 Antonio Carrillo Ledesma, Et alii

factor multiplicado  $\epsilon$ , que en este caso es  $\epsilon = (\beta/2)\beta^{-p} = 5(10)^{-3} = 0.005$ . Por lo tanto, el error relativo se expresaría como  $(0.00159/3.14159)/0.005)\epsilon \approx 0.1\epsilon$ .

Para ilustrar la diferencia entre ulps y error relativo, considere el número real x=12.35. Se aproxima por  $\tilde{x}=1.24\times 10^1$ . El error es 0.5 ulps, el error relativo es  $0.8\epsilon$ . A continuación, considere el cálculo  $8\tilde{x}$ . El valor exacto es 8x=98.8, mientras que el valor calculado es  $8\tilde{x}=9.92\times 10^1$ . El error ahora es 4.0 ulps, pero el error relativo sigue siendo  $0.8\epsilon$ . El error medido en ulps es 8 veces mayor, aunque el error relativo es el mismo.

En general, cuando la base es  $\beta$ , un error relativo fijo expresado en ulps puede oscilar en un factor de hasta  $\beta$ . Y a la inversa, como muestra la Ec. (15.4) anterior, un error fijo de 0.5 ulps da como resultado un error relativo que puede oscilar por  $\beta$ .

La forma más natural de medir el error de redondeo es en ulps. Por ejemplo, el redondeo al número de punto flotante más cercano corresponde a un error menor o igual a 0.5 ulp. Sin embargo, al analizar el error de redondeo causado por varias fórmulas, el error relativo es una mejor medida. Dado que se puede sobrestimar el efecto de redondear al número de punto flotante más cercano por el factor de oscilación de  $\beta$ , las estimaciones de error de las fórmulas serán más estrictas en máquinas con una pequeña  $\beta$ .

Cuando sólo interesa el orden de magnitud del error de redondeo, ulps y  $\epsilon$  pueden usarse indistintamente, ya que difieren como máximo en un factor de  $\beta$ . Por ejemplo, cuando un número de punto flotante tiene un error de n ulps, eso significa que el número de dígitos contaminados es  $\log_{\beta} n$ . Si el error relativo en un cálculo es  $n\epsilon$ , entonces

los dígitos contaminados 
$$\approx \log_{\beta} n$$
. (15.5)

**Dígito de Guarda** Un método para calcular la diferencia entre dos números de punto flotante es calcular la diferencia exactamente y luego redondearla al número de punto flotante más cercano. Esto es muy caro si los operandos difieren mucho en tamaño. Suponiendo que  $p=3, 2.15\times 10^{12}-1.25\times 10^{-5}$  se calcularía como

```
x = 2.15 \times 10^{12}

y = 0.0000000000000000125 \times 10^{12}

x - y = 2.14999999999999875 \times 10^{12}
```

que se redondea a  $2.15 \times 10^{12}$ . En lugar de utilizar todos estos dígitos, el Hardware de punto flotante normalmente funciona con un número fijo

antoniocarrillo@ciencias.unam.mx 729 Antonio Carrillo Ledesma, Et alii

de dígitos. Suponga que el número de dígitos que se mantiene es p, y que cuando el operando más pequeño se desplaza hacia la derecha, los dígitos simplemente se descartan (en contraposición al redondeo). Entonces  $2.15 \times 10^{12} - 1.25 \times 10^{-5}$  se convierte en

```
x = 2.15 \times 10^{12}

y = 0.00 \times 10^{12}

x - y = 2.15 \times 10^{12}
```

La respuesta es exactamente la misma que si la diferencia se hubiera calculado exactamente y luego se hubiera redondeado. Tome otro ejemplo: 10.1 - 9.93. Esto se convierte en

```
x = 1.01 \times 10^{1}
y = 0 - 99 \times 10^{1}
x - y = 0.02 \times 10^{1}
```

La respuesta correcta es 0.17, por lo que la diferencia calculada está desviada en 30 ulps y es incorrecta en todos los dígitos. ¿Qué tan grave puede ser el error?.

**Teorema 5** Usando un formato de punto flotante con parámetros  $\beta$  y p, y calculando las diferencias usando p dígitos, el error relativo del resultado puede ser tan grande como  $\beta - 1$ .

Cuando  $\beta=2$ , el error relativo puede ser tan grande como el resultado, y cuando  $\beta=10$ , puede ser 9 veces mayor. O para decirlo de otra manera, cuando  $\beta=2$ , la Ec. (15.5) muestra que el número de dígitos contaminados es  $log_2(1/\epsilon)=log_2(2^p)=p$ . Es decir, ¡todos los dígitos p del resultado son incorrectos!. Suponga que se agrega un dígito adicional para protegerse contra esta situación (un dígito de guardia). Es decir, el número más pequeño se trunca a p+1 dígitos, y luego el resultado de la resta se redondea a p dígitos. Con un dígito de guarda, el ejemplo anterior se convierte en

```
x = 1.010 \times 10^{1}
y = 0.993 \times 10^{1}
x - y = 0.017 \times 10^{1}
```

y la respuesta es exacta. Con un solo dígito de guarda, el error relativo del resultado puede ser mayor que  $\epsilon$ , como en 110-8.59.

```
x = 1.10 \times 10^{2}
y = 0.085 \times 10^{2}
x - y = 1.015 \times 10^{2}
```

Esto se redondea a 102, en comparación con la respuesta correcta de 101.41, para un error relativo de 0.006, que es mayor que  $\epsilon=0.005$ . En general, el error relativo del resultado puede ser solo un poco mayor que  $\epsilon$ . De forma más precisa:

**Teorema 6** Si x y y son números de punto flotante en un formato con parámetros  $\beta$  y p, y si la resta se realiza con p+1 dígitos (es decir, un dígito de guarda), entonces el error de redondeo relativo en el resultado es menor que  $2\epsilon$ .

Cancelación La última sección se puede resumir diciendo que sin un dígito de guarda, el error relativo cometido al restar dos cantidades cercanas puede ser muy grande. En otras palabras, la evaluación de cualquier expresión que contenga una resta (o una suma de cantidades con signos opuestos) podría resultar en un error relativo tan grande que todos los dígitos carecen de significado (Teorema (5)). Al restar cantidades cercanas, los dígitos más significativos de los operandos coinciden y se cancelan entre sí. Hay dos tipos de cancelación: catastrófica y benigna.

La cancelación catastrófica ocurre cuando los operandos están sujetos a errores de redondeo. Por ejemplo, en la fórmula cuadrática, aparece la expresión  $b^2-4ac$ . Las cantidades  $b^2$  y 4ac están sujetas a errores de redondeo ya que son el resultado de multiplicaciones de punto flotante. Suponga que están redondeados al número de punto flotante más cercano y, por lo tanto, tienen una precisión de 0.5 ulp. Cuando se restan, la cancelación puede hacer que muchos de los dígitos precisos desaparezcan, dejando principalmente dígitos contaminados por errores de redondeo. Por tanto, la diferencia puede tener un error de muchos ulps. Por ejemplo, considere b=3.34, a=1.22 y c=2.28. El valor exacto de  $b^2-4ac$  es 0.0292. Pero  $b^2$  se redondea a 11.2 y 4ac se redondea a 11.1, por lo que la respuesta final es 0.1, que es un error de 100 ulps, aunque 11.2-11.1 es exactamente igual a 10.10. La resta no introdujo ningún error, sino que expuso el error introducido en las multiplicaciones anteriores.

La cancelación benigna ocurre al restar cantidades exactamente conocidas. Si x e y no tienen error de redondeo, entonces, según el Teorema (6), si la resta se realiza con un dígito de guarda, la diferencia x-y tiene un error relativo muy pequeño (menos de  $2\epsilon$ ).

Una fórmula que presenta una cancelación catastrófica a veces se puede reorganizar para eliminar el problema. Considere nuevamente la fórmula

antoniocarrillo@ciencias.unam.mx 731 Antonio Carrillo Ledesma, Et alii

cuadrática

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$
 y  $x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$  (15.6)

Cuando  $b^2 \gg 4ac$ , entonces  $b^2 - 4ac$  no implica una cancelación y

$$\sqrt{b^2 - 4ac} \approx |b|$$
.

Pero la otra suma (resta) en una de las fórmulas tendrá una cancelación catastrófica. Para evitar esto, multiplique el numerador y denominador de  $x_1$  por  $-b - \sqrt{b^2 - 4ac}$  (y de manera similar para  $x_2$ ) para obtener

$$x_1 = \frac{-2c}{b + \sqrt{b^2 - 4ac}}$$
 y  $x_2 = \frac{-2c}{b - \sqrt{b^2 - 4ac}}$  (15.7)

Si  $b^2 \gg ac$  y b > 0, entonces calcular  $x_1$  usando la Ec. (15.6) implicará una cancelación. Por lo tanto, use la Ec. (15.7) para calcular  $x_1$  y (15.6) para  $x_2$ . Por otro lado, si b < 0, use (15.6) para calcular  $x_1$  y (15.7) para  $x_2$ .

La expresión  $x^2-y^2$  es otra fórmula que presenta una cancelación catastrófica. Es más exacto evaluarlo como

$$(x-y)(x+y)$$

A diferencia de la fórmula cuadrática, esta forma mejorada todavía tiene una resta, pero es una cancelación benigna de cantidades sin error de redondeo, no catastrófica. Según el Teorema (6), el error relativo en x-y es como máximo  $2\epsilon$ . Lo mismo ocurre con x+y. Multiplicar dos cantidades con un pequeño error relativo da como resultado un producto con un pequeño error relativo.

Errores de Redondeo y de Aritmética La aritmética que realiza una computadora es distinta de la aritmética de nuestros cursos de álgebra o cálculo. En nuestro mundo matemático tradicional consideramos la existencia de números con una cantidad infinita de cifras, en la computadora cada número representable tienen sólo un número finito, fijo de cifras (véase 2.4), los cuales en la mayoría de los casos es satisfactoria y se aprueba sin más, aunque a veces esta discrepancia puede generar problemas.

Un ejemplo de este hecho lo tenemos en el cálculo de raíces de:

$$ax^2 + bx + c = 0$$

antoniocarrillo@ciencias.unam.mx 732 Antonio Carrillo Ledesma, Et alii

cuando  $a \neq 0$ , donde las raíces se calculan comúnmente con el algoritmo Ec. (15.6) o de forma alternativa con el algoritmo que se obtiene mediante la racionalización del numerador Ec. (15.7).

Otro algoritmo que podemos implementar es el método de Newton-Raphson<sup>144</sup> para buscar raíces, que en su forma iterativa está dado por

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

en el cual se usa  $x_0$  como una primera aproximación a la raíz buscada y  $x_n$  es la aproximación a la raíz después de n iteraciones (sí se converge a ella), donde  $f(x) = ax^2 + bx + c$  y  $f'(x_i) = 2ax + b$ .

Salida del Cálculo de Raíces Para resolver el problema, usamos por ejemplo el siguiente código en Python usando programación procedimental:

```
import math
\operatorname{def} f(x, a, b, c):
  """ Evalua la Funcion cuadratica """
  return (x * x * a + x * b + c);
\operatorname{def} \operatorname{df}(x, a, b):
  """ Evalua la derivada de la funcion cuadratica """
  return (2.0 * x * a + b);
def evalua(x, a, b, c):
  """ Evalua el valor X en la función cuadratica """
  print('Raiz (%1.16f), evaluacion raiz: %1.16e' % (x, f(x, a, b, c)))
def metodoNewtonRapson(x, ni, a, b, c):
  """ Metodo Newton-Raphson x = x - f(x)/f'(x) """
  for i in range(ni):
     x = x - (f(x, a, b, c) / df(x, a, b))
  return x
def raices(A, B, C):
  """ Calculo de raices """
```

$$x_{i+1} = x_i - \frac{f(x_i)f'(x_i)}{[f'(x_i)]^2 - f(x_i)f''(x_i)}$$

que involucra la función f(x), la primera derivada f'(x) y a la segunda derivada f''(x).

antoniocarrillo@ciencias.unam.mx 733 Antonio Carrillo Ledesma, Et alii

 $<sup>^{144}\</sup>mathrm{También}$ podemos usar otros métodos, como el de Newton Raphson Modificado para acelerar la convergencia

```
if A == 0.0:
    print("No es una ecuación cuadratica")
    exit(1)
  # Calculo del discriminante
  d = B * B - 4.0 * A * C
  # Raices reales
  if d >= 0.0:
    print('\nPolinomio (\%f) X^2 + (\%f) X + (\%f) = 0\n' \% (A, B, C))
    print('\nChicharronera 1')
    X1 = (-B + \text{math.sqrt(d)}) / (2.0 * A)
    X2 = (-B - \text{math.sqrt}(d)) / (2.0 * A)
    evalua(X1, A, B, C)
    evalua(X2, A, B, C)
    print('\nChicharronera 2')
    X1 = (-2.0 * C) / (B + math.sqrt(d))
    X2 = (-2.0 * C) / (B - math.sqrt(d))
    evalua(X1, A, B, C)
    evalua(X2, A, B, C)
    # Metodo Newton-Raphson
    print("\n\nMetodo Newton-Raphson")
    x = X1 - 1.0;
    print("\nValor inicial aproximado de X1 = %1.16f" % x)
    x = metodoNewtonRapson(x, 6, A, B, C)
    evalua(x, A, B, C);
    x = X2 - 1.0;
    print("\nValor\ inicial\ aproximado\ de\ X2 = \%1.16f"\ \%\ x);
    x = metodoNewtonRapson(x, 6, A, B, C)
    evalua(x, A, B, C)
    print("\n")
  else:
    # Raices complejas
    print("Raices Complejas ...")
if name _ == '__main__':
  raices(1.0, 4.0, 1.0)
y generan la siguiente salida:
Polinomio (1.000000) X^2 + (4.000000) X + (1.000000) = 0
```

#### Chicharronera 1

Raiz (-0.2679491924311228), evaluacion raiz: -4.4408920985006262e-16 Raiz (-3.7320508075688772), evaluacion raiz: 0.00000000000000000000e+00

#### Chicharronera 2

Raiz (-0.2679491924311227), evaluacion raiz: 0.0000000000000000000+00 Raiz (-3.7320508075688759), evaluacion raiz: -5.3290705182007514e-15

## Metodo Newton-Raphson

Valor inicial aproximado de X1 = -1.2679491924311228

Valor inicial aproximado de X2 = -4.7320508075688759

En esta salida se muestra la raíz calculada y su evaluación en la ecuación cuadrática, la cual debería de ser cero al ser una raíz, pero esto no ocurre en general por los errores de redondeo. Además, nótese el impacto de seleccionar el algoritmo numérico adecuado a los objetivos que persigamos en la solución del problema planteado.

En cuanto a la implementación computacional, el paradigma de programación seleccionado depende la complejidad del algoritmo a implementar y si necesitamos reusar el código generado o no. Otras implementaciones computacionales se pueden consultar en las ligas, en las cuales se usan distintos lenguajes (C, C++, Java y Python) y diferentes paradigmas de programación ( secuencial, procedimental y orientada a objetos).

Si lo que necesitamos implementar computacionalmente es una fórmula o conjunto de ellas que generen un código de decenas de líneas, la implementación secuencial es suficiente, si es menor a una centena de líneas puede ser mejor opción la implementación procedimental y si el proyecto es grande o complejo, seguramente se optará por la programación orientada a objetos o formulaciones híbridas de las anteriores.

En última instancia, lo que se persigue en la programación es generar un código: correcto, claro, eficiente, de fácil uso y mantenimiento, que sea flexible, reusable y en su caso portable.

# 15.2 Trabajando con Punto Flotante

Hay varias trampas en las que incluso los programadores muy experimentados caen cuando escriben código que depende de la aritmética de punto flotante. En esta sección explicamos algunas cosas a tener en cuenta al trabajar con números de punto flotante, es decir, tipos de datos: float (32 bits), double (64 bits) o long double (80 bits).

**Problemas de Precisión** Como ya vimos en las secciones precedentes, los números de punto flotante se representan en el Hardware de la computadora en fracciones en base 2 (binario). Por ejemplo, la fracción decimal

0.125

tiene el valor 1/10+2/100+5/1000, y de la misma manera la fracción binaria

0.001

tiene el valor 0/2 + 0/4 + 1/8. Estas dos fracciones tienen valores idénticos, la única diferencia real es que la primera está escrita en notación fraccional en base 10 y la segunda en base 2.

Desafortunadamente, la mayoría de las fracciones decimales no pueden representarse exactamente como fracciones binarias. Como consecuencia, en general los números de punto flotante decimal que usamos en la computadora son sólo aproximados por los números de punto flotante binario que realmente se guardan en la máquina.

El problema es más fácil de entender primero en base 10. Consideremos la fracción 1/3. Podemos aproximarla como una fracción de base 10 como: 0.3 o, mejor: 0.33 o, mejor: 0.333 y así. No importa cuántos dígitos desees escribir, el resultado nunca será exactamente 1/3, pero será una aproximación cada vez mejor de 1/3.

De la misma manera, no importa cuántos dígitos en base 2 quieras usar, el valor decimal 0.1 no puede representarse exactamente como una fracción en base 2. En base 2, 1/10 es la siguiente fracción que se repite infinitamente

si nos detenemos en cualquier número finito de bits, y tendremos una aproximación. En la mayoría de las máquinas hoy en día, los double se aproximan

antoniocarrillo@ciencias.unam.mx 736 Antonio Carrillo Ledesma, Et alii

usando una fracción binaria con el numerador usando los primeros 53 bits con el bit más significativo y el denominador como una potencia de dos. En el caso de 1/10, la fracción binaria es

$$3602879701896397/2^{55}$$

que está cerca pero no es exactamente el valor verdadero de 1/10.

La mayoría de los usuarios no son conscientes de esta aproximación por la forma en que se muestran los valores. Varios lenguajes de programación como C, C++, Java y Python solamente muestra una aproximación decimal al valor verdadero decimal de la aproximación binaria almacenada por la máquina. En la mayoría de las máquinas, si fuéramos a imprimir el verdadero valor decimal de la aproximación binaria almacenada para 0.1, debería mostrar

#### 0.1000000000000000055511151231257827021181583404541015625

esos son más dígitos que lo que la mayoría de la gente encuentra útil, por lo que los lenguajes de programación mantiene manejable la cantidad de dígitos al mostrar en su lugar un valor redondeado

1/10

como

0.1

sólo hay que recordar que, a pesar de que el valor mostrado resulta ser exactamente 1/10, el valor almacenado realmente es la fracción binaria más cercana posible. Esto queda de manifiesto cuando hacemos

$$0.1 + 0.1 + 0.1$$
 ó  $0.1 * 3$  ó  $0.1 + 0.2^{145}$ 

obtendremos

### 0.300000000000000004

```
\begin{aligned} & \text{double } x = 0.1; \\ & \text{double } y = 0.1; \\ & \text{double } z = (x*10.0 + y*10.0) \; / \; 10.0 \end{aligned}
```

antoniocarrillo@ciencias.unam.mx 737 Antonio Carrillo Ledesma, Et alii

 $<sup>^{145}\</sup>mathrm{Para}$ el caso de 0.1+0.2 podemos hacer un programa que nos de 0.3, usando:

que es distinto al 0.3 que esperábamos.

Es de hacer notar que hay varios números decimales que comparten la misma fracción binaria más aproximada. Por ejemplo, los números

0.1, 0.100000000000000000001 y

#### 0.1000000000000000055511151231257827021181583404541015625

son todos aproximados por  $3602879701896397/2^{55}$ .

Notemos que esta es la verdadera naturaleza del punto flotante binario: no es un error del lenguaje de programación y tampoco es un error en tu código. Verás lo mismo en todos los lenguajes que soportan la aritmética de punto flotante de tu Hardware (a pesar de que en algunos lenguajes por omisión no muestren la diferencia, o no lo hagan en todos los modos de salida). Para una salida más elegante, quizás quieras usar el formateo de cadenas de texto para generar un número limitado de dígitos significativos.

Ejemplo, el número decimal 9.2 se puede expresar exactamente como una relación de dos enteros decimales 92/10, los cuales se pueden expresar exactamente en binario como 0b1011100/0b1010. Sin embargo, la misma proporción almacenada como un número de punto flotante nunca es exactamente igual a 9.2:

ya que se guarda como la fracción:  $5179139571476070/2^{49}$ .

Otro ejemplo, si tomamos el caso del número 0.02 y vemos su representación en el lenguaje de programación Python, tenemos que:

import decimal print(decimal.Decimal(0.02))

y el resultado es:

0.0200000000000000004163336342344337026588618755340576171875

Además, tenemos la no representabilidad de  $\pi$  (y  $\pi/2$ ), esto significa que un intento de cálculo de  $tan(\pi/2)$  no producirá un resultado de infinito, ni

antoniocarrillo@ciencias.unam.mx 738 Antonio Carrillo Ledesma, Et alii

pero

siquiera se desbordará en los formatos habituales de punto flotante. Simplemente no es posible que el Hardware de punto flotante estándar intente calcular  $tan(\pi/2)$ , porque  $\pi/2$  no se puede representar exactamente. Este cálculo en C:

```
doble pi =3.1415926535897932384626433832795; tan (pi / 2.0);
```

dará un resultado de 1.633123935319537e + 16. En precisión simple usando tanf, el resultado será -22877332.0.

De la misma manera, un intento de cálculo de  $sin(\pi)$  no arrojará cero. El resultado será 1.2246467991473532e-16 en precisión doble.

Si bien la suma y la multiplicación de punto flotante son conmutativas (a+b=b+a y  $a\times b=b\times a$ ), no son necesariamente asociativas. Es decir, (a+b)+c no es necesariamente igual a a+(b+c). Usando aritmética decimal significativa de 7 dígitos:

```
a = 1234.567, b = 45.67834, c = 0.0004 (a+b)+c: 1234.567+45.67834 \qquad 1280.24534 \text{ se redondea a } 1280.245 1280.245+0.0004 \qquad 1280.2454 \text{ se redondea a } 1280.245 a+(b+c): 45.67834+0.0004 \qquad 45.67874 1234.567+45.6787 \qquad 1280.24574 \text{ se redondea a } 1280.246
```

Tampoco son necesariamente distributivos. Es decir  $(a+b) \times c$  puede no ser lo mismo que  $a \times c + b \times c$ :

```
1234.567*3.333333 = 4115.223, 1.234567*3.333333 = 4.115223, 4115.223 + 4.115223 = 4119.338
```

1234.567 + 1.234567 = 1235.802, 1235.802 \* 3.333333 = 4119.340

antoniocarrillo@ciencias.unam.mx 739 Antonio Carrillo Ledesma, Et alii

Algunos Trucos Supongamos que necesitamos hacer el siguiente cálculo:

$$4.56 * 100$$

la respuesta es:

#### 455.9999999999994

que no es la esperada. ¿Qué podemos hacer para obtener lo que esperamos?. Por ejemplo, con Python, podemos usar:

```
from sympy import * print(nsimplify(4.56 * 100, tolerance=1e-1))
```

que nos dará el 456 esperado.

SymPy es un paquete matemático simbólico para Python, y nsimplify toma un número de punto flotante y trata de simplificarlo como una fracción con un denominador pequeño, la raíz cuadrada de un entero pequeño, una expresión que involucra constantes famosas, etc.

Por ejemplo, supongamos que algún cálculo arrojó 4.242640687119286 y sospechamos que hay algo especial en ese número. Así es como puede probar de dónde vino:

```
from sympy import * print(nsimplify(4.242640687119286))
```

que nos arrojará:

$$3 * sqrt(2)$$

Tal vez hagamos un cálculo numérico y encontremos una expresión simple para el resultado y eso sugiera una solución analítica. Creo que una aplicación más común de nsimplify podría ser ayudarte a recordar fórmulas medio olvidadas. Por ejemplo, quizás estés oxidado con tus identidades trigonométricas, pero recuerdas que  $\cos(p/6)$  es algo especial.

```
from sympy import *
print(nsimplify(cos(pi/6)))
```

que nos entregará:

O, para tomar un ejemplo más avanzado, supongamos que recordamos vagamente que la función gamma toma valores reconocibles en valores semienteros, pero no recordamos exactamente cómo. Tal vez algo relacionado con  $\pi$  o e. Puede sugerir que nsimplify incluya expresiones con  $\pi$  y e en su búsqueda.

```
from sympy import * print(nsimplify(gamma(3.5), constants=[pi, E]))
```

obteniendo:

$$15 * sqrt(pi)/8$$

También podemos darle a nsimplify una tolerancia, pidiéndole que encuentre una representación simple dentro de una vecindad del número. Por ejemplo, aquí hay una manera de encontrar aproximaciones a  $\pi$ .

```
from sympy import * print(nsimplify(pi, tolerance=1e-5))
```

obteniendo:

con una tolerancia más amplia, devolverá una aproximación más simple.

```
from sympy import * print(nsimplify(pi, tolerance=1e-2))
```

obteneindo:

finalmente, aquí hay una aproximación de mayor precisión a  $\pi$  que no es exactamente simple:

```
from sympy import * print(nsimplify(pi, tolerance=1e-7))
```

obteniendo:

$$exp(141/895 + sqrt(780631)/895)$$

antoniocarrillo@ciencias.unam.mx 741 Antonio Carrillo Ledesma, Et alii

Seno de un Googol ¿Cómo se evalúa el seno de un número grande en aritmética de punto flotante? ¿Qué significa el resultado?.

Seno de un billón Comencemos por encontrar el seno de un billón  $(10^{12})$  usando aritmética de punto flotante. Hay un par de maneras de pensar en esto. El número de punto flotante t=1.0e12 solo se puede representar con 15 o 16 cifras decimales significativas (para ser precisos, 53 bits<sup>146</sup>), por lo que podría considerarlo como un representante del intervalo de números que tienen todos la misma representación de punto flotante. Cualquier resultado que sea seno de un número en ese intervalo debe considerarse correcto.

Otra forma de ver esto es decir que t se puede representar exactamente -su representación binaria requiere 42 bits y tenemos 53 bits de precisión significativa disponibles-, por lo que esperamos sin(t) devolver el seno de exactamente un billón, con precisión total.

Resulta que la aritmética del IEEE hace lo último, calculando sin(1e12) correctamente hasta 16 dígitos.

Aquí está el resultado en Python.

```
sin(1.0e12)
-0.6112387023768895
```

y verificado por Mathematica calculando el valor con 20 decimales

```
In:= N[Sin[10^12], 20]
out= -0.61123870237688949819
```

**Reducción de alcance** tenga en cuenta que el resultado anterior no es el que obtendría si primero tomara el resto al dividir por  $2\pi$  y luego tomará el seno.

```
\sin(1.0e12 \% (2*pi))
-0,6112078499756778
```

 $<sup>^{146}</sup>$ Un doble IEEE 754 tiene 52 bits significativos, pero estos bits pueden representar 53 bits de precisión porque el primer bit de la parte fraccionaria es siempre 1, por lo que no es necesario representarlo.

Esto tiene sentido. El resultado de dividir un billón por la representación de punto flotante de  $2\pi$ , 159154943091.89536, es correcto con precisión de punto flotante total. Pero la mayor parte de esa precisión está en la parte entera. La parte fraccionaria solo tiene cinco dígitos de precisión, por lo que debemos esperar que el resultado anterior sea correcto hasta un máximo de cinco dígitos. De hecho, tiene una precisión de cuatro dígitos.

Cuando su procesador calcula, sin(1e12) no toma ingenuamente el resto por  $2\pi$  y luego calcula el seno. Si así fuera, obtendríamos cuatro cifras significativas en lugar de 16.

Empezamos diciendo que había dos maneras de ver el problema y, según la primera, devolver sólo cuatro cifras significativas sería bastante razonable. Si pensamos en el valor t como una cantidad medida, medida con la precisión de la aritmética de punto flotante (aunque casi nada se puede medir con tanta precisión), entonces todo lo que deberíamos esperar sería cuatro cifras significativas. Pero las personas que diseñaron la implementación de la función seno en su procesador hicieron más de lo que se esperaba que hicieran, encontrando el seno de un billón con total precisión. Lo hacen utilizando un algoritmo de reducción de rango que conserva mucha más precisión que hacer una división ingenuamente.

¿Seno de un Googol? ¿Qué pasa si intentamos tomar el seno de un número ridículamente grande como un Googol (10<sup>100</sup>)? Esto no funcionará porque un Googol no se puede representar exactamente como un número de punto flotante (es decir, como un doble IEEE 754). No es demasiado grande; Los números de punto flotante pueden ser tan grandes como alrededor de 10<sup>308</sup>. El problema es que un número tan grande no se puede representar con total precisión. Pero podemos representar 2<sup>333</sup> exactamente, y un Googol está entre 2<sup>332</sup> y 2<sup>333</sup>. Y sorprendentemente, la función sinusoidal de Python (o más bien la función sinusoidal integrada en el procesador AMD de mi computadora) devuelve el resultado correcto con total precisión.

```
\sin(2^{**}333)
0.9731320373846353
```

¿Cómo sabemos que esto es correcto? Lo verifiqué en Mathematica:

```
In:= \sin[2.0^{3}33]
Out= 0.9731320373846353
```

¿Cómo sabemos que Mathematica tiene razón? Podemos hacer una reducción de rango ingenua usando precisión extendida, digamos 200 decimales.

```
In:= p = N[Pi, 200]
In:= theta = x - IntegerPart[x/ (2 p)] 2 p
Out= 1.8031286178334699559384196689...
```

y verificar que el seno del valor reducido del rango sea correcto.

```
In:= Sin[theta]
Out= 0.97313203738463526...
```

Interpretación Aritmética de Intervalos debido a que los números de punto flotante tienen 53 bits de precisión, todos los números reales entre  $2^{56} - 2^2$  y  $2^{56} + 2^2$  se representan como el número de punto flotante  $2^{56}$ . Este es un rango de ancho 8, más ancho que  $2\pi$ , por lo que los senos de los números en este rango cubren los posibles valores de seno, es decir, [-1,1]. Entonces, si piensa en un número de punto flotante como una muestra del conjunto de todos los números reales con la misma representación binaria, cada valor posible de seno es un valor de retorno válido para números mayores que  $2^{56}$ . Pero si piensa que un número de punto flotante se representa exactamente a sí mismo, entonces tiene sentido preguntar por el seno de números como  $2^{333}$  y mayores, hasta los límites de la representación de punto flotante  $2^{147}$ .

No Pruebes por la Igualdad Cuando se usa Punto Flotante no es recomendable escribir código como el siguiente<sup>148</sup>:

```
double x;
double y;
...
if (x == y) \{...\}
```

 $<sup>^{-147}</sup>$ El mayor exponente de un doble IEEE es 1023, y el mayor significado es  $2-2^{-53}$  (es decir, todos los unos), por lo que el mayor valor posible de un doble es  $(2^{53}-1)2^{1024-53}$  y de hecho la expresión de Python  $\sin((2^{**}53-1)^*2^{**}(1024-53))$  devuelve el valor correcto con  $^{16}$  cifras significativas.

 $<sup>^{148}</sup>$ Sin pérdida de generalidad usamos algún lenguaje de programación en particular para mostrar los ejemplos, pero esto pasa en todos los lenguajes que usan operaciones de Punto Flotante.

La mayoría de las operaciones de punto flotante implican al menos una pequeña pérdida de precisión y, por lo tanto, incluso si dos números son iguales para todos los fines prácticos, es posible que no sean exactamente iguales hasta el último bit, por lo que es probable que la prueba de igualdad falle. Por ejemplo:

```
double x = 10;
double y = \operatorname{sqrt}(x);
y *= y;
if (x == y)
cout << "La reaiz cuadrada es exacta\n";
else
cout << x-y << "\n";
```

el código imprime: -1.778636e-015, aunque en teoría, elevar al cuadrado debería deshacer una raíz cuadrada, la operación de ida y vuelta es ligeramente inexacta. En la mayoría de los casos, la prueba de igualdad anterior debe escribirse de la siguiente manera:

```
double tolerancia = ... if (fabs(x - y) < tolerancia) \{...\}
```

Aquí la tolerancia es un umbral que define lo que está "lo suficientemente cerca" para la igualdad. Esto plantea la pregunta de qué tan cerca está lo suficientemente cerca. Esto no puede responderse en abstracto; tienes que saber algo sobre tu problema particular para saber qué tan cerca está lo suficientemente cerca en tu contexto.

Por ejemplo: ¿hay alguna garantía de que la raíz cuadrada de un cuadrado perfecto sea devuelta exactamente?, por ejemplo si hago sqrt(81.0) == 9.0, por lo visto anteriormente, la respuesta es no, pero podríamos cambiar la pregunta por 9.0\*9.0 == 81.0, esto funcionará siempre que el cuadrado esté dentro de los límites de la magnitud del punto flotante.

Por otro lado, es posible que las expectativas de las matemáticas no se cumplan en el campo del cálculo de punto flotante. Por ejemplo, se sabe que

$$(x+y)(x-y) = x^2 + y^2$$

y esta otra

$$\sin^2 \theta + \cos^2 \theta = 1$$

antoniocarrillo@ciencias.unam.mx 745 Antonio Carrillo Ledesma, Et alii

sin embargo, no se puede confiar en estos hechos cuando las cantidades involucradas son el resultado de un cálculo de punto flotante. Además, el error de redondeo puede afectar la convergencia y precisión de los procedimientos numéricos iterativos.

Aún más y sin pérdida de generalidad, si comparamos usando el lenguaje de programación Python:

en todos los casos dará un resultado de falso, además:

```
print(10.4 + 20.8 > 31.2)

print(0.8 - 0.1 > 0.7)
```

el resultado será verdadero. Por ello es siempre conveniente ver con que números trabajamos usando algo como:

```
print(format(0.1, ".17g"))
print(format(0.2, ".17g"))
print(format(0.3, ".17g"))
```

así cuando sumemos 0.1 + 0.2, podemos ver el verdadero resultado:

```
print(print(format(0.1 + 0.2, ".17g")))
```

Teniendo esto en cuenta podemos comparar usando:

```
import math print(math.isclose(0.1 + 0.2, 0.3)
```

nos dará la respuesta esperada. Podemos ajustar la tolerancia relativa usando:

```
import math print(math.isclose(0.1 + 0.2, 0.3, rel_tol = 1e-20)
```

que en este caso nos dirá que es falso pues no son iguales en los 20 primeros dígitos solicitados.

Para las comparaciones >= o <=, por ejemplo para a+b <= c se debe usar:

```
a, b, c = 0.1, 0.2, 0.3
print(math.isclose(a + b, c) or (a + c < c))
```

antoniocarrillo@ciencias.unam.mx 746 Antonio Carrillo Ledesma, Et alii

¿Cómo Compara Python los Flotantes y los Enteros? veamos los siguientes ejemplos:

```
>>> 9007199254740992 == 9007199254740992.0 True >>> 9007199254740993 == 9007199254740993.0 False >>> 9007199254740994 == 9007199254740994.0 True
```

¿Qué fue lo que salió mal?, para desmenuzar lo que pasó, veamos primero:

Representación IEEE-754 de 9007199254740992.0 comencemos con el valor del primer escenario: 9007199254740992.0, resulta que este número es en realidad es  $2^{53}$ . Eso significa que en la representación binaria sería 1 seguido de 53 ceros:

la forma normalizada será  $1.0*2^{53}$ , que se obtendrá desplazando los bits 53 lugares a la derecha, lo que nos dará el exponente 53.

Además, recordemos que la mantisa tiene solo 52 bits de ancho y estamos desplazando nuestro valor 53 bits hacia la derecha, lo que significa que se perderá el bit menos significativo (LSB) de nuestro valor original. Sin embargo, como todos los bits aquí son 0, no hay ningún daño y aún podemos representar 9007199254740992.0 con precisión. Resumamos los componentes de 9007199254740992.0:

Bit de signo: 0

Exponente: 53 + 1023 (bias) = 1076 (10000110100 en binario)

Representación IEEE-754:

Representación IEEE-754 de 9007199254740993.0 el valor del segundo escenario de prueba es 9007199254740993.0 el que es uno mayor que el valor anterior. Su representación binaria se puede obtener simplemente sumando 1 a la representación binaria del valor anterior:

#### 

nuevamente, para convertir esto a la forma normalizada, tendremos que desplazarlo a la derecha 53 bits, lo que nos dará el exponente 53.

Sin embargo, la mantisa tiene sólo 52 bits de ancho. Cuando desplazamos nuestro valor a la derecha en 53 bits, el bit menos significativo (LSB) con el valor 1 se perderá, lo que dará como resultado que la mantisa sea todo ceros. Esto conduce a los siguientes componentes:

Bit de signo: 0

Exponente: 53 + 1023 (bias) = 1076 (10000110100 en binario)

Representación IEEE-754:

Observe que la representación IEEE-754 de 9007199254740993.0 es la misma que la de 9007199254740992.0. Esto significa que en su representación en memoria del número 9007199254740993.0 en realidad se representa como 9007199254740992.0.

Esto explica por qué Python da el resultado para

9007199254740993 == 9007199254740993.0

False, porque lo ve como una comparación entre

9007199254740993 y 9007199254740992.0.

Representación IEEE-754 de 9007199254740994.0 El número en el tercer y último escenario de prueba es 9007199254740994.0, que es 1 más que el valor anterior. Su representación binaria se puede obtener sumando 1 a la representación binaria de 9007199254740993.0, dándonos:

antoniocarrillo@ciencias.unam.mx 748 Antonio Carrillo Ledesma, Et alii

esto también tiene 54 bits antes del punto binario y requiere un desplazamiento a la derecha de 53 bits para convertirlo a la forma normalizada, lo que nos da el exponente 53.

Observe que esta vez, el segundo bit menos significativo (LSB) tiene el valor 1. Por lo tanto, cuando desplazamos este número a la derecha 53 bits, se convierte en el LSB de la mantisa (que tiene 52 bits de ancho). Los componentes se ven así:

Bit de signo: 0

Exponente: 53 + 1023 (bias) = 1076 (10000110100 en binario)

Representación IEEE-754:

A diferencia de 9007199254740993.0, la representación IEEE-754 del número 9007199254740994.0 puede representarlo exactamente sin pérdida de precisión. Por tanto, el resultado de

9007199254740994 == 9007199254740994.0

es True en Python.

El estándar IEEE-754 y la aritmética de punto flotante son inherentemente complejos y comparar números de punto flotante no es sencillo. Lenguajes como C y Java implementan la promoción de tipos implícita, convirtiendo números enteros en dobles y comparándolos poco a poco. Sin embargo, esto aún puede dar lugar a resultados inesperados debido a la pérdida de precisión.

Python es único a este respecto. Tiene números enteros de precisión infinita, lo que hace que la promoción de tipos sea inviable en muchas situaciones. En consecuencia, Python utiliza su algoritmo especializado para comparar estos números, que a su vez tiene casos extremos.

Ya sea que esté utilizando C, Java, Python u otro lenguaje, es recomendable utilizar funciones de biblioteca para comparar valores de punto flotante en lugar de realizar comparaciones directas para evitar posibles errores.

## Cuando Cómputo de Alto Rendimiento no es Alto Rendimiento

A todo el mundo le importa que los códigos se ejecuten rápidamente en sus computadoras. Las mejoras de Hardware de las últimas décadas lo han hecho posible. Pero, ¿qué tan bien estamos aprovechando las aceleraciones del Hardware?

Considere estos dos ejemplos de código C++. Supongamos aquí n = 10000000.

```
void sub(int* a, int* b) {
    for (int i=0; i<n; ++i)
        a[i] = i + 1;
    for (int i=0; i<n; ++i)
        b[i] = a[i];
}

void sub(int* a, int* b) {
    for (int i=0; i<n; ++i) {
        const int j = i + 1;
        a[i] = j;
        b[i] = j;
    }
}</pre>
```

¿Cuál corre más rápido? Ambos son simples y dan resultados idénticos (suponiendo que no haya alias). Sin embargo, en las arquitecturas modernas, dependiendo de la configuración de compilación, una generalmente se ejecutará significativamente más rápido que la otra.

En particular, se esperaría que el fragmento 2 se ejecutara más rápido que el fragmento 1. En el fragmento 1, los elementos de la matriz "a", que es demasiado grande para almacenarse en caché, deben recuperarse de la memoria después de escribirse, pero esto no es necesario para Fragmento 2. La tendencia durante más de dos décadas ha sido que la velocidad de computación de los sistemas recién entregados crezca mucho más rápido que la velocidad de la memoria, y la disparidad es extrema hoy en día. El rendimiento de estos núcleos depende casi por completo de la velocidad del ancho de banda de la memoria. Así, el fragmento 2, una versión de bucle fusionado del fragmento 1, mejora la velocidad al reducir el acceso a la memoria principal.

Es poco probable que bibliotecas como C++ STL ayuden, ya que esta operación es demasiado especializada para esperar que una biblioteca la admita (especialmente la versión de bucle fusionado). Además, el compilador no puede fusionar de forma segura los bucles automáticamente sin instrucciones específicas de que los punteros no tengan alias, y aun así no se garantiza que lo haga.

Afortunadamente, los lenguajes informáticos de alto nivel desde la década de 1950 han elevado el nivel de abstracción de la programación para todos nosotros. Naturalmente, a muchos de nosotros nos gustaría simplemente implementar la lógica empresarial requerida en nuestros códigos y dejar que el compilador y el Hardware hagan el resto. Pero, lamentablemente, no siempre se puede simplemente colocar el código en una computadora y esperar que se ejecute rápidamente. Cada vez más, a medida que el Hardware se vuelve más complejo, prestar atención a la arquitectura subyacente es fundamental para obtener un alto rendimiento.

Preocúpate más por la Suma y la Resta que por la Multiplicación y la División Los errores relativos en la multiplicación y división son siempre pequeños. La suma y la resta, por otro lado, pueden resultar en una pérdida completa de precisión. Realmente el problema es la resta; la suma sólo puede ser un problema cuando los dos números que se agregan tienen signos opuestos, por lo que puedes pensar en eso como una resta. Aún así, el código podría escribirse con un "+" que sea realmente una resta.

La resta es un problema cuando los dos números que se restan son casi iguales. Cuanto más casi iguales sean los números, mayor será el potencial de pérdida de precisión. Específicamente, si dos números están de acuerdo con n bits, se pueden perder n bits de precisión en la resta. Esto puede ser más fácil de ver en el extremo: si dos números no son iguales en teoría pero son iguales en su representación de máquina, su diferencia se calculará como cero, 100% de pérdida de precisión.

Aquí hay un ejemplo donde tal pérdida de precisión surge a menudo. La derivada de una función f en un punto x se define como el límite de

$$(f(x+h)-f(x))/h$$

cuando h llega a cero. Entonces, un enfoque natural para calcular la derivada de una función sería evaluar

$$(f(x+h)-f(x))/h$$

antoniocarrillo@ciencias.unam.mx 751 Antonio Carrillo Ledesma, Et alii

para alguna h pequeña. En teoría, cuanto menor es h, mejor se aproxima esta fracción a la derivada. En la práctica, la precisión mejora por un tiempo, pero más allá de cierto punto, valores más pequeños de h resultan en peores aproximaciones a la derivada. A medida que h disminuye, el error de aproximación se reduce pero el error numérico aumenta. Esto se debe a que la resta

$$f(x+h) - f(x)$$

se vuelve problemática. Si toma h lo suficientemente pequeño (después de todo, en teoría, más pequeño es mejor), entonces f(x+h) será igual a f(x) a la precisión de la máquina. Esto significa que todas las derivadas se calcularán como cero, sin importar la función, si solo toma h lo suficientemente pequeño. Aquí hay un ejemplo que calcula la derivada de sin(x) en x=1.

```
 \begin{array}{l} {\rm cout} << {\rm std::setprecision}(15); \\ {\rm for} \ ({\rm int} \ i=1; \ i<20; \ ++i) \\ {\rm double} \ h={\rm pow}(10.0, \ -i); \\ {\rm cout} << ({\rm sin}(1.0+h) \ - \ {\rm sin}(1.0))/h << \ ''\ ''; \\ {\rm } \\ {\rm cout} << "El \ verdadero \ resultado \ es: \ ''<< {\rm cos}(1.0) << \ ''\ 'n"; \\ \end{array}
```

Aquí está la salida del código anterior. Para que la salida sea más fácil de entender, los dígitos después del primer dígito incorrecto se han reemplazado por puntos.

antoniocarrillo@ciencias.unam.mx 752 Antonio Carrillo Ledesma, Et alii

```
0.544......

0.55......

0

0

0

0

El verdadero resultado es: 0.54030230586814
```

La precisión  $^{149}$  mejora a medida que h se hace más pequeña hasta que  $h=10^{-8}$ . Pasado ese punto, la precisión decae debido a la pérdida de precisión en la resta. Cuando  $h=10^{-16}$  o menor, la salida es exactamente cero porque sin(1.0+h) es igual a sin(1.0) a la precisión de la máquina. (De hecho, 1+h equivale a 1 a la precisión de la máquina. Más sobre eso a continuación).

¿Qué haces cuando tu problema requiere resta y va a causar una pérdida de precisión? A veces la pérdida de precisión no es un problema; los dobles comienzan con mucha precisión de sobra. Cuando la precisión es importante, a menudo es posible usar algún truco para cambiar el problema de modo que no requiera resta o no requiera la misma resta con la que comenzaste.

Los Números de Punto Flotante Tienen Rangos Finitos Todos saben que los números de punto flotante tienen rangos finitos, pero esta limitación puede aparecer de manera inesperada. Por ejemplo, puede encontrar sorprendente la salida de las siguientes líneas de código

float 
$$f = 16777216$$
;  
cout  $<< f << " " << f+1 << " \n";$ 

Este código imprime el valor 16777216 dos veces. ¿Que pasó? De acuerdo con la especificación IEEE para aritmética de punto flotante, un tipo flotante tiene 32 bits de ancho. Veinticuatro de estos bits están dedicados al significado (lo que solía llamarse la mantisa) y el resto al exponente. El número 16777216 es  $2^{24}$  y, por lo tanto, a la variable flotante f no le queda precisión para representar f+1. Ocurriría un fenómeno similar para  $2^{53}$  si f fuera

 $<sup>^{149}\</sup>mathrm{Los}$  resultados anteriores se calcularon con Visual C ++ 2008. Cuando se compiló con gcc 4.2.3 en Linux, los resultados fueron los mismos, excepto los últimos cuatro números. Donde VC ++ produjo ceros, gcc produjo números negativos: -0.017 ..., - 0.17 ..., -1.7 ... y 17 ....

del tipo double porque un doble de 64 bits dedica 53 bits al significado. El siguiente código imprime 0 en lugar de 1.

```
x = 9007199254740992; // 2^53

cout << ((x+1) - x) << "\n";
```

También podemos quedarnos sin precisión al agregar números pequeños a números de tamaño moderado. Por ejemplo, el siguiente código imprime "¡Lo siento!" porque DBL\_EPSILON (definido en float.h) es el número positivo más pequeño  $\epsilon$  tal que  $1 + \epsilon! = 1$  cuando se usan tipos dobles.

```
x = 1.0;

y = x + 0.5*DBL\_EPSILON;

if (x == y)

cout << "¡Lo siento!\n";
```

De manera similar, la constante FLT\_EPSILON es el número positivo más pequeño  $\epsilon$  tal que  $1+\epsilon$  no es 1 cuando se usan tipos flotantes.

¿Por qué los Flotantes IEEE Tienen Dos Ceros: +0 y -0? Aquí hay un detalle extraño de la aritmética de punto flotante IEEE: las computadoras tienen dos versiones de 0: cero positivo y cero negativo. La mayoría de las veces, la distinción entre +0 y -0 no importa, pero de vez en cuando las versiones firmadas del cero son útiles.

Si una cantidad positiva llega a cero, se convierte en +0. Y si una cantidad negativa llega a cero, se convierte en -0. Podría pensar en +0 (respectivamente, -0) como el patrón de bits para un número positivo (negativo) demasiado pequeño para representarlo.

El estándar de punto flotante IEEE dice que 1/+0 debería ser +infinito y 1/-0 debería ser -infinito. Esto tiene sentido si interpreta +/-0 como el fantasma de un número que se desbordó dejando solo su signo. El recíproco de un número positivo (negativo) demasiado pequeño para representarlo es un número positivo (negativo) demasiado grande para representarlo.

Para demostrar esto, ejecute el siguiente código C:

```
int main() { double x = 1e-200;
```

antoniocarrillo@ciencias.unam.mx 754 Antonio Carrillo Ledesma, Et alii

```
double y = 1e-200 * x;
          printf("Reciproco de +0: %g\n", 1/y);
          y = -1e-200*x;
          printf("Reciproco de -0: \%g\n", 1/y);
     }
En Linux con gcc, la salida es:
```

Reciproco de +0: inf

Reciproco de -0: -inf

Sin embargo, hay algo acerca de los ceros firmados y las excepciones que no tiene sentido. El informe acertadamente denominado "Lo que todo informático debería saber sobre la aritmética de punto flotante" tiene lo siguiente que decir sobre los ceros con signo.

En aritmética IEEE, es natural definir  $log0 = -\infty$  y logxcomo un NaN cuando x < 0. Suponga que x representa un pequeño número negativo que se ha desbordado a cero. Gracias al cero con signo, x será negativo, por lo que loq puede devolver un NaN. Sin embargo, si no hubiera un cero con signo, la función logarítmica no podría distinguir un número negativo subdesbordado de 0 y, por lo tanto, tendría que devolver  $-\infty$ .

Esto implica que log(-0) debe ser NaN y log(+0) debe ser  $-\infty$ . Eso tiene sentido, pero eso no es lo que sucede en la práctica. La función de log devuelve  $-\infty$  para +0 y -0.

Ejecuté el siguiente código en C:

```
int main()
{
    double x = 1e-200;
    double y = 1e-200 * x;
    printf("Log de +0: %g\n", log(y));
    y = -1e-200*x;
    printf("Log de -0: \%g\n", log(y));
}
```

En Linux, el código imprime:

```
Log de +0: -inf
Log de -0: -inf
```

## Use Logaritmos para Evitar Desbordamiento y Subdesbordamiento

Las limitaciones de los números de punto flotante descritos en la sección anterior provienen de tener un número limitado de bits en el significado. El desbordamiento y el subdesbordamiento resultan de tener también un número finito de bits en el exponente. Algunos números son demasiado grandes o demasiado pequeños para almacenarlos en un número de punto flotante.

Muchos problemas parecen requerir calcular un número de tamaño moderado como la razón de dos números enormes. El resultado final puede ser representable como un número de punto flotante aunque los resultados intermedios no lo sean. En este caso, los logaritmos proporcionan una salida. Si desea calcular M/N para grandes números M y N, calcule log(M) - log(N) y aplique exp() al resultado. Por ejemplo, las probabilidades a menudo implican proporciones de factoriales, y los factoriales se vuelven astronómicamente grandes rápidamente. Para  $N>170,\ N!$  es mayor que DBL\_MAX, el número más grande que puede representarse por un doble (sin precisión extendida). Pero es posible evaluar expresiones como 200!/(190!10!) Sin desbordamiento de la siguiente manera:

```
x = \exp(\ \log Factorial(200) - \log Factorial(190) - \log Factorial(10));
```

Una función logFactorial simple pero ineficiente podría escribirse de la siguiente manera:

```
double logFactorial(int n) 
 { double sum = 0.0; 
 for (int i = 2; i <= n; ++i) sum += log((double)i); 
 return sum; 
}
```

Un mejor enfoque sería utilizar una función de registro gamma si hay una disponible. Consulte Cómo calcular las probabilidades binomiales para obtener más información.

Las Operaciones Numéricas no Siempre Devuelven Números Debido a que los números de punto flotante tienen sus limitaciones, a veces las operaciones de punto flotante devuelven "infinito" como una forma de decir "el resultado es más grande de lo que puedo manejar". Por ejemplo, el siguiente código imprime 1. # *INF* en Windows e inf en Linux.

antoniocarrillo@ciencias.unam.mx 756 Antonio Carrillo Ledesma, Et alii

$$x = DBL\_MAX;$$
 $cout << 2*x << "\n";$ 

A veces, la barrera para devolver un resultado significativo tiene que ver con la lógica en lugar de la precisión finita. Los tipos de datos de punto flotante representan números reales (a diferencia de los números complejos) y no hay un número real cuyo cuadrado sea -1. Eso significa que no hay un número significativo que devolver si el código solicita sqrt(-2), incluso con una precisión infinita. En este caso, las operaciones de punto flotante devuelven NaN. Estos son valores de punto flotante que representan códigos de error en lugar de números. Los valores NaN se muestran como 1. # IND en Windows y NAN en Linux.

Una vez que una cadena de operaciones encuentra un NaN, todo es un NaN de ahí en adelante. Por ejemplo, suponga que tiene un código que equivale a algo como lo siguiente:

if 
$$(x - x == 0)$$
  
// hacer algo

¿Qué podría impedir que se ejecute el código que sigue a la instrucción if? Si x es un NaN, entonces también lo es x-x y los NaN no equivalen a nada. De hecho, los NaN ni siquiera se igualan. Eso significa que la expresión x == x se puede usar para probar si x es un número (posiblemente infinito). Para obtener más información sobre infinitos y NaN, consulte las excepciones de punto flotante IEEE en C ++.

Trabajando con Factoriales Los cálculos de probabilidad a menudo implican tomar la razón de números muy grandes para producir un número de tamaño moderado. El resultado final puede caber dentro de un doble con espacio de sobra, pero los resultados intermedios se desbordarían. Por ejemplo, suponga que necesita calcular el número de formas de seleccionar 10 objetos de un conjunto de 200. ¡Esto es 200!/(190!10!), Aproximadamente 2.2e16. Pero 200! y 190! desbordaría el rango de un doble.

Hay dos formas de solucionar este problema. Ambos usan la siguiente regla: Use trucos algebraicos para evitar el desbordamiento.

El primer truco es reconocer que

$$200! = 200 * 199 * 198 * ... * 191 * 190!$$

antoniocarrillo@ciencias.unam.mx 757 Antonio Carrillo Ledesma, Et alii

y así

$$200!/(190!10!) = 200 * 199 * 198 * ... * 191/10!$$

esto ciertamente funciona, pero está limitado a factoriales.

Una técnica más general es usar logaritmos para evitar el desbordamiento: tome el logaritmo de la expresión que desea evaluar y luego exponga el resultado. En este ejemplo

$$log(200!/(190!10!)) = Log(200!) - log(190!) - log(10!)$$

si tiene un código que calcula el logaritmo de los factoriales directamente sin calcular primero los factoriales, puede usarlo para encontrar el logaritmo del resultado que desea y luego aplicar la función exp.

Calcular Inverso del Factorial Dado un número positivo x, ¿cómo se puede encontrar un número n tal que n! = x, o tal que  $n! \approx x$ ?. El Software matemático tiende a trabajar con la función gamma en lugar de factoriales porque  $\Gamma(x)$  extiende x! en números reales, con la relación  $\Gamma(x+1) = x!$ . El lado izquierdo a menudo se toma como una definición del lado derecho cuando x no es un número entero positivo.

No sólo preferimos trabajar con la función gamma, sino que es más fácil trabajar con el logaritmo de la función gamma para evitar el desbordamiento.

El siguiente código Python encuentra el inverso del logaritmo de la función gamma utilizando el método de bisección. Este método requiere un límite superior e inferior. Si solo pasamos valores positivos, 0 sirve como límite inferior. Podemos encontrar un límite superior probando potencias de 2 hasta obtener algo lo suficientemente grande.

```
from scipy.special import gammaln from scipy.optimize import bisect import math as mt def inverse_log_gamma(logarg): assert(logarg > 0) a = 1 b = 2 while b < logarg: b = b*2 return bisect(lambda x: gammaln(x) - logarg, a, b) def inverse_factorial(logarg):
```

antoniocarrillo@ciencias.unam.mx 758 Antonio Carrillo Ledesma, Et alii

Tenga en cuenta que la función inverse\_factorial toma el logaritmo del valor cuyo inverso del factorial desea calcular. Por ejemplo:

```
print(inverse_factorial(mt.log(mt.factorial(42))))) regresa 42.
```

Trabajar en la escala logarítmica nos permite trabajar con números mucho mayores. El factorial de 171 es mayor que el mayor número de doble precisión de IEEE, por lo que inverse\_factorial no podría devolver ningún valor mayor que 171 si pasáramos x en lugar de log x. Pero al tomar log x como argumento, podemos calcular factoriales inversos de números tan grandes que el factorial se desborda.

Por ejemplo, supongamos que queremos encontrar el valor de m tal que m! es el factorial más cercano a  $\sqrt{(2024!)}$ . Podemos usar el código.

$$print(inverse\_factorial(gammaln(2025)/2))$$

para encontrar m=1112 aunque 1112! es del orden de  $10^{2906}$ , mucho mayor que el mayor número de punto flotante representable, que es del orden de  $10^{308}$ .

Cuando n! = x no tiene una solución exacta, existe alguna opción sobre qué valor de n devolver como inverso del factorial de x. El código anterior minimiza la distancia desde n! a x en una escala logarítmica. Es posible que desee modificar el código para minimizar la distancia en una escala lineal o encontrar la n más grande con n! < x, o la n más pequeña con n! > x dependiendo de su aplicación.

$$u(z) = \frac{e^z - 1 - z}{z^2}$$

para valores pequeños de z, digamos  $z=10^{-8}$ .

El código Python

<sup>&</sup>lt;sup>150</sup>Este ejemplo proviene de Lloyd N. Trefethen and J. A. C. Weideman. The Exponentially Convergent Trapezoid Rule. SIAM Review. Vol. 56, No. 3. pp. 385-458.

```
import numpy as np def f(z): return (np.exp(z) - 1 - z)/z**2 print(f(1e-8))
```

imprime:

-0.607747099184471.

Ahora suponga que sospecha dificultades numéricas y calcula su resultado con 50 decimales usando bc - $l^{151}$  usando:

imprime:

## .50000001666666670833333341666666600000000000000000

Esto sugiere que el cálculo original era completamente erróneo. ¿Qué está sucediendo?

Para z pequeña,

$$e^z \approx 1 + z$$

y así perdemos precisión al evaluar directamente el numerador en la definición de u (en nuestro ejemplo, perdimos toda precisión).

El teorema del valor medio del análisis complejo dice que el valor de una función analítica en un punto es igual al promedio continuo de los valores sobre un círculo centrado en ese punto. Si aproximamos este promedio tomando el promedio de 16 valores en un círculo de radio 1 alrededor de nuestro punto, obtenemos una precisión total. El código Python

```
\begin{aligned} \text{def g(z):} & N = 16 \\ & \text{ws} = \text{z} + \text{np.exp(2j*np.pi*np.arange(N)/N)} \\ & \text{return sum(f(ws))/N} \\ & \text{print(g(1e-8))} \end{aligned}
```

<sup>&</sup>lt;sup>151</sup>El comando be (Basic Calculator) de la línea de comandos de Linux/Unix, es una calculadora con grandes posibilidades de uso.

imprime:

```
0.50000000166666668 + 8.673617379884035e-19j
```

que se aparta del resultado calculado con bc en el decimosexto decimal. En un nivel alto, evitamos dificultades numéricas al promediar puntos alejados de la región difícil.

**Logit inverso** A continuación, veamos el cálculo  $f(x) = e^x/(1+e^x)$ . (Los estadísticos llaman a esto la función "logit inverso" porque es el inverso de la función que ellos llaman la función "logit".) El enfoque más directo sería calcular exp(x)/(1+exp(x)). Veamos dónde se puede romper eso.

```
double x = 1000;
double t = \exp(x);
double y = t/(1.0 + t);
```

Imprimir y da -1. # IND. Esto se debe a que el cálculo de t se desbordó, produciendo un número mayor que el que cabía en un doble. Pero podemos calcular el valor de y fácilmente. Si t es tan grande que no podemos almacenarlo, entonces 1+t es esencialmente lo mismo que t y la relación es muy cercana a 1. Esto sugiere que descubramos qué valores de x son tan grandes que f(x) será igual 1 a la precisión de la máquina, luego solo devuelva 1 para esos valores de x para evitar la posibilidad de desbordamiento. Esta es nuestra siguiente regla: No calcules un resultado que puedas predecir con precisión.

El archivo de encabezado float.h tiene un  $DBL\_EPSILON$  constante, que es la precisión doble más pequeña que podemos agregar a 1 sin recuperar 1. Un poco de álgebra muestra que si x es más grande que  $-log(DBL\_EPSILON)$ , entonces f(x) será igual a 1 a la precisión de la máquina. Así que aquí hay un fragmento de código para calcular f(x) para valores grandes de x:

```
const double x_max = -log(DBL_EPSILON); if (x > x_max) return 1.0;
```

El código provisto en esta sección calcula siete funciones que aparecen en las estadísticas. Cada uno evita problemas de desbordamiento, subdesbordamiento o pérdida de precisión que podrían ocurrir para grandes argumentos negativos, grandes argumentos positivos o argumentos cercanos a cero:

antoniocarrillo@ciencias.unam.mx 761 Antonio Carrillo Ledesma, Et alii

- LogOnePlusX calcúlese como log(1+x) como en ejemplo antes visto
- ExpMinusOne calcúlese como  $e^x 1$
- Logit calcúlese como log(x/(1-x))
- Logit Inverse calcúlese como  $e^x/(1+e^x)$  como se discutió en el ejemplo último
- LogLogitInverse calcúlese como  $log(e^x/(1+e^x))$
- LogitInverseDifference calcúlese como LogitInverse(x)-LogitInverse(y)
- LogOnePlusExpX calcúlese como log(1 + exp(x))
- ComplementaryLogLog calcúlese como log(-log(1-x))
- ComplementaryLogLogInverse calcúlese como 1.0 exp(-exp(x))

Las soluciones presentadas aquí parecen innecesarias o incluso incorrectas al principio. Si este tipo de código no está bien comentado, alguien lo verá y lo "simplificará" incorrectamente. Estarán orgullosos de todo el desorden innecesario que eliminaron. Y si no prueban valores extremos, su nuevo código parecerá funcionar correctamente. Las respuestas incorrectas y los NaN solo aparecerán más tarde.

**Log**  $(1 + \mathbf{x})$  Ahora veamos el ejemplo del cálculo de log(x+1). Considere el siguiente código:

```
double x = 1e-16;
double y = log(1 + x)/x;
```

En este código y=0, aunque el valor correcto sea igual a 1 para la precisión de la máquina.

¿Qué salió mal? los números de doble precisión tienen una precisión de aproximadamente 15 decimales, por lo que 1+x equivale a 1 para la precisión de la máquina. El registro de 1 es cero, por lo que y se establece en cero. Pero para valores pequeños de x, log(1+x) es aproximadamente x, por lo que log(1+x)/x es aproximadamente 1. Eso significa que el código anterior para calcular log(1+x)/x devuelve un resultado con 100% de error relativo. Si x no es tan pequeño que 1+x es igual a 1 en la máquina, aún podemos tener problemas. Si x es moderadamente pequeño, los bits en x no se pierden totalmente al calcular 1+x, pero algunos sí. Cuanto más se acerca x a 0, más bits se pierden. Podemos usar la siguiente regla: Utilice aproximaciones analíticas para evitar la pérdida de precisión.

La forma favorita de aproximación de los analistas numéricos es la serie de potencia. ¡La serie de potencia para

$$log(1+x) = x + x^2/2! + x^3/3! + \dots$$

Para valores pequeños de x, simplemente devolver x para log(1+x) es una mejora. Esto funcionaría bien para los valores más pequeños de x, pero para algunos valores no tan pequeños, esto no será lo suficientemente preciso, pero tampoco lo hará directamente el cálculo del log(1+x).

La Precisión Arbitraria no es una Panacea Tener acceso a aritmética de precisión arbitraria no necesariamente hace que desaparezcan las dificultades de cálculo numérico. Aún necesitas saber lo que estás haciendo. Antes de ampliar la precisión, hay que saber hasta dónde ampliarla. Si no lo amplía lo suficiente, sus respuestas no serán lo suficientemente precisas. Si lo extiendes demasiado, desperdicias recursos. Tal vez esté bien desperdiciar recursos, por lo que extiende la precisión más de lo necesario. ¿Pero hasta qué punto es más de lo necesario?

Como ejemplo, considere el problema de encontrar valores de n tales que tan(n) > n, uno de los valores es:

$$k = 1428599129020608582548671.$$

Verifiquemos que tan(k) > k. Usaremos bc porque admite precisión arbitraria. Nuestra k es un número de 25 dígitos, así que digamos a bc que queremos trabajar con 30 decimales para tener un poco de espacio<sup>152</sup>. bc no

 $<sup>^{152}</sup>bc$ automáticamente le otorga un poco más de espacio del que solicita, pero le pediremos aún más

tiene una función tangente, pero sí la tiene s() para seno y c() para coseno, por lo que calcularemos la tangente como seno sobre coseno.

```
bc -l

escala = 30

k = 1428599129020608582548671

s(k)/c(k) - k
```

Esperamos que esto devuelva un valor positivo, pero en su lugar devuelve

```
-1428599362980017942210629.31...
```

Entonces, ¿está mal la hipótesis? ¿O es bc ignorar nuestra solicitud? Resulta que ninguna de las dos cosas es cierta.

No se puede calcular directamente la tangente de un número grande. Utiliza la reducción de rango para reducirlo al problema de calcular la tangente de un ángulo pequeño donde funcionarán sus algoritmos. Dado que la tangente tiene un período  $\pi$ , reducimos  $k \mod \pi$  calculando  $k-[k/\pi]\pi$ . Es decir, restamos tantos múltiplos de  $\pi$  como podamos hasta obtener un número entre 0 y  $\pi$ . Volviendo a bc, calculamos

$$pi = 4*a(1)$$
 k/pi

esto regresa

454737226160812402842656.500000507033221370444866152761

y entonces calculamos la tangente de

```
\begin{array}{l} t = 0,\!500000507033221370444866152761 * pi \\ = 1,\!570797919686740002588270178941 \end{array}
```

Como t es ligeramente mayor que  $\pi/2$ , la tangente será negativa. No podemos tener tan(t) mayor que k porque tan(t) ni siquiera es mayor que 0. Entonces, ¿dónde se estropearon las cosas?.

El cálculo de  $\pi$  tuvo una precisión de 30 cifras significativas y el cálculo de  $k/\pi$  tuvo una precisión de 30 cifras significativas, dado nuestro valor de  $\pi$ . Hasta ahora ha funcionado según lo prometido.

antoniocarrillo@ciencias.unam.mx 764 Antonio Carrillo Ledesma, Et alii

El valor calculado de  $k/\pi$  tiene una precisión de 29 cifras significativas, 23 antes del decimal y 6 después. Entonces, cuando tomamos la parte fraccionaria, solo tenemos seis cifras significativas y eso no es suficiente. Ahí es donde las cosas van mal. Obtenemos un valor  $[k/\pi]$  que es mayor que 0.5 en el séptimo decimal, mientras que el valor exacto es menor que 0.5 en el vigésimo quinto decimal. Necesitábamos 25-6=19 cifras más significativas.

Ésta es la principal dificultad del cálculo en punto flotante: restar números casi iguales pierde precisión. Si dos números coinciden en m decimales, puedes esperar perder m cifras significativas en la resta. El error en la resta será pequeño en relación con los datos de entrada, pero no pequeño en relación con el resultado.

Observe que calculamos  $k/\pi$  hasta 29 cifras significativas y, dado ese resultado, calculamos la parte fraccionaria exactamente . Calculamos con precisión  $[k/\pi]\pi$ , pero perdimos precisión cuando calculamos y restamos ese valor de k.

Nuestro error al calcular  $k-[k/\pi]\pi$  fue pequeño en relación con k, pero no en relación con el resultado final. Nuestra k es del orden de  $10^{25}$  y el error en nuestra resta es del orden de  $10^{-7}$ , pero el resultado es del orden de 1. No hay ningún error en bc. Llevó a cabo todos los cálculos con la precisión anunciada, pero no tuvo suficiente precisión de trabajo para producir el resultado que necesitábamos.

## 15.3 Aritmética de Baja Precisión

La popularidad de la aritmética de baja precisión para el cómputo de alto rendimiento se ha disparado desde el lanzamiento en 2017 de la GPU Nvidia Volta. Los núcleos tensores de media precisión de Volta ofrecieron una enorme ganancia de rendimiento 16 veces mayor que la doble precisión para operaciones clave. Y el rendimiento del Hardware está mejorando aún más: el FP16 con núcleo tensor Nvidia H100 es 58 veces más rápido que el FP64 estándar.

Esta sorprendente aceleración ciertamente llama la atención. Sin embargo, en el cálculo científico, la aritmética de baja precisión suele considerarse insegura para los códigos de modelado y simulación. De hecho, a veces se puede aprovechar una precisión más baja, comúnmente en una configuración de "precisión mixta" en la que sólo partes del cálculo se realizan con baja precisión. Sin embargo, en general, cualquier precisión menor que el doble se considera inadecuada para modelar fenómenos físicos complejos

con fidelidad.

En respuesta, los desarrolladores han creado herramientas para medir la seguridad de la aritmética de precisión reducida en códigos de aplicación. Algunas herramientas pueden incluso identificar qué variables o matrices se pueden reducir de forma segura a una precisión menor sin perder precisión en el resultado final. Sin embargo, el uso de estas herramientas a ciegas, sin el respaldo de algún tipo de proceso de razonamiento, puede resultar peligroso.

Un ejemplo ilustrará esto:

El método del gradiente conjugado para la resolución y optimización de sistemas lineales y el método de Lanczos, estrechamente relacionado, para la resolución de problemas de valores propios mostraron una gran promesa tras su invención a principios de los años cincuenta. Sin embargo, se consideraban inseguros debido a errores de redondeo catastróficos en la aritmética de punto flotante, que son aún más pronunciados a medida que se reduce la precisión del punto flotante.

No obstante, Chris Paige demostró en su trabajo pionero en la década de 1970 que el error de redondeo, aunque sustancial, no excluye la utilidad de los métodos cuando se utilizan correctamente. El método del gradiente conjugado se ha convertido en un pilar del cómputo científico.

Teniendo en cuenta que ninguna herramienta podría llegar a este hallazgo sin un cuidadoso análisis matemático de los métodos. Una herramienta detectaría inexactitudes en el cálculo pero no podría certificar que estos errores no puedan perjudicar el resultado final.

Algunos podrían proponer en cambio un enfoque puramente basado en datos: simplemente pruebe con baja precisión en algunos casos de prueba; si funciona, utilice baja precisión en producción. Sin embargo, este enfoque está lleno de peligros: es posible que los casos de prueba no capturen todas las situaciones que podrían encontrarse en producción.

Por ejemplo, uno podría probar un código de aerodinámica sólo en regímenes de flujo suaves, pero las series de producción pueden encontrar flujos complejos con gradientes pronunciados, que la aritmética de baja precisión no puede modelar correctamente. Los artículos académicos que prueban métodos y herramientas de baja precisión deben evaluarse rigurosamente en escenarios desafiantes del mundo real como este.

Lamentablemente, los equipos de ciencia computacional frecuentemente no tienen tiempo para evaluar sus códigos para un uso potencial de aritmética de menor precisión. Las herramientas ciertamente podrían ayudar. Además, las bibliotecas que encapsulan métodos de precisión mixta pueden ofrecer beneficios a muchos usuarios. Una gran historia de éxito son los solucionadores lineales densos de precisión mixta, basados en el sólido trabajo teórico de Nick Highnam y sus colegas, que han llegado a bibliotecas como: Lu, Hao; Matheson, Michael; Wang, Feiyi; Joubert, Wayne; Ellis, Austin; Oles, Vladyslav. "OpenMxP-Opensource Mixed Precision Computing,".

Entonces la respuesta final es "depende". Cada nuevo caso debe examinarse cuidadosamente y tomar una decisión basada en alguna combinación de análisis y pruebas.

Sí bien, NVIDIA lidera el mercado de las GPU para inteligencia artificial (IA) con una cuota de mercado aproximada del 80%, pero no es en absoluto la única empresa que tiene en su porfolio chips de vanguardia para IA. La compañía californiana Cerebras posee, de hecho, los procesadores para este escenario de uso más complejos que existen. Su chip WSE-2, por ejemplo, aglutina nada menos que 2.6 billones de transistores contabilizados en la escala numérica larga y 850,000 núcleos optimizados para IA.

Cerebras entrega a sus clientes estos procesadores integrados en una plataforma para IA conocida como CS-2, y precisamente uno de ellos es la compañía de Emiratos Árabes G42. Esta última está construyendo seis superordenadores para IA capaces de superar la barrera de la exaescala que aglutinan una gran cantidad de sistemas CS-2. Y según la CIA algunas de estas máquinas irán a parar a las grandes tecnológicas chinas. No obstante, esto no es todo. Y es que Cerebras dió a conocer en 2024 un procesador para IA aún más potente que su WSE-2.

El procesador WSE-3 Cerebras ya tiene listo su procesador WSE-3 (Wafer Scale Engine 3), un producto que, como podemos intuir, está llamado a suceder al también ambicioso WSE-2.

Ambos procesadores se fabrican a partir de una oblea completa de silicio, lo que permite a Cerebras integrar muchos más bloques funcionales y núcleos en la lógica que una GPU convencional como las que fabrican NVIDIA, AMD o Huawei.

Y es que aglutina 4 billones de transistores, tiene una superficie de 46, 225 mm<sup>2</sup>, integra nada menos que 900,000 núcleos optimizados para IA y tiene

antoniocarrillo@ciencias.unam.mx 767 Antonio Carrillo Ledesma, Et alii

una potencia de cálculo, según Cerebras, de 125 petaflops.

Según Cerebras su procesador WSE-3 es el doble de potente que el WSE-2. De hecho, de acuerdo con las especificaciones que ha publicado rinde como 62 GPU H100 de NVIDIA trabajando al unísono, y no debemos pasar por alto que este procesador de la compañía liderada por Jensen Huang es el más potente que tiene hasta que se produzca el lanzamiento de la GPU H200.

Sea como sea Cerebras entrega sus procesadores WSE-3 integrados en un superordenador conocido como CS-3 que es capaz de entrenar grandes modelos de IA con hasta 24 billones de parámetros. El mapa de memoria externa de este superordenador oscila entre 1.5 TB y 1.2 PB, un espacio de almacenamiento descomunal que permite almacenar modelos de lenguaje masivos en un único espacio lógico.

Según informan, el chip WSE-3 optimizado para la IA es capaz de entrenar hasta 24,000 millones de parámetros, lo que también equivaldría a un rendimiento máximo de IA de 125 petaflops.