
Lezione 7

Standard IEEE-754

per le operazioni floating-point

<http://www.dii.unisi.it/~giorgi/didattica/arc1>

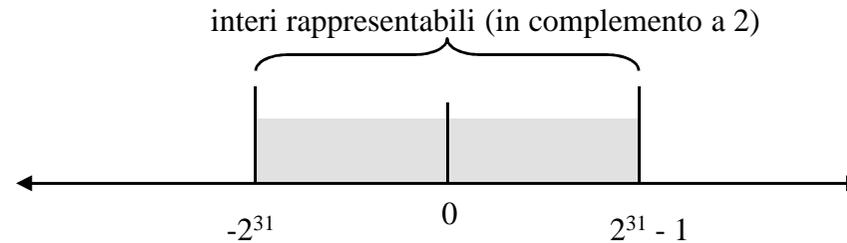
Floating-Point (breve schema)

- Nelle operazioni floating point dobbiamo rappresentare
 - numero con decimali, e.g., 3.1416
 - numeri molto piccoli, e.g., .000000001
 - numeri molto grandi, e.g., 3.15576×10^9
- Rappresentazione:
 - segno, mantissa, esponente: $(-1)^{\text{segno}} \times \text{mantissa} \times 2^{\text{esponente}}$
 - piu' bit di mantissa danno piu' accuratezza
 - piu' bit di esponente danno un intervallo piu' ampio
- IEEE 754 floating point standard:
 - singola precisione ("float" del C) → uso 32 bit per memorizzare:
1 bit segno, 8 bit esponente, 23 bit mantissa**
 - doppia precisione ("double" del C) → uso 64 bit per memorizzare:
1 bit segno, 11 bit esponente, 52 bit mantissa**

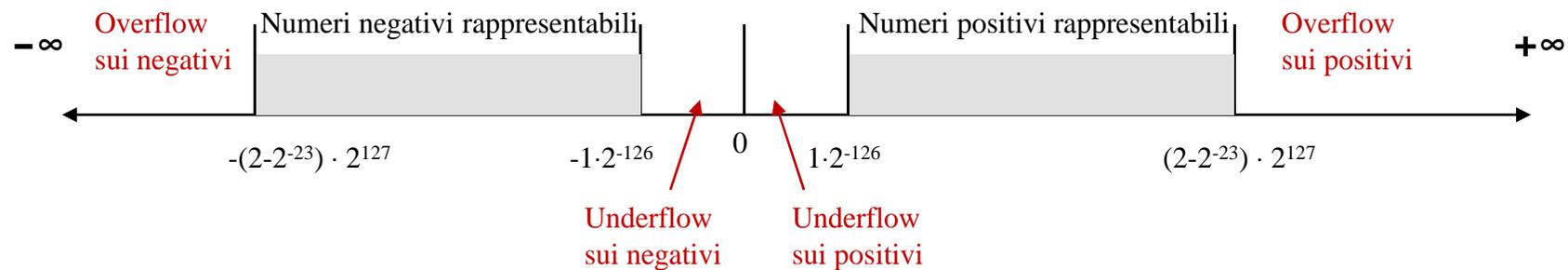
**Nota: utilizzando il "trucco" della normalizzazione della mantissa nella forma 1.xxxxx le cifre (binarie) effettivamente rappresentate sono rispettivamente 24 e 53

Intervalli rappresentabili

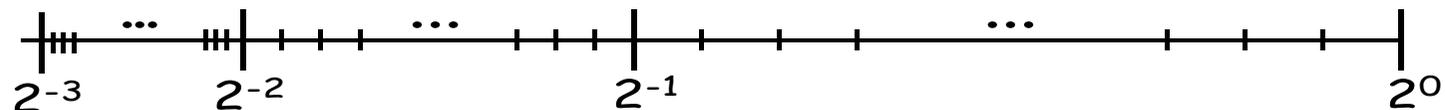
- Nel caso di interi a 32 bit:



- Nel caso di floating-point a 32 bit (IEEE-754)



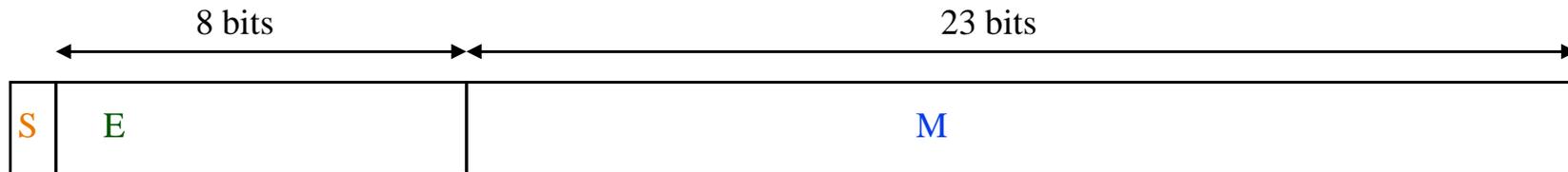
Nota: il numero di valori rappresentati con 32 bit e' circa lo stesso nel caso INT e in quello FP... la densita' dei punti non e' costante, quindi.



Esempio

- Rappresentare -0.75 in formato IEEE single precision
 - rappresentazione decimale: $-0.75 = -3/4 = -3 \times 2^{-2}$
 - rappresentazione binaria: $-11 \times 2^{-2} = -1.1 \times 2^{-1}$
 - rappresentazione binaria normalizzata: -1.1×2^{-1}
 - segno=negativo
→ $S=1$
 - mantissa_normalizzata=1.1000 0000 0000 0000 0000 000
→ $M = 1000\ 0000\ 0000\ 0000\ 0000\ 000$
 - esponente=-1
→ $E = \text{esponente} + \text{polarizzazione} = -1 + 127 = 126 \rightarrow 01111110$
- IEEE single precision: $10111111010000000000000000000000$
- Nota: la mantissa e' stata normalizzata. Per questo tutti i bit piu' significativi risultano spostati verso SINISTRA.

Floating-Point: single precision (32 bits)



Esempi con mantissa ed esponente positivi/negativi:

$+1.1010001 \cdot 2^{+10100} \rightarrow 0 \ 10010011 \ 101000100000000000000000$
 $-1.1010001 \cdot 2^{+10100} \rightarrow 1 \ 10010011 \ 101000100000000000000000$
 $+1.1010001 \cdot 2^{-10100} \rightarrow 0 \ 01101011 \ 101000100000000000000000$
 $-1.1010001 \cdot 2^{-10100} \rightarrow 1 \ 01101011 \ 101000100000000000000000$

Esempi di numeri "notevoli":

$0 \rightarrow 0 \ 00000000 \ 000000000000000000000000$
 $1 \rightarrow 0 \ 01111111 \ 000000000000000000000000$
minore rappresentabile (2^{-126}) $\rightarrow 0 \ 00000001 \ 000000000000000000000000$
maggiore rappres. ($(2-2^{-23}) \cdot 2^{127}$) $\rightarrow 0 \ 11111110 \ 111111111111111111111111$
 $\infty \rightarrow 0 \ 11111111 \ 000000000000000000000000$
NaN $\rightarrow 0 \ 11111111 \ \text{XXXXXXXXXXXXXXXXXXXXXXXXXX}$

Tabella di riepilogo codifica FP nei processori INTEL

Da: Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 1: Basic Architecture, Sep. 2008.

URL: <http://download.intel.com/design/processor/manuals/253665.pdf>

Class		Sign	Biased Exponent	Significand	
				Integer	Fraction
Positive	+ ∞	0	11..11	1	00..00
	+ Normals	0	11..10	1	11..11
	
		0	00..01	1	00..00
	+ Denormals	0	00..00	0	11..11
.		.	.	.	
0		00..00	0	00..01	
+ Zero	0	00..00	0	00..00	
Negative	- Zero	1	00..00	0	00..00
	- Denormals	1	00..00	0	00..01
	
		1	00..00	0	11..11
	- Normals	1	00..01	1	00..00
	
1		11..10	1	11..11	
- ∞	1	11..11	1	00..00	
NaNs	SNaN	X	11..11	1	0X..XX ²
	QNaN	X	11..11	1	1X..XX
	QNaN	1	11..11	1	10..00
	Floating-Point Indefinite				

Single-Precision ← 8 Bits → ← 23 Bits →

Double-Precision ← 11 Bits → ← 52 Bits →

Double Extended-Precision ← 15 Bits → ← 63 Bits →

Note:

- (1) il bit della parte intera e' implicito e non viene quindi memorizzato
- (2) La parte frazionaria di SNaN deve essere diversa da zero con il bit piu' significativo pari a 0

SNaN= Signaling NaN (solleva eccezione)

QNaN=Quiet NaN (non solleva eccezione)

QNaN-FP= "INDETERMINATO" → nessuno dei due operandi era un NaN: lo ha generato l'operazione

I numeri denormalizzati possono essere generati ma non vengono memorizzati

E' prevista una modalita' DAZ =Denormal-Are-Zero

E' supportata una modalita' a 80 bit (non IEEE-754)

Complessita' del Floating Point

- Le operazioni sono leggermente piu' complesse (v. testo)
- Oltre alla condizione di overflow si puo' avere "underflow"
 - (codifica: esponente con bit tutti a zero) → numeri denormalizzati
- L'accuratezza puo' essere un grosso problema
 - IEEE 754 prevede 2 ulteriori bit: 1 *guard* bit e 1 *round* bit
 - 4 modalita' per effettuare l'arrotondamento
(*round to $+\infty$* , *round to $-\infty$* , *round to 0*, *round to nearest representable*)
 - altre complessita' (es. numeri denormalizzati)
- Implementare lo standard puo' essere "ingannevole"
 - Non usare lo standard puo' essere anche peggio
 - vedere il testo per la descrizione dell'80x86 e del Pentium bug!

Esempi di arrotondamento

numero	r.to 0	r.to nearest	r.to +inf	r.to -inf
1.00000003	1.0	1.0	1.00000012	1.0
1.00000007	1.0	1.00000012	1.00000012	1.0
-1.00000003	-1.0	-1.0	-1.0	-1.00000012
-1.00000007	-1.0	-1.00000012	-1.0	-1.00000012

Ricapitolazione

- L'aritmetica del calcolatore e' limitata da precisione finita
- I pattern di bit non hanno particolare significato ma esistono standard per:
 - complemento a due
 - IEEE 754 floating point
- Le istruzioni determinano il "significato" dei pattern di bit
- Le prestazioni e l'accuratezza sono importanti: ci sono diverse complessita' nelle macchine reali
- Lo standard e' stato rivisto e il documento pubblicato nel Giugno 2008 come IEEE 754-2008
<http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=4610935&isnumber=4610934>
 - Fra le novita' il supporto per floating point a 128-bit
 - Aritmetica floating point decimale (oltre che binaria)

EPSILON di macchina (MACHEPS = ϵ_{mach})

- E' la differenza fra 1 e "1+"
(1+ \equiv il piu' piccolo numero rappresentabile maggiore di 1)
 - E' anche chiamato "precisione di macchina"
 - Da' una maggiorazione dell'errore di arrotondamento

• Esempio, IEEE-754 Single-Precision

$$\begin{aligned} 1 &\equiv (0 \text{ 01111111 } 000 \text{ 0000 } 0000 \text{ 0000 } 0000 \text{ 0000})_{\text{IEEE-754-SP}} \\ &= (1.000 \text{ 0000 } 0000 \text{ 0000 } 0000 \text{ 0000})_2 \\ &= (1)_{10} \end{aligned}$$

$$\begin{aligned} 1+ &\equiv (0 \text{ 01111111 } 000 \text{ 0000 } 0000 \text{ 0000 } 0000 \text{ 0001})_{\text{IEEE-754-SP}} \\ &= (1.000 \text{ 0000 } 0000 \text{ 0000 } 0000 \text{ 0001})_2 \\ &\approx (1.000000119209289\dots)_{10} \end{aligned}$$

$$\begin{aligned} \rightarrow \epsilon_{mach} &= (0 \text{ 01101000 } 000 \text{ 0000 } 0000 \text{ 0000 } 0000 \text{ 0000})_{\text{IEEE-754-SP}} \\ &= (0.000 \text{ 0000 } 0000 \text{ 0000 } 0000 \text{ 0001})_2 \\ &= (2^{-23})_{10} \approx (0.000000119209289\dots)_{10} \approx (1.19209289\dots \cdot 10^{-7})_{10} \end{aligned}$$

- In generale, per un tipo floating-point rappresentato in base b con p cifre di mantissa (normalizzato con una cifra $\neq 0$ prima della virgola), ovvero con $p-1$ cifre dopo la virgola

$$\rightarrow \epsilon_{mach} = b^{-(p-1)} = b^{1-p} \quad (\text{es. nel caso IEEE-754 single-precision } p=24)$$

Precisazioni su MACHEPS

- Alcuni usano come definizione: "il piu' piccolo numero positivo rappresentabile, tale che aggiunto ad 1 lo rende diverso da 1"
 - tale definizione e' diversa dalla precedente perche' il valore viene a dipendere dalla modalita' di arrotondamento usata
- Nel caso di IEEE-754 Single-Precision:
 - Per arrotondamenti al successivo (round to even) questa e' $2^{-24} + 2^{-47}$ (poco piu' di meta' rispetto alla precedente definizione)
 - Per arrotondamenti verso $+\infty$ questo e' il piu' piccolo numero rappresentabile, ovvero 2^{-149} (un numero denormalizzato)
 - Per arrotondamenti verso $-\infty$ questo coincide con la precedente definizione
- Il seguente programma C non determina MACHEPS bensì un numero distante non più di un fattore 2 da MACHEPS, usando una ricerca lineare

Calcolo di MACHEPS

```
#include <stdio.h>

int main( int argc, char **argv )
{
    float machEps = 1.0f;

    printf( "current Epsilon, 1 + current Epsilon\n" );
    do {
        printf( "%G\t%.20f\n", machEps, (1.0f + machEps) );
        machEps /= 2.0f;
        /* Se sommando l'epsilon successivo ho 1 --> ho trovato epsilon */
    }
    while ((float)(1.0 + (machEps/2.0)) != 1.0);

    printf( "\nCalculated Machine epsilon: %G\n", machEps );
    return 0;
}
```

MACHEPS Output

current Epsilon, 1 + current Epsilon

1	2.00000000000000000000
0.5	1.50000000000000000000
0.25	1.25000000000000000000
0.125	1.12500000000000000000
0.0625	1.06250000000000000000
0.03125	1.03125000000000000000
0.015625	1.01562500000000000000
0.0078125	1.00781250000000000000
0.00390625	1.00390625000000000000
0.00195312	1.00195312500000000000
0.000976562	1.00097656250000000000
0.000488281	1.00048828125000000000
0.000244141	1.00024414062500000000
0.00012207	1.00012207031250000000
6.10352E-05	1.00006103515625000000
3.05176E-05	1.00003051757812500000
1.52588E-05	1.00001525878906250000
7.62939E-06	1.00000762939453125000
3.81470E-06	1.00000381469726562500
1.90735E-06	1.00000190734863281250
9.53674E-07	1.00000095367431640625
4.76837E-07	1.00000047683715820312
2.38419E-07	1.00000023841857910156

Calculated Machine epsilon: 1.19209E-07

Nota: in MATLAB

“eps” o “eps(‘double’)” restituisce 2^{-52}

“eps(‘single’)” restituisce 2^{-23}

MACHEPS di alcune macchine

Da: Kalbasi, K., "Can you trust your computer? [digital arithmetic]," *Potentials, IEEE*, vol.9, no.2, pp.15-18, Apr 1990.

URL: <http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=52995&isnumber=1912>

Computer	Rounding or Chopping	base	cifre mantissa	E_{min}	E_{max}	MACHEPS
CDC CYBER 170	R	2	48	-976	1'071	3.55×10^{-15}
CDC CYBER 205	C	2	47	-28'626	28'718	1.42×10^{-14}
Cray-1	C	2	48	-8'192	8'191	7.11×10^{-15}
DEC VAX (single)	R	2	24	-127	127	5.96×10^{-8}
DEC VAX (double)	R	2	56	-1'023	1'023	1.11×10^{-16}
HP-11C, 15C	R	10	10	-99	99	5.00×10^{-10}
IBM 3033 (single)	C	16	6	-64	63	9.54×10^{-7}
IBM 3033 (double)	C	16	14	-64	63	2.22×10^{-16}
IBM/PC (single)	R	2	24	-126	127	5.96×10^{-8}
IBM/PC (double)	R	2	53	-1'022	1'023	1.11×10^{-16}
PRIME 850 (single)	C	2	23	-128	127	2.38×10^{-7}
PRIME 850 (double)	C	2	47	-32'896	32'639	1.42×10^{-14}

L'epsilon di macchina dipende da diversi fattori, incluso la modalita' di arrotondamento

FORMATO IEEE-754

Esercizio

- **Compito 2/11/07 domanda 1**