

CSC 373 Winter 2020 Prof. Lytinen

Representation of floating point numbers

Reading: Bryant and O'Hallaron, section 2.4

- **Motivation**

- What is the difference between integer and floating point numbers?
- How does their representation affect their range (smallest and largest possible values)?
- When does one use `int`, and when `double` or `float`? What are the advantages/disadvantages of each?

Example of floating point behavior

```
#include <stdio.h>
```

```
int main() {  
    double x = 1.2;  
    printf("x = : %0.50f\n\n", x);  
}
```

```
$ gcc -o test test.c ./test
```

```
x = 1.19999999999999999995559107901499373838305473327636719
```

Binary floating point

For our purposes, a binary floating point number is of the form

$(0|1)^+.(0|1)^* \mid (0|1)^*.(0|1)^+$

where

| means or

* means zero or more

+ means one or more

Examples of binary floating point numbers:

1.

.1

10.1

11111.0000101

Not binary floating point:

10

Meaning of a binary floating point number

- Each 1 represents a power of 2
- Each digit to the right of the '.' corresponds to a negative power of 2
 - The first digit to the right is 2^{-1} , etc.
- Each digit to the left of '.' is a non-negative power of 2, starting with 2^0
- Simple examples

Binary	Base 10
1000.0	8
0.1	0.5
0.01	0.25

- Example: 1001.101

$$2^3 + 2^0 + 2^{-1} + 2^{-3}$$

$$8 + 1 + .5 + .125 = 9.625$$

- **Exercise** What is 1011.011_2 in base 10?

Binary floating point	1	0	1	1	.	0	1	1
Each digit represents	2^3	2^2	2^1	2^0		2^{-1}	2^{-2}	2^{-3}

$$1011.011_2 = 11.375$$

Scientific notation

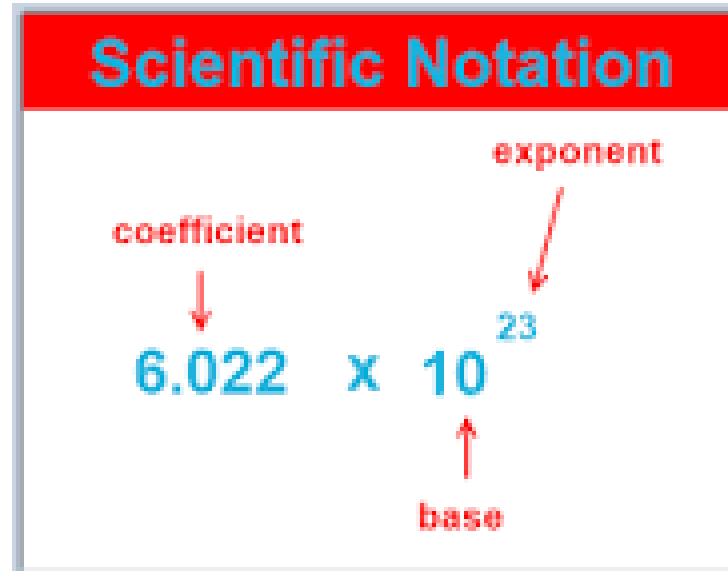
<https://www.boundless.com/physics/textbooks/boundless-physics-textbook/the-basics-of-physics-1/significant-figures-and-order-of-magnitude-33/scientific-notation-201-6207/>

For base 10:

$(\text{Coefficient}) \times 10^{\text{exponent}}$

$1 \leq \text{coefficient} < 10$

(or $-1 \geq \text{coefficient} > -10$
for negative numbers)



Note that 0 must have a special value: 0×10^0

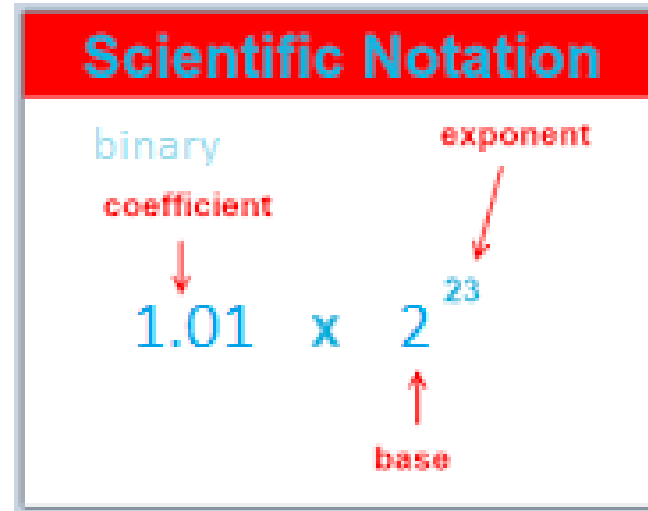
Binary scientific notation

$$(\text{Coefficient}) * 2^e$$

e must be an integer

$$1_2 \leq \text{coefficient} < 10_2$$

$$(1_{10} \leq \text{coefficient} < 2_{10})$$



Or for negative numbers:

$$-1_2 \geq \text{coefficient} > -10_2$$

For our purposes, we will still write the exponent in base 10.

Also, a special format is needed to represent 0

$$0 * 2^0$$

Converting a binary floating point number into binary scientific notation

- The binary floating point number must be non-zero
- Adjust the location of the point so that it appears in the right place (i.e., just to the right of the leftmost 1)
- Adjust the value of the exponent (e) accordingly. If the point has been moved to the left, then e is a positive number. If it's been moved to the right, e is a negative number. If the point has not been moved, then $e=0$.
- Examples:

$$111.01 = 1.1101 * 2^2$$

$$0.00101 = 1.01 * 2^{-3}$$

$$1.001 = 1.001 * 2^0$$

Exercises

Binary Floating Point	Binary Scientific Notation
1101.011	$1.101011 * 2^3$
.01011	$1.011 * 2^{-2}$
-1010100	$-1.0101 * 2^6$
101.01	$1.0101 * 2^2$
1	$1 * 2^0$

Note: 2s complement is not used in floating point numbers

Exercises

Decimal Fractional	Binary Floating Point	Binary Scientific Notation
1/4	.01	$1.0 * 2^{-2}$
1 7/16	1.0111	$1.0111 * 2^0$
3/8	0.011	$1.1 * 2^{-2}$
5/8	.101	$1.01 * 2^{-1}$
-4 5/16	-100.0101	$-1.00101 * 2^2$
3 15/16	11.1111	$1.11111 * 2^1$
-7	-1111	$-1.111 * 2^3$

Note: 2s complement is not used in floating point numbers

IEEE Floating Point Representation (Institute of Electrical and Electronics Engineers)

We modify our definition of binary floating point scientific notation slightly as follows:

$$(-1)^s * C * 2^E$$

s indicates whether number is negative or non-negative (1 = negative)

C is a binary floating point number; where $1 \leq C < 10_2$ ($1 \leq C < 2_{10}$).

Note the coefficient cannot be negative, because the sign is indicated by **s**.

Alternate way to represent negative integers

In Binary Scientific Notation, the exponent must be an integer (could be +, -, or 0)

We could use 2s complement but there is an alternative, which uses a ***Bias***.

Counting: 2s complements vs biased

In general, if we use n bits to represent an integer, then the **bias** is $2^{n-1}-1$. *Example: 4 bits (bias = 7)*

Number	2s compl	w/bias
7	0111	1110
6	0110	1101
5	0101	1100
4	0100	1011
3	0011	1010
2	0010	1001
1	0001	1000

Number	2s compl	w/bias
0	0000	0111
-1	1111	1110
-2	1110	1101
-3	1101	1000
-4	1100	0111
-5	1101	0110
-6	1110	0001

In biased form, we'll reserve 0000 and 1111 to have special meanings

Counting: 8-bit 2s complements vs biased

In general, if we use n bits to represent an integer, then the ***bias*** is $2^{n-1}-1$. For 8-bit numbers, the bias is 2^7-1 (127)

Number	2s compl	w/bias
3	0x03	0x82
2	0x02	0x81
1	0x01	0x80
0	0x00	0x7f
-1	0xff	0x7e
-2	0xfe	0x7d
-3	0xfd	0x7c

In biased form, we'll reserve 0x00 and 0xff to have special meanings

IEEE Encoding

- Leftmost bit indicates sign
- Remaining bits are divided between `exp` and `frac` according to size of datatype
- $\text{exp} = E + \text{the bias}$, which depends on the data size (in bits)
- $\text{frac} = C - 1$
- We can represent a range of large (positive) and small (negative) exponents; for 32-bit, the exponent is between 11111110_2 and 00000001_2 (11111111 and 00000000 have special meanings)

Data type	# bits in s	# bits in exp	bias	# bits in frac
double	1	11	$2^{10}-1$	52
float	1	8	$2^7 - 1$ (127)	23

Simple example

Represent $\frac{1}{2}$ in IEEE 32-bit floating point

$$\frac{1}{2} = 0.1_2 = 1.0 * 2^{-1}$$

$$B = 1.0$$

$$E = -1$$

$$\text{sign} = 0$$

$$\text{frac} = 1.0 - 1 = .0$$

$$\text{exp} = E + \text{bias} = -1 + 127 = 126$$

s	exp	frac
0	01111110	000000000000000000000000
1 bit	8 bits	23 bits

Simple example

Represent -2.25_{10} in IEEE 32-bit floating point

$$-2.25_{10} = -10.01_2 = -1.001 * 2^1$$

$$C = -1.001$$

$$E = 1$$

$$\text{sign} = 1$$

$$\text{frac} = 1.001 - 1 = .001$$

$$\text{exp} = E + \text{bias} = 1 + 127 = 128$$

s	exp	frac
1	10000000	001000000000000000000000
1 bit	8 bits	23 bits

Simple example

Represent $\frac{1}{2}$ in IEEE 64-bit floating point

$$\frac{1}{2} = 0.1_2 = 1.0 * 2^{-1}$$

$$C = 1.0$$

$$E = 0$$

$$\text{sign} = 0$$

$$\text{frac} = 1.0 - 1 = .0$$

$$\text{exp} = E - \text{bias} = -1 + 1023 = 1022 = 01111111110$$

s	exp	frac
0	01111111110	0000000000000000000000...000
1 bit	11 bits	52 bits

Simple example

Represent -2.5_{10} in IEEE 32-bit and 64-bit floating point

$$-2.5_{10} = -10.1_2 = -1.01 * 2^1$$

$$C = -1.001$$

$$E = 1$$

$$\text{sign} = 1$$

$$\text{frac} = 1.01 - 1 = .01$$

$$\text{exp (32-bit)} = E + \text{bias} = 1 + 127 = 128$$

$$\text{exp} = E - \text{bias} = 1 + 1023 = 1024$$

32-bit	s	exp	Frac
	1	10000000	0100000000000000000000...000
	1 bit	8 bits	23 bits

64-bit	s	exp	frac
	1	100000000000	0100...000
	1 bit	11 bits	54 bits

Exercises

Fill in the blanks.

Decimal	Binary Scientific Notation	IEEE 32-bit	IEEE 64-bit
4.5	$1.001 * 2^2$	0x40900000	0x4012000000000000
-2.375	$-1.0011 * 2^1$	0xc0180000	
	$-1.11 * 2^0$		
		0xbf000000	
			0x3ff4000000000000
-1.625	$-1.101 * 2^0$	1 01111111 10100000....	0xbfd00000

Exercise solutions

Fill in the blanks.

Decimal	Binary Scientific Notation	IEEE 32-bit	IEEE 64-bit
4.5	$1.001 * 2^2$	0x40900000	0x4012000000000000
-2.25			
	$-1.11 * 2^0$		
		0xbf000000	
			0x3ff4000000000000

Code: print the bits of a 32-bit IEEE floating point (in hex)

```
void print_float(float *fptr) {
    printf("0x");
    void *vptr = (void *) fptr; // pointer to void is typeless
    int i;
    for (i=7; i>=0; i-=1) {
        // shift 4 bits at a time, because 4 bits = 1 hex digit
        int x = (*((int *) vptr) >> 4*i) & 15;
        printf("%x", x);
    }
    printf("\n");
}
```

Range of doubles

Largest number 0 11111111110 1111...111111111111111111

exp = 1111111110 - bias = 0111111111

E = 0111111111 = 1023 = $2^{10}-1$

largest M is approximately 2

largest number = $2 * 2^{1023} = 2^{1024} =$

Largest number is $1.797693134862315907729305190789 * 10^{308}$

Smallest number 0 0000000001 0000 ... 00000000000000000001 =

Smallest M is 1.00000000000000000000000001

smallest E is 0000000001 – bias = -0111111110

Smallest number is approximately 2^{-1023}

= $1.1125369292536006915451163586662 * 10^{-308}$

Denormalized Values

Condition: $\text{exp} = 000\dots 0$

Value

Exponent value $E = -\text{Bias} + 1$

Significand value $M = 0.x_1x_2x_3x_4x_5x_6x_7x_8x_9x_{10}x_{11}x_{12}x_{13}x_{14}x_{15}x_{16}x_{17}x_{18}x_{19}x_{20}x_{21}x_{22}x_{23}x_{24}x_{25}x_{26}x_{27}x_{28}x_{29}x_{30}x_{31}x_{32}x_{33}x_{34}x_{35}x_{36}x_{37}x_{38}x_{39}x_{40}x_{41}x_{42}x_{43}x_{44}x_{45}x_{46}x_{47}x_{48}x_{49}x_{50}x_{51}x_{52}x_{53}x_{54}x_{55}x_{56}x_{57}x_{58}x_{59}x_{60}x_{61}x_{62}x_{63}x_{64}x_{65}x_{66}x_{67}x_{68}x_{69}x_{70}x_{71}x_{72}x_{73}x_{74}x_{75}x_{76}x_{77}x_{78}x_{79}x_{80}x_{81}x_{82}x_{83}x_{84}x_{85}x_{86}x_{87}x_{88}x_{89}x_{90}x_{91}x_{92}x_{93}x_{94}x_{95}x_{96}x_{97}x_{98}x_{99}x_{100}x_{101}x_{102}x_{103}x_{104}x_{105}x_{106}x_{107}x_{108}x_{109}x_{110}x_{111}x_{112}x_{113}x_{114}x_{115}x_{116}x_{117}x_{118}x_{119}x_{120}x_{121}x_{122}x_{123}x_{124}x_{125}x_{126}x_{127}x_{128}x_{129}x_{130}x_{131}x_{132}x_{133}x_{134}x_{135}x_{136}x_{137}x_{138}x_{139}x_{140}x_{141}x_{142}x_{143}x_{144}x_{145}x_{146}x_{147}x_{148}x_{149}x_{150}x_{151}x_{152}x_{153}x_{154}x_{155}x_{156}x_{157}x_{158}x_{159}x_{160}x_{161}x_{162}x_{163}x_{164}x_{165}x_{166}x_{167}x_{168}x_{169}x_{170}x_{171}x_{172}x_{173}x_{174}x_{175}x_{176}x_{177}x_{178}x_{179}x_{180}x_{181}x_{182}x_{183}x_{184}x_{185}x_{186}x_{187}x_{188}x_{189}x_{190}x_{191}x_{192}x_{193}x_{194}x_{195}x_{196}x_{197}x_{198}x_{199}x_{200}x_{201}x_{202}x_{203}x_{204}x_{205}x_{206}x_{207}x_{208}x_{209}x_{210}x_{211}x_{212}x_{213}x_{214}x_{215}x_{216}x_{217}x_{218}x_{219}x_{220}x_{221}x_{222}x_{223}x_{224}x_{225}x_{226}x_{227}x_{228}x_{229}x_{230}x_{231}x_{232}x_{233}x_{234}x_{235}x_{236}x_{237}x_{238}x_{239}x_{240}x_{241}x_{242}x_{243}x_{244}x_{245}x_{246}x_{247}x_{248}x_{249}x_{250}x_{251}x_{252}x_{253}x_{254}x_{255}x_{256}x_{257}x_{258}x_{259}x_{260}x_{261}x_{262}x_{263}x_{264}x_{265}x_{266}x_{267}x_{268}x_{269}x_{270}x_{271}x_{272}x_{273}x_{274}x_{275}x_{276}x_{277}x_{278}x_{279}x_{280}x_{281}x_{282}x_{283}x_{284}x_{285}x_{286}x_{287}x_{288}x_{289}x_{290}x_{291}x_{292}x_{293}x_{294}x_{295}x_{296}x_{297}x_{298}x_{299}x_{300}x_{301}x_{302}x_{303}x_{304}x_{305}x_{306}x_{307}x_{308}x_{309}x_{310}x_{311}x_{312}x_{313}x_{314}x_{315}x_{316}x_{317}x_{318}x_{319}x_{320}x_{321}x_{322}x_{323}x_{324}x_{325}x_{326}x_{327}x_{328}x_{329}x_{330}x_{331}x_{332}x_{333}x_{334}x_{335}x_{336}x_{337}x_{338}x_{339}x_{340}x_{341}x_{342}x_{343}x_{344}x_{345}x_{346}x_{347}x_{348}x_{349}x_{350}x_{351}x_{352}x_{353}x_{354}x_{355}x_{356}x_{357}x_{358}x_{359}x_{360}x_{361}x_{362}x_{363}x_{364}x_{365}x_{366}x_{367}x_{368}x_{369}x_{370}x_{371}x_{372}x_{373}x_{374}x_{375}x_{376}x_{377}x_{378}x_{379}x_{380}x_{381}x_{382}x_{383}x_{384}x_{385}x_{386}x_{387}x_{388}x_{389}x_{390}x_{391}x_{392}x_{393}x_{394}x_{395}x_{396}x_{397}x_{398}x_{399}x_{400}x_{401}x_{402}x_{403}x_{404}x_{405}x_{406}x_{407}x_{408}x_{409}x_{410}x_{411}x_{412}x_{413}x_{414}x_{415}x_{416}x_{417}x_{418}x_{419}x_{420}x_{421}x_{422}x_{423}x_{424}x_{425}x_{426}x_{427}x_{428}x_{429}x_{430}x_{431}x_{432}x_{433}x_{434}x_{435}x_{436}x_{437}x_{438}x_{439}x_{440}x_{441}x_{442}x_{443}x_{444}x_{445}x_{446}x_{447}x_{448}x_{449}x_{450}x_{451}x_{452}x_{453}x_{454}x_{455}x_{456}x_{457}x_{458}x_{459}x_{460}x_{461}x_{462}x_{463}x_{464}x_{465}x_{466}x_{467}x_{468}x_{469}x_{470}x_{471}x_{472}x_{473}x_{474}x_{475}x_{476}x_{477}x_{478}x_{479}x_{480}x_{481}x_{482}x_{483}x_{484}x_{485}x_{486}x_{487}x_{488}x_{489}x_{490}x_{491}x_{492}x_{493}x_{494}x_{495}x_{496}x_{497}x_{498}x_{499}x_{500}x_{501}x_{502}x_{503}x_{504}x_{505}x_{506}x_{507}x_{508}x_{509}x_{510}x_{511}x_{512}x_{513}x_{514}x_{515}x_{516}x_{517}x_{518}x_{519}x_{520}x_{521}x_{522}x_{523}x_{524}x_{525}x_{526}x_{527}x_{528}x_{529}x_{530}x_{531}x_{532}x_{533}x_{534}x_{535}x_{536}x_{537}x_{538}x_{539}x_{540}x_{541}x_{542}x_{543}x_{544}x_{545}x_{546}x_{547}x_{548}x_{549}x_{550}x_{551}x_{552}x_{553}x_{554}x_{555}x_{556}x_{557}x_{558}x_{559}x_{560}x_{561}x_{562}x_{563}x_{564}x_{565}x_{566}x_{567}x_{568}x_{569}x_{570}x_{571}x_{572}x_{573}x_{574}x_{575}x_{576}x_{577}x_{578}x_{579}x_{580}x_{581}x_{582}x_{583}x_{584}x_{585}x_{586}x_{587}x_{588}x_{589}x_{590}x_{591}x_{592}x_{593}x_{594}x_{595}x_{596}x_{597}x_{598}x_{599}x_{600}x_{601}x_{602}x_{603}x_{604}x_{605}x_{606}x_{607}x_{608}x_{609}x_{610}x_{611}x_{612}x_{613}x_{614}x_{615}x_{616}x_{617}x_{618}x_{619}x_{620}x_{621}x_{622}x_{623}x_{624}x_{625}x_{626}x_{627}x_{628}x_{629}x_{630}x_{631}x_{632}x_{633}x_{634}x_{635}x_{636}x_{637}x_{638}x_{639}x_{640}x_{641}x_{642}x_{643}x_{644}x_{645}x_{646}x_{647}x_{648}x_{649}x_{650}x_{651}x_{652}x_{653}x_{654}x_{655}x_{656}x_{657}x_{658}x_{659}x_{660}x_{661}x_{662}x_{663}x_{664}x_{665}x_{666}x_{667}x_{668}x_{669}x_{670}x_{671}x_{672}x_{673}x_{674}x_{675}x_{676}x_{677}x_{678}x_{679}x_{680}x_{681}x_{682}x_{683}x_{684}x_{685}x_{686}x_{687}x_{688}x_{689}x_{690}x_{691}x_{692}x_{693}x_{694}x_{695}x_{696}x_{697}x_{698}x_{699}x_{700}x$

xxx...x: bits of frac

Cases

exp = 000...0, frac = 000...0

Represents value 0

Note distinct values $+0$ and -0

exp = 000...0, frac != 000...0

Numbers very close to 0.0

Example: 0 00000000 000000000000000000000001 =

$$2^{-126} * 2^{-52} = 2^{-178} = 2.6101217871994098106841176437026 * 10^{-54}$$

Rounding

In base 10, not all numbers can be expressed in a finite decimal notation

Example: $1/3 = 0.333333333333333333333333$

The same is true for binary

$1/5 = ???$

Only numbers that are sums of integral powers of 2 (positive or negative) are exactly representable

Other numbers must be rounded

$$1/5 = 1/8 + 1/16 + 1/64 + 1/128 + \dots = 0.001100110011\dots_2$$

IEEE 32-bit representation:

$$0.001100110011\dots = 1.100110011\dots * 2^{-3}$$

32-bit IEEE: Exp = $127 - 3 = 124 = 01111100$ Frac = $100110011\dots$ 0x3e4ccccd

Rounding Example

```
int main() {  
    float frac = .6;  
    printf("%.20f\n", frac);  
    int i;  
    for (i=1; i<10000; i++) {  
        frac += (3.+13./((float)i));  
    }  
    for (i=1; i<10000; i++) {  
        frac -= (3.+13./((float)i));  
    }  
    printf("%.20f\n", frac);  
}
```

```
[slytinen@cdmlinux ieee]$ ./round  
0.600000002384185791016  
1.44581007957458496094
```