

exponent: E

Example 1 (1):

b = 2
p = 3
L = -1
U = 1

start enumerating possibilities:

+-	m	E	
+-	0.00	-1	-> 0
+-	0.00	0	-> 0
+-	0.00	+1	-> 0
+-	0.01	-1	-> 0.001
+-	0.01	0	-> 0.01
+-	0.01	+1	-> 0.1
+-	0.10	-1	-> 0.01
+-	0.10	0	-> 0.1
+-	0.10	+1	-> 1.0

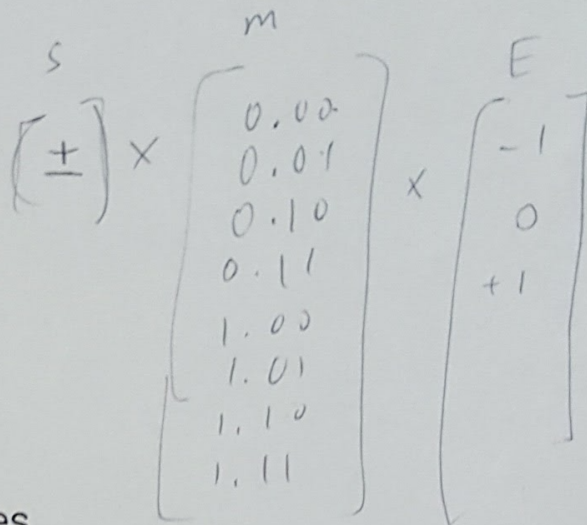
duplicates!

In general, number of possibilities

$$2 * b^p * (U - L + 1)$$

but

- lots of duplicates
- non-unique representation



$$2 * 8 * 3 = 48$$

Normalization

- require the leading digit to be non-zero
- so mantissa, m
- $1 \leq m < b$
- nice because:
 - representation is now unique*
 - don't waste digits on any leading 0's
 - for binary base, leading digit must be 1
 - so don't need to store it, just assume number is 1.d1d2..dp
 - gain an extra bit of precision!

$$b=10 \quad (10.00)_{10} = (10)_{10}$$

$$b=2 \quad (10.00)_2 = (2)_{10}$$

$$b=3 \quad (10.00)_3 = (3)_{10}$$

Properties

- finite and discrete system

- finite: how many (normalized) numbers can be represented?

count them:

$$2 * (b - 1) * b^{-(p-1)} * (U - L + 1) + 1$$

$$\text{Ex. 1 } 2 * 1 * 2^2 * 3 + 1 = 25$$

- what's the smallest (positive) normalized number? or "underflow level (UFL)"

$$\text{Ex. 1 } \text{UFL} = .5$$

$$1.0 \dots 0 * b^{-L} = b^{-L}$$

OFL

- what's the biggest normalized number? or "overflow level (OFL)"

$$\begin{aligned} & (b-1) \cdot (b-1) \dots (b-1) * b^U \\ & = (b - b^{-(p-1)}) * b^U \\ & = (1 - b^{-p}) * b^{(U+1)} \end{aligned}$$

$$\text{Ex. 1 } \text{OFL} = (1 - 2^{-3}) * 2^2 = (1 - \frac{1}{8}) * 4 = 4 - \frac{1}{2} = 3.5$$

Example 1 (2):

$$\begin{aligned} b &= 2 \\ p &= 3 \\ L &= -1 \\ U &= 1 \end{aligned}$$

- number of normalized

$$\begin{aligned} & 2 * (b - 1) * b^{-(p-1)} * (U - L + 1) + 1 \\ & = 2 * (2 - 1) * 2^{-(3-1)} * (1 - -1 + 1) + 1 \\ & = 2 * 1 * 4 * 3 + 1 \\ & = 25 \quad \checkmark \end{aligned}$$

- UFL

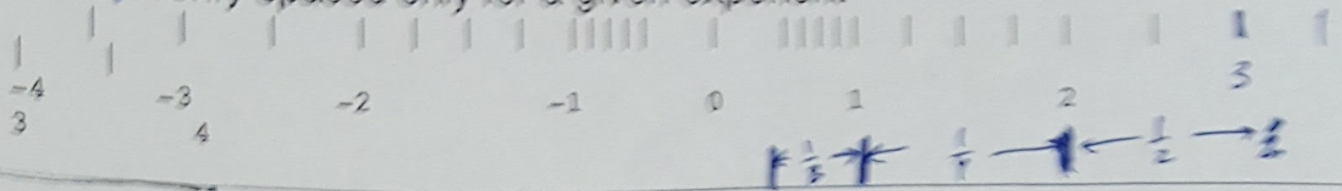
$$\begin{aligned} & b^{-L} \\ & = 2^{-(-1)} \\ & = .5 \quad \checkmark \end{aligned}$$

- OFL

$$\begin{aligned} & (1 - b^{-p}) * b^{(U+1)} \\ & = (1 - 2^{-3}) * 2^2 \\ & = 3.5 \quad \checkmark \end{aligned}$$

PICTURE of representable numbers

- note, evenly spaced only for a given exponent



Rounding

- floating point system is discrete!
- not all real numbers representable
- those that are called "machine numbers"
- others must be "rounded"

$x \leftarrow f1(x)$

- leads to "rounding error" or "roundoff error"

- How to round?

1. Chop - truncate digits - "round to zero"
2. Round to nearest
 - in case of tie go to even

Example: Rounding

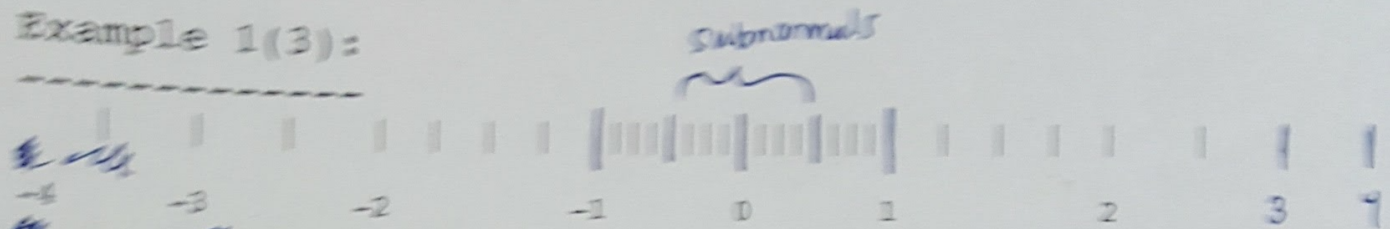
Number	Chop	Round to nearest
1.649	1.6	1.6
1.650	1.6	1.6 (tie - round to even)
1.651	1.6	1.7
1.699	1.6	1.7
1.749	1.7	1.7
1.750	1.7	1.8 (tie - round to even)
1.751	1.7	1.8
1.799	1.7	1.8

Machine Precision

Subnormals

- normalized numbers: gap between 0 and b^L
- fill in by allowing denormalized or subnormal numbers
- can make use of capacity for non-normalized numbers by allowing leading 0's
- though precision won't be full precision, since have leading 0's

Example 1(3):



- allows 6 new numbers around 0
- new smallest number is $(0.01)_2 * 2^{-1} = (0.125)_{10}$
- called "gradual underflow" because we gradually lose precision
- implementation: reserved value of exponent field
- leading bit not stored

Exceptional values

- Inf

- dividing finite number by 0
- exceeding OFL

- NaN

- undefined operation $0/0$, Inf/Inf , $0 * \text{Inf}$
- implemented through reserved values of exponent field

Floating Point Math

- adding or subtracting
- match exponents first
- must shift smaller number

Rounding

- floating point system is discrete!
- not all real numbers representable
- those that are called "machine numbers"
- others must be *rounded*

$x \leftarrow fl(x)$

- leads to "rounding error" or "roundoff error"

- How to round?

1. Chop - truncate digits - "round to zero"

2. Round to nearest

- in case of tie go to even

Example: Rounding

Number	Chop	Round to nearest
1.649	1.6	1.6
1.650	1.6	1.6 (tie - round to even)
1.651	1.6	1.7
1.699	1.6	1.7
1.749	1.7	1.7
1.750	1.7	1.8 (tie - round to even)
1.751	1.7	1.8
1.799	1.7	1.8

EXAMPLE : !!!!! warning: don't compare fp numbers with == !!!!!
octave-online.net

```
>> 4/3-1 == 1/3
```

```
ans = 0
```

```
>> single((4/3-1))==single(1/3)
```

```
ans = 1
```

```
>> (4/3-1)-1/3
ans = -5.5511e-17
```

right way to compare:

```
>> abs((4/3 - 1) - 1/3) <= 1e-16
ans = 1
```

Machine Precision

eps_mach ϵ_{mach}

- characterizes accuracy
- "machine epsilon", "machine precision", "unit roundoff"
- depends on rounding rule

$\bar{0} . \bar{1} \bar{2} \bar{3} \dots \bar{(p-1)} \bar{p} \dots$

 $x \ x \ x \ \dots \ x$

- chop: (chop everything at and after b^p position)
 $b^{-(p-1)} = b^{(1-p)}$ lose up to this amount
- round: (lose up to half of chop)
 $1/2 \ b^{(1-p)}$

- ^{bound} tells us the max possible relative error in representation

$$\frac{|fl(x) - x|}{|x|} \leq \text{eps_mach}$$

(for normalized #'s)

- check:

$$\begin{aligned} &\leq \text{eps_mach} * b^e / |x| \\ &= \text{eps_mach} * b^e / (m * b^e) \\ &= \text{eps_mach} / m \\ &\leq \text{eps_mach} \end{aligned}$$

because mantissa is at least 1

$$m \geq 1$$

- alternative characterization, ϵ_{mach} smallest # s.t.

$$fl(1 + \epsilon_{mach}) > 1$$

Examples:

- (Ex. 1) ϵ_{mach} (chop, nearest) = .25, .125
- IEEE SP ϵ_{mach} (nearest) = $2^{-24} \approx 10^{-7}$ (about 7 decimal digits of precision)
- IEEE DP ϵ_{mach} (nearest) = $2^{-53} \approx 10^{-16}$ (about 16 decimal digits of precision)

Floating Point Math

+,-

- adding or subtracting +,-
- match exponents first
- must shift smaller number
- if the sum (or diff) contains more than p digits, then the ones smaller than p will be lost
- smallest number may be lost completely

x

- multiplication ok
- mult mantissas and sum exponents
- still need to round though, because product will generally have more digits (up to 2p)

- division (also need to round)

Example

$$\begin{array}{r}
 1.23 * 10^5 \\
 + 1.00 * 10^4 \quad (10^3, 10^2)
 \end{array}$$

at this point smaller # totally lost

$$\begin{array}{r}
 (1.23 * 10^5) \\
 + (1.00 * 10^4) \Rightarrow (1.23 * 10^5) \\
 \hline
 1.33 * 10^5 \\
 \\
 (1.23 * 10^5) \\
 + (1.00 * 10^2) \Rightarrow (1.23 * 10^5) \\
 \hline
 1.23 * 10^5
 \end{array}$$

- can also get overflow or underflow
- underflow often ok - 0 is good approximation
- overflow more serious problem - can't approximate the number in question

- IEEE standard gives us

$$x \text{ flop } y = \text{fl}(x \text{ op } y)$$

as long as overflow doesn't occur

$$\text{op} = +, -, \times, \div$$

- + and * commutative but *not* associative!

- Ex: for $\text{eps} < \text{eps_mach}$, and $2 \text{ eps} > \text{eps_mach}$

$$(1 + \text{eps}) + \text{eps} = 1$$

$$1 + (\text{eps} + \text{eps}) = 1 + 2 \text{ eps} > 1$$

Rounding Error Analysis

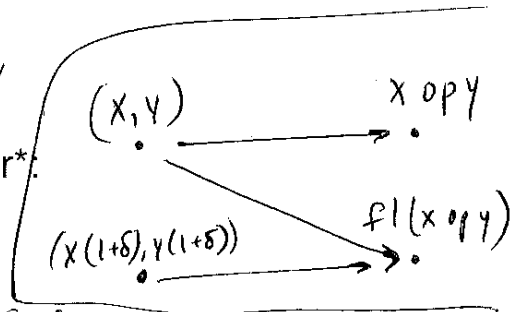
Basic idea is:

$$\text{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \delta),$$

$$|\delta| \leq \text{eps_mach}, \text{ and } \text{op} = +, -, *, /$$

rearranging, get bound on relative *forward error*:

$$\frac{|\text{fl}(x \text{ op } y) - (x \text{ op } y)|}{|(x \text{ op } y)|} = |\delta| \leq \text{eps_mach}$$



or, can interpret in terms of *backward error* (with $\text{op} = +$):

$$\text{fl}(x + y) = (x + y)(1 + \delta) = x(1+\delta) + y(1+\delta)$$

Example: Compute $x(y+z)$

$$\text{fl}(y+z) = (y+z)(1+d_1), \quad |d_1| \leq \text{eps_mach}$$

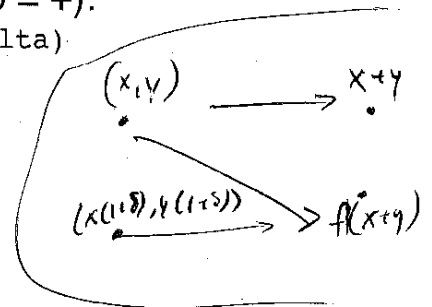
and

$$\text{fl}(x(y+z)) = (x(y+z)(1+d_1))(1+d_2), \quad |d_2| \leq \text{eps_mach}$$

$$= x(y+z)(1+d_1+d_2+d_1d_2)$$

$$\approx \approx x(y+z)(1+d_1+d_2)$$

$$= x(y+z)(1+d), \quad |d| = |d_1 + d_2| \leq 2 \text{ eps_mach}$$



- pessimistic bound

- typical, multiples of eps_mach accumulate

- but in practice this is generally ok

Cancellation

problems can arise when subtracting two very close numbers
 - result is exactly representable, but
 - e.g., if the numbers differ by rounding error, this can basically leave rounding error only after subtracting

Examples

$$\begin{array}{r}
 x = 1.92403 * 10^2 \\
 - y = 1.92275 * 10^2 \\
 \hline
 0.00128 * 10^2 = .128 = 1.28 * 10^{-1}
 \end{array}$$

- only 3 significant digits in the result

BAD: computing *small quantity* as a difference of *large quantities*

$$e^x = 1 + x + x^2/2 + x^3/3! + \dots, \text{ for } x < 0$$

Example: Quadratic formula

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

$$b = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$\begin{array}{l}
 0.05010 x^2 - 98.78 x + 5.015 \\
 \text{roots } \approx 1971.605916, \quad \underline{\text{answer to 10 digits}} \\
 \quad \quad \quad 0.05077069387
 \end{array}$$

$$\begin{array}{l}
 b^2 - 4ac = 9757 - 1.005 = 9756 \quad \underline{\text{answer to 4 digits}} \\
 \text{sqrt}(") = 98.77 \\
 \text{roots: } (98.78 \pm 98.77) / 0.1002 = 1972, 0.09980
 \end{array}$$

subtraction of two *close* numbers (cancellation error), followed by division by *small* number (amplification)

Examples of Floating Point Issues

Finite precision

Calculating e in floating point. The limit converges very slowly. By the time we hit machine epsilon for $1/n$, we still don't have a very accurate estimate. On the other hand, the series formula converges up to full precision fairly quickly.

1. As a limit:

$$e = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n.$$

```
n=20;
for i=1:n
    e_lim_bad(i) = (1+1/10^i)^(10^i);
end
```

2. Series formula for e :

$$e^x = 1 + x + x^2/2 + \dots + x^n/n!$$
$$e^1 = 1 + 1 + 1/2 + \dots + 1/n!$$

```
n=20;
e_series(1) = 1;
for i=2:n
    e_series(i) = e_series(i-1) + 1 / factorial(i-1);
end
```

Overflow

Calculating

$$s = \frac{x}{\sqrt{1+x^2}}.$$

Note that $s = \left(\sqrt{1 + \frac{1}{x^2}}\right)^{-1}$, so $s \rightarrow 1$ as $x \rightarrow \infty$. Computing s directly in the first form given is problematic.

As x gets too large, the intermediate computation of x^2 will overflow, even though the final value we are interested in computing is close to 1. See below. Instead, the second form is better suited to computation.

(single precision)

```
res = 0;
for i=1:4
    res(i) = single(10^i)/sqrt(1+single(10^i)^2)
end
```

(double precision)

```
for i=0:10:1000
    res(i) = 10^i/sqrt(1+(10^i)^2)
end
```

Cancellation

Example 1: Variance

. Subtracting two large numbers to get at a small difference is a bad idea. There won't be much if any precision to represent the small number we are interested in computing. As an example, let's consider two ways of calculating the variance of a set of numbers.

Method 1: Calculate the mean first. Then calculate variances as sum of squares of distance to mean.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad [\text{mean}]$$
$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad [\text{variance}]$$

Method 2: Use this (mathematically equivalent) formula:

$$\sigma^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \quad [\text{variance}]$$

Method 2 requires only one pass over the data, where as Method 1 requires two passes. Which method is better numerically? If \bar{x} is very large, but σ^2 is small, then Method 2 is very bad. The precision we have for σ^2 will depend on $\sigma^2/(n\bar{x}^2)$. In the first formula the dependence is on σ^2/\bar{x} .

Example 2: Series with alternating sign

. Consider the series

$$e^{-x} = 1 - x + \frac{x^2}{2} - \frac{x^3}{3!} + \dots + (-1)^n \frac{x^n}{n!} + \dots$$

Computing this directly involves subtracting values of relatively large magnitude to compute a relatively small answer. This will result in loss of precision.

A better approach to computing this would be

$$e^{-x} = \frac{1}{e^x} = \frac{1}{1 + x + \frac{x^2}{2} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!} + \dots}$$