

Ex 1 Non associativity in the presence of round-off.

Solution: This is done in a class demo and the diary should be posted. Here are the commands.

```
n = 10000;  
a = randn(n,1);  b = randn(n,1);  c = randn(n,1);  
t = ((a+b)+c == a+(b+c));  
sum(t)
```

Right-hand side in 3rd line returns 1 for each instance when the two numbers are the same.

Ex 2 Find machine epsilon in matlab.

Solution:

```
u = 1;
```

```

for i=0:999
    fprintf(1,' i = %d , u = %e \n',i,u)
    if (1.0 +u == 1.0) break, end
    u = u/2;
end
u = u*2

```



4 Proof of Lemma: If $|\delta_i| \leq \underline{u}$ and $n\underline{u} < 1$ then

$$\prod_{i=1}^n (1 + \delta_i) = 1 + \theta_n \quad \text{where} \quad |\theta_n| \leq \frac{n\underline{u}}{1 - n\underline{u}}$$

Solution:

The proof is by induction on n .

1) Basis of induction. When $n = 1$ then the product reduces to $1 + \delta_i$ and so we can take

$\theta_n = \delta_n$ and we know that $|\delta_n| \leq \underline{u}$ from the assumptions and so

$$|\theta_n| \leq \underline{u} \leq \frac{\underline{u}}{1 - \underline{u}},$$

as desired.

2) Induction step. Assume now that the result as stated is true for n and consider a product with $n + 1$ terms: $\prod_{i=1}^{n+1} (1 + \delta_i)$. We can write this as $(1 + \delta_{n+1}) \prod_{i=1}^n (1 + \delta_i)$ and from the induction hypothesis we get:

$$\prod_{i=1}^{n+1} (1 + \delta_i) = (1 + \theta_n)(1 + \delta_{n+1}) = 1 + \theta_n + \delta_{n+1} + \theta_n \delta_{n+1}$$

with θ_n satisfying the inequality $\theta_n \leq (n\underline{u})/(1 - n\underline{u})$. We call θ_{n+1} the quantity $\theta_{n+1} = \theta_n + \delta_{n+1} + \theta_n \delta_{n+1}$, and we have

$$\begin{aligned} |\theta_{n+1}| &= |\theta_n + \delta_{n+1} + \theta_n \delta_{n+1}| \\ &\leq \frac{n\underline{u}}{1 - n\underline{u}} + \underline{u} + \frac{n\underline{u}}{1 - n\underline{u}} \times \underline{u} = \frac{n\underline{u} + \underline{u}(1 - n\underline{u}) + n\underline{u}^2}{1 - n\underline{u}} = \frac{(n + 1)\underline{u}}{1 - n\underline{u}} \\ &\leq \frac{(n + 1)\underline{u}}{1 - (n + 1)\underline{u}} \end{aligned}$$

This establishes the result with n replaced by $n + 1$ as wanted and completes the proof. \square

Supplemental notes: Floating Point Arithmetic

In most computing systems, real numbers are represented in two parts: A mantissa and an exponent. If the representation is in the base β then:

$$x = \pm (.d_1 d_2 \cdots d_m)_{\beta} \beta^e$$

- $.d_1 d_2 \cdots d_m$ is a fraction in the base- β representation
- e is an integer - can be negative, positive or zero.
- Generally the form is normalized in that $d_1 \neq 0$.

Example: In base 10 (for illustration)

1. 1000.12345 can be written as

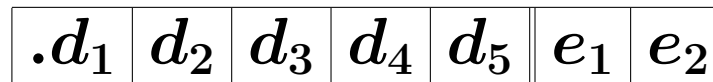
$$0.100012345_{10} \times 10^4$$

2. 0.000812345 can be written as

$$0.812345_{10} \times 10^{-3}$$

➤ Problem with floating point arithmetic: we have to live with limited precision.

Example: Assume that we have only 5 digits of accuracy in the mantissa and 2 digits for the exponent (excluding sign).



Try to add $1000.2 = .10002e+03$ and $1.07 = .10700e+01$:

$$1000.2 = \boxed{.1} \boxed{0} \boxed{0} \boxed{0} \boxed{2} \boxed{0} \boxed{4} ; \quad 1.07 = \boxed{.1} \boxed{0} \boxed{7} \boxed{0} \boxed{0} \boxed{0} \boxed{1}$$

First task: align decimal points. The one with smallest exponent will be (internally) rewritten so its exponent matches the largest one:

$$1.07 = 0.000107 \times 10^4$$

Second task: add mantissas:

$$\begin{array}{r} 0.10002 \\ + 0.000107 \\ \hline = 0.100127 \end{array}$$

Third task:

round result. Result has 6 digits - can use only 5 so we can

➤ Chop result:

.1	0	0	1	2
----	---	---	---	---

 ;

➤ Round result:

.1	0	0	1	3
----	---	---	---	---

 ;


Fourth task:

Normalize result if needed (not needed here)

result with rounding:

.1	0	0	1	3	0	4
----	---	---	---	---	---	---

 ;

 Redo the same thing with $7000.2 + 4000.3$ or $6999.2 + 4000.3$.

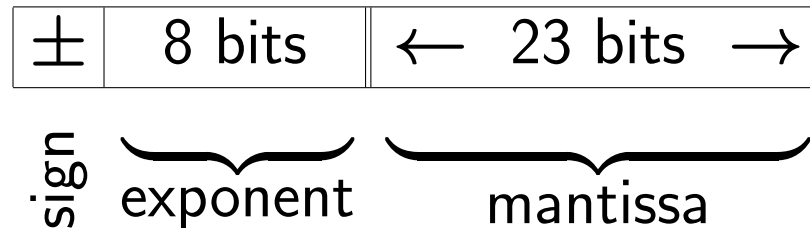
Some More Examples

- Each operation $fl(x \odot y)$ proceeds in 4 steps:
1. Line up exponents (for addition & subtraction).
 2. Compute temporary exact answer.
 3. Normalize temporary result.
 4. Round to nearest representable number (round-to-even in case of a tie).

	.40015 e+02	.40010 e+02	.41015 e-98
+	.60010 e+02	.50001 e-04	-.41010 e-98
temporary	1.00025 e+02	.4001050001e+02	.00005 e-98
normalize	.100025e+03	.400105 \oplus e+02	.00050 e-99
round	.10002 e+03	.40011 e+02	.00050 e-99
note:	round to even	round to nearest \oplus =not all 0's	too small: unnormalized
	exactly halfway between values	closer to upper value	exponent is at minimum

The IEEE standard

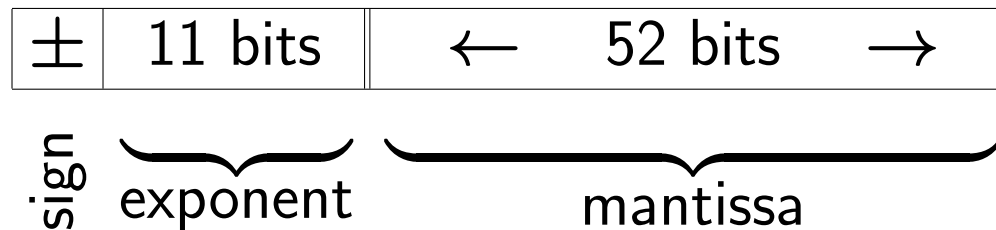
32 bit (Single precision) :




- Number is scaled so it is in the form $1.d_1d_2\dots d_{23} \times 2^e$ - but leading one is not represented.
- e is between -126 and 127.
- [Here is why: Internally, exponent e is represented in “biased” form: what is stored is actually $c = e + 127$ - so the value c of exponent field is between 1 and 254. The values $c = 0$ and $c = 255$ are for special cases (0 and ∞)]

64 bit

(Double precision) :



- Bias of 1023 so if e is the actual exponent the content of the exponent field is $c = e + 1023$
- Largest exponent: **1023**; Smallest = -1022.
- $c = 0$ and $c = 2047$ (all ones) are again for 0 and ∞
- Including the hidden bit, mantissa has total of 53 bits (52 bits represented, one hidden).
- In single precision, mantissa has total of 24 bits (23 bits represented, one hidden).

 6 Take the number 1.0 and see what will happen if you add $1/2, 1/4, \dots, 2^{-i}$. Do not forget the hidden bit!

Hidden bit (Not represented)

Expon. ↓ ← 52 bits →

e	1	1	0	0	0	0	0	0	0	0	0	0	0
e	1	0	1	0	0	0	0	0	0	0	0	0	0
e	1	0	0	1	0	0	0	0	0	0	0	0	0
.....													
e	1	0	0	0	0	0	0	0	0	0	0	0	1
e	1	0	0	0	0	0	0	0	0	0	0	0	0

(Note: The 'e' part has 12 bits and includes the sign)

➤ Conclusion

$$fl(1 + 2^{-52}) \neq 1 \text{ but: } fl(1 + 2^{-53}) == 1_{n+1}$$

Special Values

- Exponent field = 0000000000 (smallest possible value)
No hidden bit. All bits == 0 means exactly zero.
- Allow for unnormalized numbers,
leading to gradual underflow.
- Exponent field = 1111111111 (largest possible value)
Number represented is "Inf" "-Inf" or "NaN".