

ΕΚΠΑ Σχολή Θετικών Επιστημών
Τμήμα Μαθηματικών

Αριθμητική Γραμμική Άλγεβρα

Αριθμητική Κινητής Υποδιαστολής στη Julia

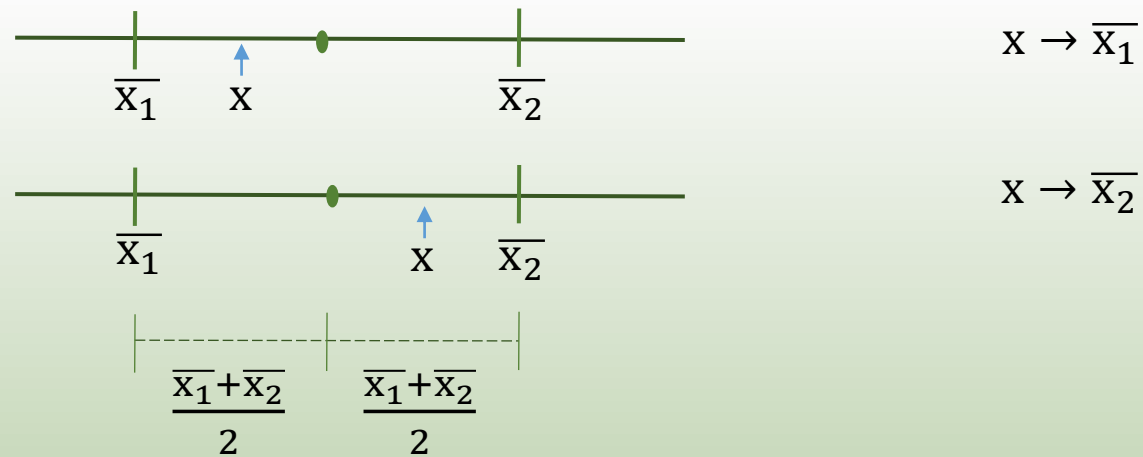
Οκτώβριος 2021

Ιορδάνης Ιωάννης

Καθηγήτρια: Μαριλένα Μητρούλη

Μετατροπή Πραγματικού Αριθμού x σε Αριθμό Κινητής Υποδιαστολής \bar{x}

Στρογγύλευση



Μορφή Αριθμού Κινητής Υποδιαστολής

$$\bar{x} = \sigma \cdot s \cdot 2^e$$

Diagram illustrating the components of the floating-point number format:

- σ is labeled "πρόσημο" (sign).
- s is labeled "mantissa" (mantissa).
- e is labeled "εκθέτης" (exponent).

Παράδειγμα : $\bar{x} = 1 \cdot 1.57075 \cdot 2^1 = 3.1415$

Μορφή Αριθμού Κινητής Υποδιαστολής

$$\bar{x} = \sigma \cdot s \cdot 2^e$$

πρόσημο : $\sigma = \begin{cases} 1, \text{θετικός} \\ -1, \text{αρνητικός} \end{cases}$

όρια εκθέτη : $m, M: m \leq e < M$ (καθορίζουν το εύρος των αριθμών \bar{x})

mantissa : $s = 1. (b_1 b_2 \dots b_d)_2 = 1 + b_1 \frac{1}{2^1} + b_2 \frac{1}{2^2} + \dots + b_d \frac{1}{2^d}, 1 \leq s < 2$

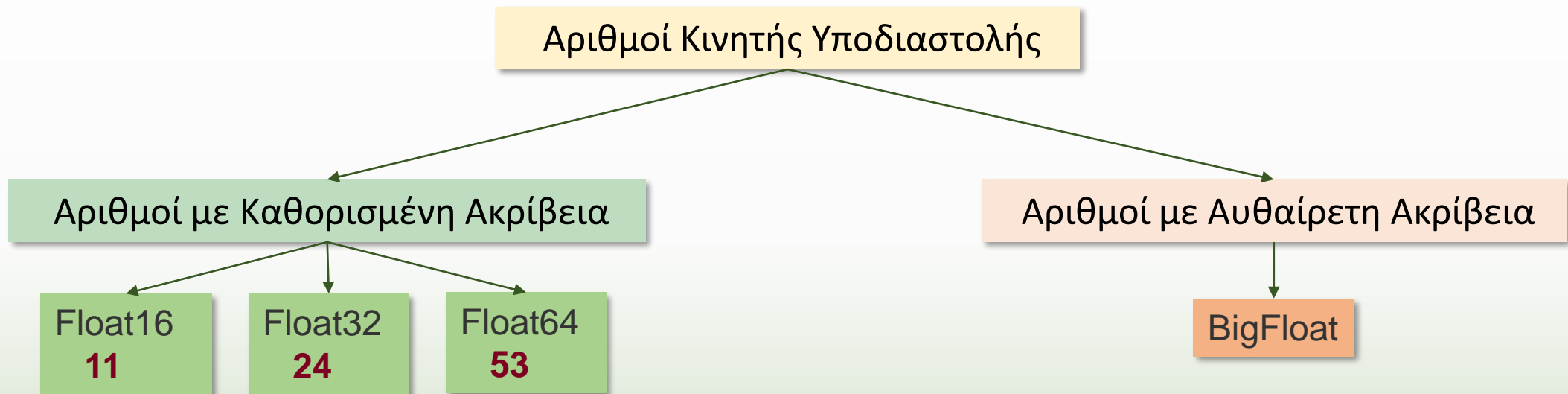
πλήθος ψηφίων σε mantissa: $t=1+d$ (καθορίζουν την ακρίβεια των αριθμών \bar{x})

Παράδειγμα

Είναι εφικτή η προσέγγιση του $x = 3.1415$ με τον $\bar{x} = 1 \cdot 1.57075 \cdot 2^1$;

Απώλεια στην ακρίβεια ($t=11$) $1.57075 \rightarrow 1.1001001000 \rightarrow 1.5703125$

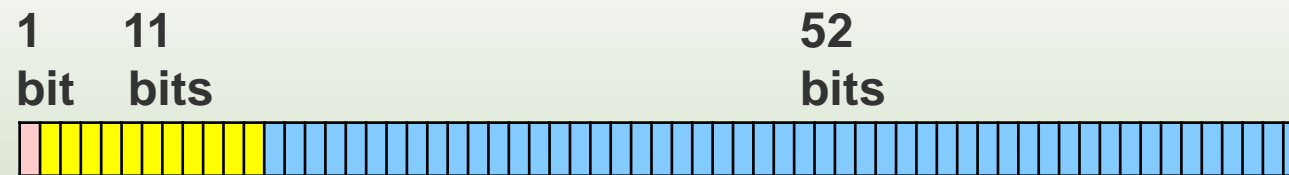
Τύποι (Κατηγορίες) Αριθμών Κινητής Υποδιαστολής



Τύποι Αριθμών Κινητής Υποδιαστολής με καθορισμένη ακρίβεια

Τύπος (Julia)	Τίτλος	Σύνολο (bits)	Πρόσημο (bit)	Εκθέτης (bits)	Ελάχιστος εκθέτης m (τιμή)	Μέγιστος εκθέτης M (τιμή)	Bias (τιμή)	$d=t-1$ (bits)	Ακρίβεια t (bits)
Float16	Μισή	16	1	5	-14	15	15	10	11
Float32	Απλή	32	1	8	-126	127	127	23	24
Float64	Διπλή	64	1	11	-1022	1023	1023	52	53

Δυαδική μορφή Float64 Αριθμού



Αλγόριθμος αποθήκευσης

Πρόσημο : Αν στο bit αποθηκεύεται η τιμή 0, ο \bar{x} είναι θετικός, αλλιώς ο \bar{x} είναι αρνητικός $(-1)^{\text{τιμή}}$

Εκθέτης : Αποθηκεύεται η τιμή $e+\text{bias}$, όπου bias σταθερά που εξαρτάται από τον τύπο του αριθμού

Mantissa : Αποθηκεύονται μόνο τα d bits. Δεν αποθηκεύεται το ακέραιο μέρος που είναι πάντα 1.

Μεγαλύτερο t , μεγαλύτερη ακρίβεια

Ακριβής αποθήκευση του $x=3.1415$ αν $\bar{x} = 1 \cdot 1.57075 \cdot 2^1$

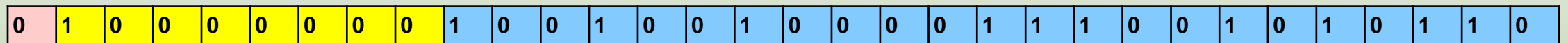
Με τι ακρίβεια αποθηκεύεται η mantissa;

Ακρίβεια t=11 (Float16)



0	16	0.5703125
- 15 (bias)	+ 1 (το ακέραιο μέρος της mantissa)	
-----	-----	
= 1	= 1.5703125	

Ακρίβεια t=24 (Float32)



0	128	0.570749998
- 127	+ 1	
-----	-----	
= 1	= 1.570749998	

Μετατροπή σε τύπο Αριθμού Κινητής Υποδιαστολής (Julia)

`x = 3.1415`

`x1 = Float16(x) = Float16(3.14)`

`(Float64(x1), typeof(x1), Float64(sign(x1)), exponent(x1), Float64(significand(x1))) =
(3.140625, Float16, 1.0, 1, 1.5703125)`

`x2 = Float32(x) = 3.1415f0`

`(Float64(x2), typeof(x2), Float64(sign(x2)), exponent(x2), Float64(significand(x2))) =
(3.14149999618, Float32, 1.0, 1, 1.5707499980926514)`

`x3 = Float64(x) = 3.1415`

`(Float64(x3), typeof(x3), sign(x3), exponent(x3), significand(x3)) =
(3.1415, Float64, 1.0, 1, 1.57075)`

ακρίβεια στο Float16: **precision(x1)** = 11

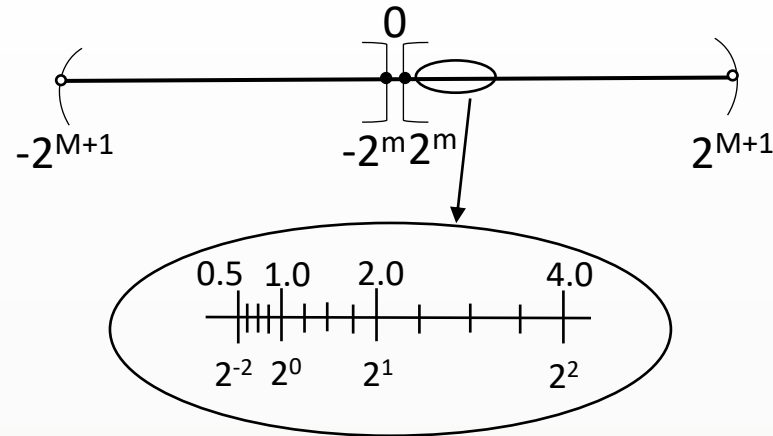
ακρίβεια στο Float32: **precision(x2)** = 24

ακρίβεια στο Float64: **precision(x3)** = 53

precision: επιστροφή των t ψηφίων της mantissa

Τυπολόγιο Αριθμών Κινητής Υποδιαστολής

Το σύνολο των αριθμών x χωρίζεται σε διαστήματα της μορφής $[2^e, 2^{e+1})$



1) Για κάθε διάστημα $[2^e, 2^{e+1})$:

- Εύρος διαστήματος $= 2^{e+1} - 2^e = 2^e$ **(1)** (κάθε διάστημα είναι διπλάσιο από το προηγούμενο)
- Πλήθος των \bar{x} στο διάστημα $= 2^d = 2^{t-1}$ **(2)** (το πλήθος σε κάθε διάστημα είναι σταθερό)
- Απόσταση μεταξύ 2 διαδοχικών αριθμών \bar{x} στο ίδιο διάστημα $= \frac{2^e}{2^d} = 2^{e-d}$ **(3)**

2) Έψιλον της μηχανής $\epsilon_{\text{mach}} = 2^{-d}$ **(4)**

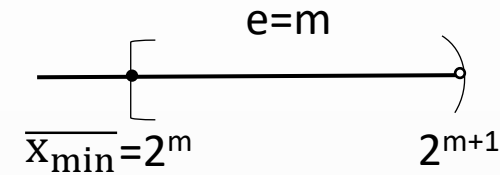
Απόσταση των στοιχείων στο $[2^0, 2^{0+1}) = [1, 2)$ ($e=0$)

Απόσταση του 1 με τον επόμενο διαδοχικό του

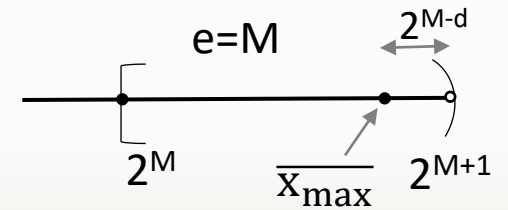
Τυπολόγιο Αριθμών Κινητής Υποδιαστολής

3) Όρια εκθέτη : $m, M: m \leq e \leq M$

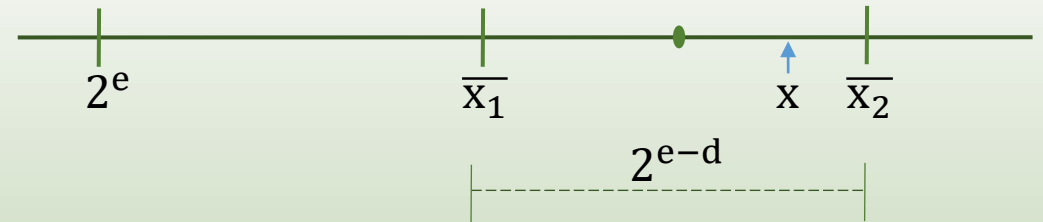
- Ελάχιστος θετικός αριθμός : $\bar{x}_{\min} = 2^m$ (5)



- Μέγιστος θετικός αριθμός : $\bar{x}_{\max} = 2^{M+1} - 2^{M-d} = 2^M - 2^{M-d} + 2^M$ (6)
(αφαίρεση απόστασης 2^{M-d} από πάνω όριο)



4) Στρογγύλευση του πραγματικού αριθμού x στον \bar{x} όπου ο \bar{x} είναι ο πλησιέστερος αριθμός ως προς x



Στο διάστημα $[2^e, 2^{e+1})$

- απόσταση μεταξύ x και \bar{x} : $|\bar{x} - x| \leq \frac{2^{e-d}}{2}$ (7)

- $\bar{x} \geq 2^e$ (8)

Σχετικό σφάλμα: $\frac{|\bar{x}-x|}{|x|} \leq \frac{2^{e-d}}{2 \cdot 2^e} = \frac{2^{-d}}{2} = \frac{\epsilon_{mach}}{2} = 2^{-t}$ (9)

μοναδιαίο σφάλμα στρογγύλευσης

Απόσταση μεταξύ 2 διαδοχικών αριθμών, έψιλον μηχανής (Julia) και διαστήματα

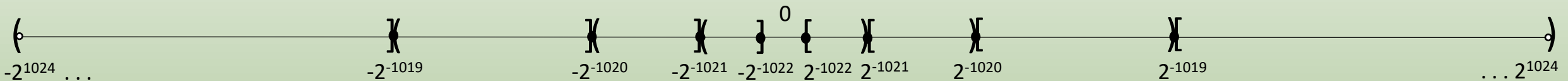
Απόσταση μεταξύ 2 διαδοχικών αριθμών κινητής υποδιαστολής στο $[2^e, 2^{e+1})$

απόσταση : $2.0^{(e-d)} = 4.440892098500626e-16$ (e=1, d=52) διάστημα $[2^1, 2^2)=[2,4)$

έψιλον μηχανής : $2.0^{(e-d)} = 2.220446049250313e-16$ (e=0, d=52) διάστημα $[2^0, 2^1)=[1,2)$

απόσταση στο διάστημα [2,4)	eps(Float64(3)) = 4.440892098500626e-16
έψιλον της μηχανής	eps(Float64) = 2.220446049250313e-16
έψιλον της μηχανής	eps() = 2.220446049250313e-16

Διαστήματα ευθείας αριθμών μηχανής (Float64)



Πλήθος αριθμών ανά διάστημα (Float64) : $\text{Int}(2.0^{(\text{precision}(\text{Float64})-1)})=4.503.599.627.370.496$

Αριθμοί Κινητής Υποδιαστολής με αυθαίρετη ακρίβεια (BigFloat)

Διερεύνηση προβλημάτων με: - δεδομένα που είναι εκτός των ορίων των Float64 αριθμών
- απαιτήσεις μέγιστης ακρίβειας

- Αριθμοί με αυθαίρετη ακρίβεια
- Προεπιλεγμένη ακρίβεια :256 bits
- Ελάχιστη ακρίβεια :2 bits
- Μέγιστη ακρίβεια : $2^{31} - 1 = 2,147,483,647$ bits
- Απαιτούμενα bits του εκθέτη :31 bits

ty=BigFloat

m = exponent(floatmin(ty)) = -1,073,741,824

M = exponent(floatmax(ty)) = 1,073,741,822

Int(round(log(2, -m + M + 1))) = 31

Εντολή για αυθαίρετη ακρίβεια: **setprecision(BigFloat, t)**

Δεν ισχύει η εντολή **bitstring**

Μετατροπή και ακρίβεια με BigFloat

pistring =

```
"3.141592653589793238462643383279502884197169399375105820974944592307816406286208998628034825342117067982148086513282306647"
```

aF32 = Float32(pi) = 3.1415927f0

aF64 = Float64(pi) = 3.141592653589793

aBF = **BigFloat(pi)** = 3.141592653589793238462643383279502884197169399375105820 9749445923078164062862

setprecision(BigFloat, 300) = 300

aBF300 = **aBF * BigFloat(1.0)** = 3.141592653589793238462643383279502884197169399 3751058209749445923078164062861 **980294536250318**

aBF300=**BigFloat(pi)** = 3.141592653589793238462643383279502884197169399 3751058209749445923078164062862089986280348248

Μετατροπή και ακρίβεια με BigFloat

```
typeof(aF32)    = Float32
typeof(aF64)    = Float64
typeof(aBF)     = BigFloat
typeof(aBF300)  = BigFloat
precision(aBF) = 256
precision(aBF300) = 300
```

```
exponent(aF32)  = 1
exponent(aF64)  = 1
exponent(aBF)   = 1
exponent(aBF300) = 1
```

```
significand(aF32) = 1.5707964f0
```

```
significand(aF64) = 1.5707963267948966
```

```
significand(aBF) = 1.570796326794896619231321691639751442098584699687552910487472296153908203143099
```

```
significand(aBF300) = 1.5707963267948966192313216916397514420985846996875529104874722961539082031431044993140174124
```


Σωστή και Λανθασμένη χρήση με BigFloat

```
a1 = 1.2
typeof(a1) = Float64
a2 = BigFloat(a1) = 1.1999999999999999555910790149937383830547332763671875
typeof(a2) = BigFloat
```

```
b1 = 2.4
typeof(b1) = Float64
b2 = BigFloat(b1) = 2.399999999999999911182158029987476766109466552734375
typeof(b2) = BigFloat
```

```
c2 = a2 + b2      = 3.5999999999999998667732370449812151491641998291015625
```

```
a22 = BigFloat("3.2")      = 1.200000000000000000000000000000000000000000000000000000000000000000000000000000007
a22 = BigFloat(string(a1)) = 1.200000000000000000000000000000000000000000000000000000000000000000000000000000007
b22 = BigFloat(string(b1)) = 2.400000000000000000000000000000000000000000000000000000000000000000000000000000014
c22 = a22 + b22           = 3.600000000000000000000000000000000000000000000000000000000000000000000000000000021
```


Υπολογισμός του επόμενου του μέγιστου θετικού

$m = 2^{1024.0} = \text{Inf}$

$m1 = \text{floatmax}(\text{Float64}) = 1.7976931348623157e308$

Εκτός ορίων σε Float64

$m2 = \text{eps}(\text{Float64}(m1)) = 1.99584030953472e292$

$m1 + m2 = \text{Inf}$

$\text{setprecision}(\text{BigFloat}, 53) = 53$

$m3 = \text{BigFloat}(0) = 0.0$

$\text{nextmax1} = m1 + m2 + m3 = \text{Inf}$

$\text{nextmax1} = m1 + m3 + m2 = 1.7976931348623159e+308$

1η λύση με Bigfloat

$\text{nextmax1} = m3 + m1 + m2 = 1.7976931348623159e+308$

$\text{typeof}(\text{nextmax1}) = \text{BigFloat}$

$\text{nextmax2} = 2^{\text{BigFloat}(1024)} = 1.7976931348623159e+308$

2η λύση με BigFloat

Έλεγχος

$2^{\text{BigInt}(1024)} = 179769313486231590772930519078902473361797697894230657273430081157732675805500963132708477322407536021120113879871393357658789768814416622492847430639474124377767893424865485276302219601246094119453082952085005768838150682342462881473913110540827237163350510684586298239947245938479716304835356329624224137216$

$q1 = 2^{1023.0} = 8.98846567431158e307$

mantissa του $2^{1023.0}$: $q2 = q1 / 10^{307.0} = 8.98846567431158$

Λύση με Float64

$q3 = 2 \cdot q2 = 17.97693134862316$

mantissa του $2^{1024.0}$: $\text{nextmaxmant64} = q3 / 10 = 1.7976931348623162$