

Fonctionnement des Ordinateurs

TP4 - Nombres flottants

B. QUOTIN
Faculté des Sciences
Université de Mons

Résumé

L'objectif de cette séance d'exercices est de renforcer votre compréhension de la représentation des nombres flottants et du standard IEEE 754.

Table des matières

1	Motivation	1
1.1	Erreur suite à une addition	1
1.2	Swamping lors de l'addition	1
1.3	Somme de fractions	1
2	Virgule flottante	3
2.1	Représentation IEEE754	3
2.2	Limites de la représentation IEEE754 (format réduit)	4
2.3	Normalisation d'un nombre	5
2.4	Conversion vers IEEE754 (format réduit)	5
2.5	Biais adéquat	6
2.6	Erreur absolue, erreur relative et epsilon machine	6
2.7	Arrondis	7
2.8	Addition	8
2.9	Soustraction	11

Q3) Programme en C ou Java

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

Pouvez-vous écrire un programme donnant un comportement similaire en Python ?

Q7) 0 10000000 110000000000000000000000

.....

.....

.....

.....

Q8) 0 11111111 011111111111111111111111

.....

.....

.....

.....

Q9) 1 00000000 100000000000000000000000

.....

.....

.....

.....

2.2 Limites de la représentation IEEE754 (format réduit)

Donnez les plus grand et plus petit nombres positifs non nuls normalisés représentables avec les tailles de mantisse M , d'exposant E et les biais B suivants.

Q10) $M = 4, E = 3, B = 4$

.....

.....

.....

Q11) $M = 8, E = 4, B = 8$

.....

.....

.....

Question identique pour les nombres dénormalisés.

Q12) $M = 4, E = 3, B = 4$

.....

.....

.....

Q18) 64,17

.....

.....

.....

Q19) 0,03125

.....

.....

.....

Q20) -0,015625

.....

.....

.....

2.5 Biais adéquat

Soit une taille d'exposant E , donner la valeur du biais permettant d'équilibrer au mieux le nombre d'exposants positifs et négatifs.

Q21) $E = 3$

.....

.....

.....

Q22) $E = 8$

.....

.....

.....

2.6 Erreur absolue, erreur relative et epsilon machine

Soit une taille de mantisse $M = 4$, une taille d'exposant $E = 3$ et un biais $B = 4$. Fournissez la formule donnant l'epsilon machine, ϵ_M , la borne supérieure sur l'erreur relative. Donnez sa valeur pour $M = 4$.

Q23) epsilon machine (ϵ_M)

.....

Calculez l'erreur d'approximation causée lors de la représentation des nombres suivants. Fournissez l'erreur absolue (Δ_x) et l'erreur relative (ϵ_x). Comparez cette dernière à l'epsilon machine. Note : les nombres ci-dessous sont identiques à ceux de la Section 2.4.

Q24) 1,5

.....

.....

.....

Q25) 0,625

.....

.....

.....

Q26) 3,2

.....

.....

.....

Q27) 64,17

.....

.....

.....

Q28) 0,03125

.....

.....

.....

Q29) 0,015625

.....

.....

.....

2.7 Arrondis

En utilisant la méthode *round-to-nearest-even* (*arrondi au plus proche pair*), arrondissez les nombres dont la représentation binaire est donnée ci-dessous, de façon à ne garder que la partie entière.

Q30) 101.000000

.....

Q31) 1.010101

.....

Q47) $x = 1001\ 010000; y = 1000\ 100000$

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....