

Rapport TP 1

encadré par

Stéphane GRAILLAT

Arithmétique en virgule Flottante et Analyse des Erreurs

réalisé par

Clément APAVOU

Arthur ZUCKER

Mathématiques Appliquées et Informatique

3^e année de cycle ingénieur

23 octobre 2020

Paris, France

Exercice 1 : Compensated Horner-Scheme

Question 1 : Schéma de Horner classique

1.1.a Le schéma de Horner pour l'évaluation d'un polynôme est donnée dans le listing suivant :

```

1 function res = Horner(p,x)
2 p = fliplr(p);
3 n = length(p);
4 a = p;
5 s(n) = a(n);
6 for i = (n-1) : (-1) : 1
7     p(i) = s(i+1) * x;
8     s(i) = p(i) + a(i);
9 end
10 res = s(1);
11 end

```

La fonction `fliplr` permet d'inverser l'ordre des coefficients du vecteur afin d'obtenir les mêmes résultats que pour la fonction `polyval`. La représentation d'un polynôme dans matlab est probablement inverse.

1.1.b Effectuons une comparaison des performance de l'évaluation d'Horner sur le polynôme $p = 1 + 2x^2 + 3x^3 + 4x^4 + 5x^5 + 6x^6 + 7x^7 + 8x^8 + 9x^9$ évalué en $x = \pi$

```

1 format long
2 p = [1 2 3 4 5 6 7 8 9];
3 x = pi;
4 res2 = Horner (p,x)

```

```
res2 = 2.041366687877934e+04
```

La fonction `polyval` fournie par matlab donne un résultat similaire :

```
1 res1 = polyval(p,x)
```

```
res1 = 2.041366687877934e+04
```

Question 2 : Error-Free Transformations (EFT)

Les algorithmes de EFT permettent d'effectuer des opérations entre deux flottants en prenant en compte l'erreur. L'implémentation des algorithmes dit "EFT" est la suivante :

```
1 function [x,y] = FastTwoSum(a,b)
2 x = a+b;
3 y = (a-x)+b;
4 end
```

```
1 function [x,y] = TwoSum(a,b)
2 x = a+b;
3 z = x-a;
4 y = (a-(x-z))+(b-z);
5 end
```

```
1 function [x, y] = Split(a)
2 factor = 2^(27) +1;
3 c = factor * a;
4 x = c - (c - a);
5 y = a-x;
6 end
```

```
1 function [x,y] = TwoProduct(a,b)
2 x = (a*b);
3 [a1,a2]=Split(a);
4 [b1,b2]=Split(b);
5 y = (a2*b2 - (((x-(a1*b1))-(a2*b1))-(a1*b2)));
6 end
```

Question 3 : Schéma de Horner compensé

La version compensé du schéma de Horner permet d'effectuer une évaluation de Horner en prenant compte l'erreur grâce aux algorithmes de EFT. Une version compensée de l'évaluation d'Horner est implémentée dans le listing suivant :

```

1 function res = CompHorner(p,x)
2 p = fliplr(p);
3 n = length(p);
4 a = p;
5 s(n) = a(n);
6 r(n) = 0;
7 for i = (n-1) : (-1) : 1
8     [p(i),pi(i)] = TwoProduct(s(i+1),x);
9     [s(i),si(i)] = TwoSum(p(i),a(i));
10    r(i) = (r(i+1)*x+(pi(i)+si(i)));
11 end
12 res = s(1) + r(1);
13 end

```

En reprenant l'exemple précédent, on trouve l'évaluation suivante :

```

1 res3 = CompHorner(p,x)

```

```

res3 = 2.041366687877935e+04

```

Question 4 : Schéma de Horner en valeur exacte

En utilisant la fonction sym de la bibliothèque Symbolic Math Toolbox, la version exacte de l'évaluation de Horner est donnée par le code suivant :

```

1 function res = sHorner(p,x)
2 p = sym(p,'f')
3 x = sym(x,'f')
4 p = fliplr(p);
5 n = length(p);
6 a = p;
7 s(n) = a(n);
8 for i = (n-1) : (-1) : 1
9     p(i) = s(i+1) * x;
10    s(i) = p(i) + a(i);
11 end
12 res = s(1);
13 end

```

Question 5 : Conditionnement

La fonction `condp` calcul le conditionnement de p pour une évaluation en x . On utilise ici la fonction `sym` de la bibliothèque Symbolic Math Toolbox pour avoir une valeur exacte.

```

1 function condp = condp(p,x)
2 p=sym(p,'f');
3 x= sym(x,'f');
4 s1 = 0;
5 s2 = 0;
6 for i=1:length(p)
7     s1 = s1 + abs(p(i))*(abs(x)^i);
8     s2 = s2 + p(i)*(x^i);
9 end
10 s2 = abs(s2);
11 condp = s1/s2;
12 end

```

Question 6 : Comparaison de la précision des méthodes

Afin de tester les fonctions précédemment implémenté, nous les avons testé sur les polynômes $p_n(x) = (x - 1)^n$ pour $n = [3, \dots, 42]$ en prenant le flottant $x = 1.333$, et avons obtenu les résultats suivants :

```

1 n = 3:42;
2 x = 1.333;
3 for i=1:length(n)
4     a = poly(ones(1,n(i)));
5     res1 = Horner(a,x);           % Classic horner
6     res2 = sHorner(a,x);         % Thorical value
7     res3 = CompHorner(a,x);      % Compensated Horner
8     er1(i) = abs(res2-res1)/abs(res2);
9     er2(i) = abs(res2-res3)/abs(res2);
10    c(i) = condp(a,x);
11    %c(i) = abs((x+1)/(x-1))^i;
12    cond1(i) = 2*n(i)*eps*c(i);
13    cond2(i) = eps + 4*n(i)*n(i)*eps*eps*c(i);
14 end

```

Sur le graphe 1, les majorations théorique pour les deux schéma sont bien respectées.

Représentation des courbes d'erreurs relatives en fonction du conditionnement pour l'évaluation d'Horner

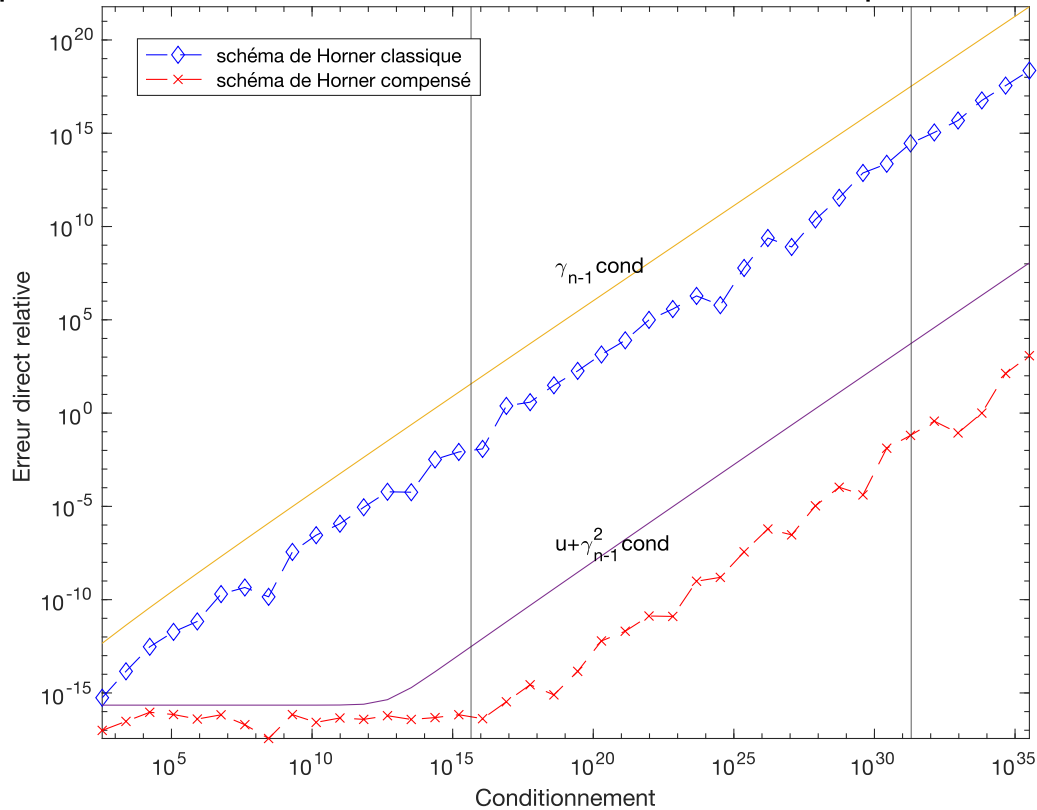


FIGURE 1 – Comparaison des erreurs pour l'évaluation d'Horner

Exercice 2 : Précisions des algorithmes de sommation

Question 1 : Les différents algorithmes de sommation

2.1.a L'algorithme de sommation récursif classique donne le code suivant :

```
1 function res = Sum(p)
2 n = length(p);
3 o = 0;
4 for i = 1:n
5     o = o + p(i);
6 end
7 res = o;
8 end
```

2.1.b L'algorithme de sommation de Kahan correspond au code qui suit :

```

1 function res = SCompSum(p)
2 n=length(p);
3 o = 0;
4 e = 0;
5 for i = 1:n
6     y = p(i) + e;
7     [o, e] = FastTwoSum(o, y);
8 end
9 res = o;
10 end

```

2.1.c L'algorithme de double compensation de la somme de Priest est implémenté avec le code suivant :

```

1 function res = DCompSum(p)
2 n=length(p);
3 [out,idx] = sort(abs(p), 'descend');
4 p = p(idx);
5 s = 0;
6 c = 0;
7 for i = 1:n
8     [y, u] = FastTwoSum(c, p(i));
9     [t, v] = FastTwoSum(y, s);
10    z = u + v;
11    [s, c] = FastTwoSum(t, z);
12 end
13 res = s;
14 end

```

2.1.d La méthode de somme compensée de Rump est présentée ci-après :

```

1 function res = CompSum(p)
2 n=length(p);
3 pi = zeros(1,n);
4 o = zeros(1,n);
5 pi(1) = p(1);
6 o(1) = 0;
7 q = zeros(1,n);
8 for i = 2:n
9     [pi(i), q(i)] = TwoSum(pi(i-1), p(i));
10    o(i) = o(i-1) + q(i);

```

```
11 end
12 res = pi(n)+o(n);
13 end
```

Question 2 : Comparaison de la précision

Étudions la précision des différentes méthodes présentées ci-dessus en fonction du conditionnement de la somme.

Représentation des courbes d'erreurs relatives en fonction du conditionnement pour la sommation

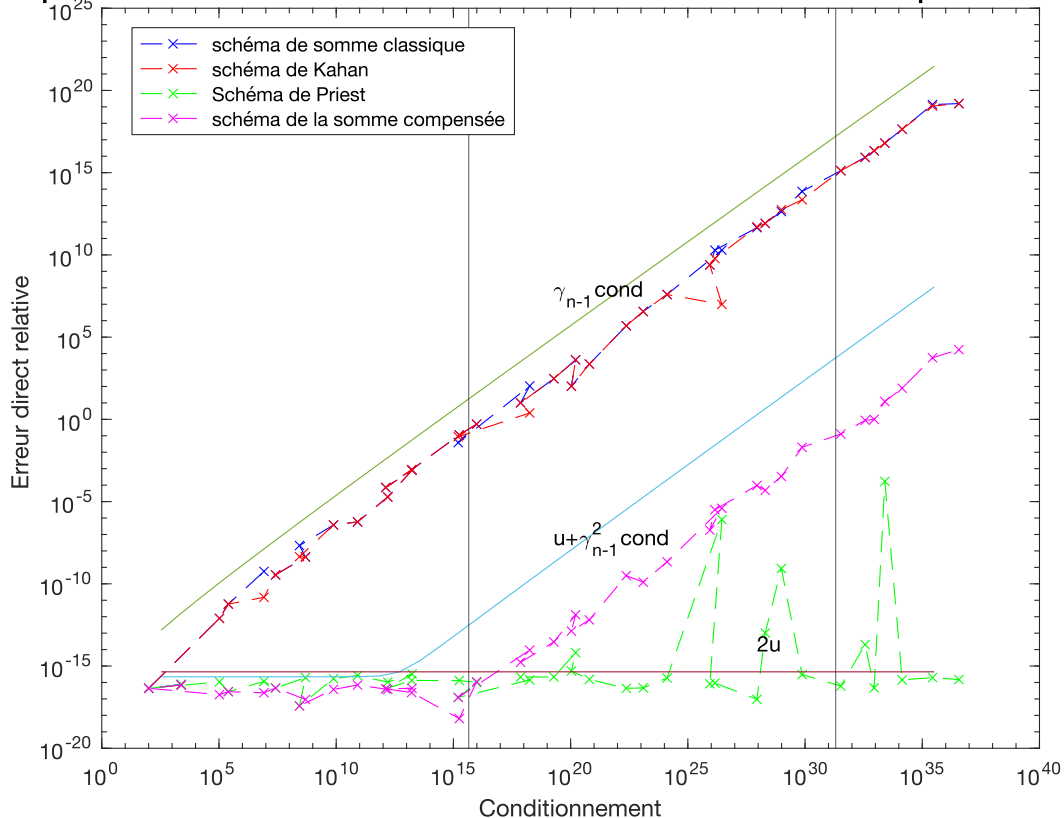


FIGURE 2 – Comparaison des erreurs sur les méthodes de sommations

Comme pour l'évaluation d'Horner, la méthode de sommation compensée est la plus intéressante (sans prendre en compte la méthode de Priest qui requière de trier le tableau d'entrée), on voit que son erreur commence à réellement perdre en précision lorsque le conditionnement dépasse $\frac{1}{u}$ avec u la double précision de l'ordinateur utilisé pour faire les calculs, ici $u = 2^{-52}$.

On observe juste quelques valeurs d'erreur un peu anormal pour DCompSum, qui est censée

rester en dessous de la droite d'équation $2u$ selon le théorème du cours. Ces valeurs sont probablement liées à un mauvais calcul au sein de la fonction.

L'algorithme de sommation de Kahan n'est pas très intéressant puisqu'il donne à peu près les mêmes résultats de performances que l'algorithme de sommation classique.