



HAL
open science

Contribution à l'étude des erreurs d'arrondi en arithmétique à virgule flottante

Michèle Pichat

► **To cite this version:**

Michèle Pichat. Contribution à l'étude des erreurs d'arrondi en arithmétique à virgule flottante. Modélisation et simulation. Institut National Polytechnique de Grenoble - INPG; Université Joseph-Fourier - Grenoble I, 1976. tel-00287209

HAL Id: tel-00287209

<https://theses.hal.science/tel-00287209>

Submitted on 11 Jun 2008

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

THESE

présentée à

UNIVERSITE SCIENTIFIQUE ET MEDICALE DE GRENOBLE
INSTITUT NATIONAL POLYTECHNIQUE DE GRENOBLE

POUR OBTENIR LE GRADE DE
DOCTEUR ES SCIENCES MATHEMATQUES

Michèle PICHAT

CONTRIBUTION
A L'ETUDE DES ERREURS D'ARRONDI
EN ARITHMETIQUE A VIRGULE FLOTTANTE

Thèse soutenue en 1976 devant la Commission d'Examen :

Président : J. KUNTZMANN

Examineurs { N. GASTINEL
A. HOCQUENGHEM
J. VIGNES
J.H. WILKINSON

UNIVERSITE SCIENTIFIQUE ET MEDICALE DE GRENOBLE

Monsieur Michel SOUTIF : Président

Monsieur Gabriel CAU : Vice-Président

MEMBRES DU CORPS ENSEIGNANTS DE L'U.S.M.G.

PROFESSEURS TITULAIRES

MM.	ANGLES D'AURIAC Paul	Mécanique des Fluides
	ARNAUD Paul	Chimie
	AUBERT Guy	Physique
	AYANT Yves	Physique Approfondie
Mme	BARBIER Marie-Jeanne	Electrochimie
MM.	BARBIER Jean-Claude	Physique Expérimentale
	BARBIER Reynold	Géologie Appliquée
	BARJON Robert	Physique Nucléaire
	BARNOUD Fernand	Biosynthèse de la Cellulose
	BARRA Jean-René	Statistiques
	BARRIE Joseph	Clinique Chirurgicale
	BEAUDOING André	Clinique de Pédiatrie et Puériculture
	BERNARD Alain	Mathématiques Pures
Mme	BERTRANDIAS Françoise	Mathématiques Pures
MM.	BERTRANDIAS Jean-Paul	Mathématiques Pures
	BEZES Henri	Pathologie Chirurgicale
	BLAMBERT Maurice	Mathématiques Pures
	BOLLIET Louis	Informatique (I.U.T. B)
	BONNET Georges	Electrotechnique
	BONNET Jean-Louis	Clinique Ophtalmologique
	BONNET-EYMARD Joseph	Clinique Gastro-entérologique
Mme	BONNIER Marie-Jeanne	Chimie Générale
MM.	BOUCHERLE André	Chimie et Toxicologie
	BOUCHEZ Robert	Physique nucléaire
	BOUSSARD Jean-Claude	Mathématiques Appliquées
	BRAVARD Yves	Géographie
	CABANEL Guy	Clinique Rhumatologique et Hydrologique
	CALAS François	Anatomie
	CARLIER Georges	Biologie Végétale
	CARRAZ Gilbert	Biologie Animale et Pharmacodynamie
	CAU Gabriel	Médecine Légale et Toxicologie
	CAUQUIS Georges	Chimie Organique
	CHABAUTY Claude	Mathématiques Pures
	CHARACHON Robert	Clinique Oto-Rhino-Laryngologique
	CHATEAU Robert	Clinique de Neurologie
	CHIBON Pierre	Biologie Animale
	COEUR André	Pharmacie Chimique et Chimie Analytique
	CONTAMIN Robert	Clinique Gynécologique
	COUDERC Pierre	Anatomie Pathologique
	CRAYA Antoine	Mécanique
Mme	DEBELMAS Anne-Marie	Matière Médicale
MM.	DEBELMAS Jacques	Géologie Générale
	DEGRANGE Charles	Zoologie
	DELORMAS Pierre	Pneumo-Phtisiologie
	DEPORTES Charles	Chimie Minérale
	DESRE Pierre	Métallurgie

MM.	DESSAUX Georges	Physiologie Animale
	DODU Jacques	Mécanique Appliquée (I.U.T. A)
	DOLIQUE Jean-Michel	Physique des Plasmas
	DREYFUS Bernard	Thermodynamique
	DUCROS Pierre	Cristallographie
	DUGOIS Pierre	Clinique de Dermatologie et Syphiligraphie
	GAGNAIRE Didier	Chimie Physique
	GALLISSOT François	Mathématiques Pures
	GALVANI Octave	Mathématiques Pures
	GASTINEL Noël	Analyse Numérique
	GAVEND Michel	Pharmacologie
	GEINDRE Michel	Electroradiologie
	GERBER Robert	Mathématiques Pures
	GERMAIN Jean-Pierre	Mécanique
	GIRAUD Pierre	Géologie
	JANIN Bernard	Géographie
	KAHANE André	Physique Générale
	KLEIN Joseph	Mathématiques Pures
	KOSZUL Jean-Louis	Mathématiques Pures
	KRAVTCHENKO Julien	Mécanique
	KUNTZMANN Jean	Mathématiques Appliquées
	LACAZE Albert	Thermodynamique
	LACHARME Jean	Biologie Végétale
Mme	LAJZEROWICZ Janine	Physique
MM.	LAJZEROWICZ Joseph	Physique
	LATREILLE René	Chirurgie Générale
	LATURAZE Jean	Biochimie Pharmaceutique
	LAURENT Pierre-Jean	Mathématiques Appliquées
	LEDRU Jean	Clinique Médicale B
	LLIBOUTRY Louis	Géophysique
	LOISEAUX Pierre	Sciences Nucléaires
	LONGEQUEUE Jean-Pierre	Physique Nucléaires
	LOUP Jean	Géographie
Melle	LUTZ Elisabeth	Mathématiques Pures
	MALGRANGE Bernard	Mathématiques Pures
	BOUTET DE MONVEL Louis	Mathématiques Pures
	MALINAS Yves	Clinique Obstétricale
	MARTIN-NOEL Pierre	Seméiologie médicale
	MAZARE Yves	Clinique Médicale A
	MICHEL Robert	Minéralogie et Pétrographie
	MICOUD Max	Clinique Maladies Infectieuses
	MOURIQUAND Claude	Histologie
	MOUSSA André	Chimie Nucléaire
	MULLER Jean-Michel	Thérapeutique (Néphrologie)
	NEEL Louis	Physique du solide
	OZENDA Paul	Botanique
	PAYAN Jean-Jacques	Mathématiques Pures
	PEBAY-PEYROULA Jean-Claude	Physique
	RASSAT André	Chimie Systématique
	RENARD Michel	Thermodynamique
	RINALDI Renaud	Physique
	DE ROUGEMONT Jacques	Neuro-Chirurgie
	SEIGNEURIN Raymond	Microbiologie et Hygiène
	SENGEL Philippe	Zoologie
	SIBILLE Robert	Construction Mécanique (I.U.T. A)
	SOUTIF Michel	Physique Générale
	TANCHE Maurice	Physiologie

MM.	TRAYNARD Philippe	Chimie Générale
	VAILLANT François	Zoologie
	VALENTIN Jacques	Physique Nucléaire
	VAUQUOIS Bernard	Calcul Electronique
Mme	VERAIN Alice	Pharmacie Galénique
MM.	VERAIN André	Physique
	VEYRET Paul	Géographie
	VIGNAIS Pierre	Biochimie médicale
	YOCCOZ Jean	Physique Nucléaire Théorique

PROFESSEURS ASSOCIES

MM.	CLARK Gilbert	Spectrométrie Physique
	CRABBE Pierre	CERMO
	ENGLMAN Robert	Spectrométrie Physique
	HOLTZBERG Frédéric	Basses Températures
	ROST Ernest	Sciences Nucléaires

PROFESSEURS SANS CHAIRE

Melle	AGNIUS-DELORD Claudine	Physique Pharmaceutique
	ALARY Josette	Chimie Analytique
MM.	AMBROISE-THOMAS Pierre	Parasitologie
	BELORIZKY Elie	Physique
	BENZAKEN Claude	Mathématiques Appliquées
	BIAREZ Jean-Pierre	Mécanique
	BILLET Jean	Géographie
	BOUCHET Yves	Anatomie
	BRUGEL Lucien	Energétique (I.U.T. A)
	BUISSON René	Physique (I.U.T. A)
	CONTE René	Physique (I.U.T. A)
	DEPASSEL Roger	Mécanique des Fluides
	GAUTHIER Yves	Sciences Biologiques
	GAUTRON René	Chimie
	GIDON Paul	Géologie et Minéralogie
	GLENAT René	Chimie organique
	GROULADE Joseph	Biochimie Médicale
	HACQUES Gérard	Calcul Numérique
	HOLLARD Daniel	Hématologie
	HUGONOT Robert	Hygiène et Médecine Préventive
	IDELMAN Simon	Physiologie Animale
	JOLY Jean-René	Mathématiques Pures
	JULLIEN Pierre	Mathématiques Appliquées
Mme	KAHANE Josette	Physique
MM.	KUHN Gérard	Physique (I.U.T. A)
	LE ROY Philippe	Mécanique (I.U.T. A)
	LUU DUC Cuong	Chimie Organique
	MAYNARD Roger	Physique du Solide
	PELMONT Jean	Biochimie
	PERRIAUX Jean-Jacques	Géologie et Minéralogie
	PSISTER Jean-Claude	Physique du Solide
Melle	PIERY Yvette	Physiologie Animale
MM.	RAYNAUD Hervé	M.I.A.G.
	REBECQ Jacques	Biologie (CUS)
	REVOL Michel	Urologie
	REYMOND Jean-Charles	Chirurgie Générale
	RICHARD Lucien	Biologie Végétale
Mme	RINAUDO Marguerite	Chimie Macromoléculaire
MM.	ROBERT André	Chimie Papetière

MM.	SARRAZIN Roger	Anatomie et Chirurgie
	SARROT-REYNAUD Jean	Géologie
	SIROT Louis	Chirurgie Générale
Mme	SOUTIF Jeanne	Physique Générale
MM.	STREGLITZ Paul	Anesthésiologie
	VIALON Pierre	Géologie
	VAN CUTSEM Bernard	Mathématiques Appliquées

MAITRES DE CONFERENCES ET MAITRES DE CONFERENCES AGREGES

MM.	AMBLARD Pierre	Dermatologie
	ARMAND Gilbert	Géographie
	ARMAND Yves	Chimie (I.U.T. A)
	BACHELOT Yvan	Endocrinologie
	BARGE Michel	Neuro chirurgie
	BARJOLLE Michel	MIAG
	BEGUIN Claude	Chimie organique
Mme	BERIEL Hélène	Pharmacodynamie
MM.	BOST Michel	Pédiatrie
	BOUCHARLAT Jacques	Psychiatrie adultes
Mme	BOUCHE Liane	Mathématiques (C.U.S.)
MM.	BRODEAU François	Mathématiques (I.U.T. B)
	BUTEL Jean	Orthopédie
	CHAMBAZ Edmond	Biochimie médicale
	CHAMPETIER Jean	Anatomie et Organogénèse
	CHARDON Michel	Géographie
	CHERADAME Hervé	Chimie Papetière
	CHIAVERINA Jean	Biologie Appliquée (EFP)
	COHEN-ADDAD Jean-Pierre	Spectrométrie Physique
	COLOMB Maurice	Biochimie Médicale
	CONTAMIN Charles	Chirurgie thoracique et cardio-vasculaire
	CORDONNIER Daniel	Néphrologie
	COULOMB Max	Radiologie
	CROUZET Guy	Radiologie
	CYROT Michel	Physique du solide
	DELOBEL Claude	M.I.A.G.
	DENIS Bernard	Cardiologie
	DOUCE Roland	Physiologie Végétale
	DUSSAUD René	Mathématiques (C.U.S.)
Mme	ETERRADOSSI Jacqueline	Physiologie
MM.	FAURE Jacques	Médecine Légale
	FAURE Gilbert	Urologie
	FONTAINE Jean-Marc	Mathématiques Pures
	GAUTIER Robert	Chirurgie Générale
	GENSAC Pierre	Botanique
	GIDON Maurice	Géologie
	GROS Yves	Physique (I.U.T. A)
	GUITTON Jacques	Chimie
	HICTER Pierre	Chimie
	IVANES Marcel	Electricité
	JALBERT Pierre	Histologie
	KOLODIE Lucien	Hématologie
	KRAKOWIAK Sacha	Mathématiques Appliquées
	LE NOC Pierre	Bactériologie-virologie
	LEROY Philippe	I.U.T. A
	MACHE Régis	Physiologie Végétale
	MAGNIN Robert	Hygiène et Médecine Préventive
	MALLION Jean-Michel	Médecine du Travail
	MARECHAL Jean	Mécanique (I.U.T. A)
	MARTIN-BOUYER Michel	Chimie (C.U.S.)

M.	MICHOULIER Jean	Physique (I.U.T. A)
Mme	MINIER Colette	Physique (I.U.T. A)
MM.	NEGRE Robert	Mécanique (I.U.T. A)
	NEMOZ Alain	Thermodynamique
	NOUGARET Marcel	Automatique (I.U.T.A)
	PARAMELLE Bernard	Pneumologie
	PECCOUD François	Analyse (I.U.T. B)
	PEFFEN René	Métallurgie (I.U.T. A)
	PERRET Jean	Neurologie
	PERRIER Guy	Géophysique - Glaciologie
	PHELIP Xavier	Rhumatologie
	RACHAIL Michel	Médecine Interne
	RACINET Claude	Gynécologie et Obstétrique
	RAMBAUD André	Hygiène et Hydrologie
	RAMBAUD Pierre	Pédiatrie
	Mme	RENAUDET Jacqueline
MM.	ROBERT Jean-Bernard	Chimie-Physique
	ROMIER Guy	Mathématiques (I.U.T. B)
	SHOM Jean-Claude	Chimie générale
	STOEBNER Pierre	Anatomie pathologique
	VROUSOS Constantin	Radiologie

MAITRE DE CONFERENCES ASSOCIES

M.	COLE Antony	Sciences Nucléaires
----	-------------	---------------------

CHARGE DE FONCTIONS DE MAITRE DE CONFERENCES

M.	JUNIEN-LAVILLAVROY Paul	O.R.L.
----	-------------------------	--------

Fait à SAINT MARTIN D'HERES,
DECEMBRE 1975.

Président : M. NEEL Louis
 Vice-Présidents : M. BENOIT Jean
 M. BONNETAIN Lucien

PROFESSEURS TITULAIRES

MM.	BENOIT Jean	Radioélectricité
	BESSON Jean	Electrochimie
	BLOCH Daniel	Physique du solide
	BONNETAIN Lucien	Chimie Minérale
	BONNIER Etienne	Electrochimie et Electrometallurgie
	BRISSENEAU Pierre	Physique du solide
	BUYLE-BODIN Maurice	Electronique
	COUMES André	Radioélectricité
	FELICI Noël	Electrostatique
	LESPINARD Georges	Mécanique
	MOREAU René	Mécanique
	PARIAUD Jean-Charles	Chimie-Physique
	PAUTHENET René	Physique du solide
	PERRET René	Servomécanismes
	POLOUJADOFF Michel	Electrotechnique
	SILBERT Robert	Mécanique des Fluides

PROFESSEURS ASSOCIES

MM.	RUPPERSBERG Albert, Henner	Chimie
	ROUXEL Roland	Automatique

PROFESSEURS SANS CHAIRE

MM.	BLIMAN Samuel	Electronique
	BOUVARD Maurice	Génie Mécanique
	COHEN Joseph	Electrotechnique
	DURAND Francis	Métallurgie
	FOULARD Claude	Automatique
	LACOME Jean-Louis	Géophysique
	LANCIA Roland	Electronique
	VEILLON Gérard	Informatique Fondamentale & Appliquée
	ZADWORNY François	Electronique

MAITRES DE CONFERENCES

MM.	ANCEAU François	Mathématiques Appliquées
	BOUDOURIS Georges	Radioélectricité
	CHARTIER Germain	Electronique
	GUYOT Pierre	Chimie Minérale
	IVANES Marcel	Electrotechnique
	JOUBERT Jean-Claude	Physique du solide
	MORET Roger	Electrotechnique Nucléaire
	PIERRARD Jean-Marie	Hydraulique
	ROBERT François	Analyse Numérique
	SABONNADIÈRE Jean-Claude	Informatique Fondamentale & Appliquée
Mme	SAUCIER Gabrièle	Informatique Fondamentale & Appliquée

MAITRE DE CONFERENCES ASSOCIE

M. LANDAU Ioan Automatique

CHERCHEURS DU C.N.R.S. (Directeurs et Maîtres de Recherche)

MM. FRUCHART Robert Directeur de Recherche
ANSARA Ibrahim Maître de Recherche
CARRE René Maître de Recherche
DRIOLE Jean Maître de Recherche
MATHIEU Jean-Claude Maître de Recherche
MUNIER Jacques Maître de Recherche

Je tiens à exprimer ma profonde reconnaissance à :

Monsieur le Professeur GASTINEL, qui a dirigé cette thèse et qui tout au long de ce travail, m'a constamment encouragée, et soutenue de son intérêt, de ses conseils. Par ses précieuses indications, ses qualités humaines, son enthousiasme, sa haute compétence, Monsieur GASTINEL m'a permis de mener à bien ce travail. Je voudrais qu'il trouve ici l'expression de ma respectueuse gratitude.

Monsieur le Professeur KUNTZMANN, qui m'a fait détacher à l'Institut de Mathématiques Appliquées de GRENOBLE, qu'il a créé et sans cesse développé, et qui a bien voulu accepter aujourd'hui de présider le Jury de cette thèse.

Monsieur le Professeur WILKINSON, qui a bien voulu examiner ce travail: Que Monsieur WILKINSON me permette de lui exprimer combien je suis sensible à l'honneur qu'il me fait en acceptant de faire partie du Jury.

Monsieur le Professeur HOCQUENGHEM, qui par la bienveillante attention qu'il m'a toujours accordée et par ses conseils et ses encouragements, m'a aidée dans la poursuite de ce travail.

Monsieur le Professeur VIGNES, pour les suggestions apportées à une partie de cette étude, et pour l'attention qu'il a bien voulu m'accorder.

J'adresse enfin mes vifs remerciements à Mademoiselle DAVID et à Mademoiselle LAURENT, du département de Mathématiques du C.N.A.M, pour le soin apporté à la dactylographie de cette thèse, et à Monsieur IGLESIAS, du Service tirage de GRENOBLE, pour la réalisation de ce document.

INTRODUCTION

La représentation, sur un calculateur, d'un nombre réel par un élément d'un sous-ensemble fini S de \mathbb{R} , l'utilisation résultante d'une arithmétique approchant l'arithmétique réelle, génèrent des erreurs d'arrondi. Nous nous sommes attachés au problème de la minimisation de telles erreurs :

- Nous avons d'abord comparé la précision de différents schémas de sommation . Ce premier aspect est un problème de "stratégie". Sa raison d'être en est la non-associativité des opérations de l'arithmétique approchée. Par exemple, il a été établi que :

a) Si dans un calcul de somme de série alternée à terme général décroissant, on effectue la sommation des termes, non dans l'ordre des indices croissants mais dans l'ordre inverse, tous les chiffres de la somme approchée sont significatifs.

b) Si l'on compare les deux schémas suivants de sommation de N nombres flottants positifs x_i ($1 \leq i \leq N$) :

Schéma 1 : A x_1 est ajouté x_2 , puis au résultat x_3 , et ainsi de suite (schéma usuel).

Schéma 2 : x_1 et x_2 , x_3 et x_4 , ..., x_{N-1} et x_N sont ajoutés 2 à 2 puis ces résultats intermédiaires sont ajoutés les uns aux autres selon le schéma 1.

L'erreur d'arrondi dans le schéma 2 est inférieure ou égale à l'erreur d'arrondi dans le schéma 1 augmentée d'une unité dans son poids le plus faible.

- Nous avons ensuite étudié le problème de la correction des erreurs :

Nous avons, au niveau de l'opération élémentaire -addition, soustraction, multiplication, division- exprimé l'erreur, sous la forme d'un produit de composition sur S , par des opérations de l'arithmétique approchée ; la mise en oeuvre d'une telle représentation pour la correction d'un algorithme général a ensuite été effectuée. Les différentes arithmétiques flottantes ont été successivement étudiées. Nous avons testé les résultats sur des exemples .

- Nous avons enfin étudié le problème du contrôle des erreurs d'arrondi dans le cadre d'une utilisation parallèle de deux arithmétiques flottantes. Cette approche du contrôle d'erreur, plus proche de la réalité que l'arithmétique d'intervalles, a donné des résultats concernant plusieurs algorithmes usuels. Par exemple, deux suites d'opérations flottantes ont été définies pour le calcul des inverses C_1 et C_2 d'une M-matrice A (ou d'une matrice tridiagonale symétrique définie positive) qui vérifient $C_1 \leq A^{-1} \leq C_2$.

Le chapitre I introduit les définitions nécessaires à notre étude. Est rappelée d'abord la notion d'arrondi de \mathbb{R} sur un sous-ensemble de lui-même, puis sont considérés différents types particuliers de sous-ensembles de \mathbb{R} .

Le but du chapitre II est d'exprimer les erreurs d'arrondi élémentaires en arithmétique à virgule flottante, à l'aide de cette même arithmétique.

Les opérations sont successivement étudiées, selon le schéma suivant :

- expressions, par produits de composition sur S , de l'erreur lorsque l'arithmétique flottante satisfait à certains axiomes. Sont ainsi généralisés ou complétés les travaux de T. J. DEKKER [5], D. E. KNUTH [15], O. MOLLER [28], G. W. VELTKAMP [40] .
- étude des schémas définissant, sur les ordinateurs usuels, les opérations flottantes ; en particulier, sont considérés pour l'addition qui admet une grande variété de définitions, la troncature sans, puis avec, chiffre de garde, et l'arrondi sur chiffre de garde ; dans chaque cas est examinée la double possibilité de représentation des nombres négatifs, par signe et valeur absolue ou en notation de complément à 2.

Les chapitres III et IV traitent des applications des résultats du chapitre précédent à la correction d'algorithmes. La première partie du chapitre III est consacrée à la définition et à l'étude de la convergence d'un algorithme pour le calcul, par des opérations flottantes et avec une précision maximum, d'une somme algébrique ou d'un produit scalaire.

La correction des erreurs d'arrondi dans un algorithme plus général est ensuite abordée sous plusieurs aspects :

- les possibilités d'obtention du problème perturbé dans l'analyse a posteriori des erreurs d'arrondi.
- l'utilisation de méthodes de raffinement itératif d'un résultat approché, la détermination exacte d'une fonction de l'erreur étant assurée.
- la résolution numérique de problèmes mal conditionnés.
- l'amélioration de la solution calculée par une méthode itérative d'un système d'équations linéaires.

L'objet du chapitre V est l'étude de la stabilité numérique des schémas de sommation . En particulier ,un théorème général de comparaison des schémas de sommation séquentiel et dichotomique est donné . L'étude du calcul numérique de la somme d'une série alternée à terme général décroissant est effectuée ; deux méthodes sont en outre proposées pour doubler la précision obtenue par l'algorithme de calcul défini plus haut.

Dans le chapitre VI, nous posons le problème suivant :

$(x,y) \in S^2$ $\star \in \{ + , - , \times , / \}$, l'opération approchée sur S de \star , $\textcircled{\star}$, peut légitimement associer à $x \star y$, l'un ou l'autre des 2 éléments de S encadrant ce nombre ; une suite G de N opérations sur \mathbb{R} étant donnée, il est alors possible d'associer à G , 2^N suites différentes d'opérations approchées. Peut-on déterminer, dans cet ensemble, deux suites calculant des solutions encadrant la solution exacte ?

Quel ensemble de solutions est obtenu lors de l'application des 2^N suites précédentes ?

Une réponse est apportée à la première question pour les algorithmes de la sommation, du produit scalaire, de l'évaluation numérique d'un polynôme, de l'inversion d'une M-matrice et d'une matrice tridiagonale, symétrique définie positive ainsi que, dans un cadre différent, pour la résolution d'un système triangulaire à matrice inverse positive.

Des théorèmes relatifs à l'ensemble des solutions calculées sont établis pour la somme et le produit scalaire.

Dans le chapitre VII, on présente une étude relative à la répartition de l'erreur d'arrondi élémentaire sur son domaine de définition, lorsqu'on considère cette erreur comme variable aléatoire.

En Annexe est donné un tableau des caractéristiques des arithmétiques à virgule flottante des calculateurs usuels .

CHAPITRE I

ARRONDIS SUR R

Les paragraphes de ce chapitre introduisent successivement les définitions nécessaires pour l'étude des arrondis.

I - ARRONDIS DEFINIS DE \mathbb{R} DANS UN SOUS-ENSEMBLE FINI DE \mathbb{R}

Ce paragraphe reprend, dans le cas des ensembles de réels associés à un sous-ensemble fini, les définitions formulées par U. KULISCH [18] pour l'étude générale des arrondis dans les ensembles ordonnés ; la terminologie choisie est toutefois un peu différente de celle de [18] en particulier pour la notion d'arrondi optimal.

Désignons par E , l'ensemble \mathbb{R} ou un intervalle donné de \mathbb{R} . Soit S un sous-ensemble fini, non vide, de E .

Définition I. 1.

On dira qu'une application $\square : E \rightarrow S$ est un arrondi de E dans S , si les deux axiomes suivants sont vérifiés :

$$\begin{aligned}
x \leq y &\implies \square x \leq \square y && \forall (x,y) \in E^2 \quad (\text{monotonie}) \\
\square x &= x && \forall x \in S
\end{aligned}$$

Définition I. 2.

Un arrondi $\square : E \rightarrow S$ sera dit G - dirigé (respectivement D - dirigé), si l'axiome suivant est vérifié :

$$\begin{aligned}
\square x &\leq x && , \quad \forall x \in E \\
(\text{resp. } \square x &\geq x && , \quad \forall x \in E)
\end{aligned}$$

S étant donné, l'arrondi G - dirigé (respectivement D - dirigé) de E dans S , s'il existe, est unique : Il sera noté ∇ (resp. Δ) [cf.18]. Par extension, $x \in E$ étant donné, dès que le nombre $\text{Max} \{z \mid z \in S, z \leq x\}$ (resp. $\text{Min} \{z \mid z \in S, z \geq x\}$) existe, nous le noterons ∇x (resp. Δx). Remarquons que la définition de l'arrondi entraîne, pour $x \in E$ donné, $\square x = \nabla x \vee \Delta x$, l'un au moins des nombres ∇x , Δx étant défini, puisque S est non vide.

Compositions d'arrondis

Nous rappellerons le théorème suivant de U. KULISCH, que nous énoncerons uniquement pour des ensembles réels :

Théorème [U. KULISCH].

Si T et S sont deux sous-ensembles finis de E vérifiant $T \subseteq S \subseteq E$ et tels que soient définis les arrondis G - dirigés $\nabla : E \rightarrow T$, $\nabla_1 : E \rightarrow S$, $\nabla_2 : S \rightarrow T$ (resp. les arrondis D - dirigés $\Delta : E \rightarrow T$, $\Delta_1 : E \rightarrow S$, $\Delta_2 : S \rightarrow T$), alors, pour tout élément x de E on a :

$$\nabla x = \nabla_2 (\nabla_1 x) \qquad \Delta x = \Delta_2 (\Delta_1 x)$$

On peut aussi énoncer, ∇ , ∇_1 , ∇_2 , Δ , Δ_1 , Δ_2 ayant les mêmes définitions que ci-dessus :

Lemme I.1.

Si T et S sont deux sous-ensembles finis de E vérifiant $T \subseteq S \subseteq E$, on a, pour $x \in E - T$, l'égalité $\Delta_2(\nabla_1 x) = \nabla_2(\Delta_1 x)$ si et seulement si un et un seul des arrondis G - dirigé et D - dirigé de x dans S , $\nabla_1 x$ ou $\Delta_1 x$, est élément de T .

Soit $x \in E - T$:

a) Si aucun des arrondis dirigés $\nabla_1 x$, $\Delta_1 x$ de x dans S , n'est élément de T , l'inégalité $\nabla x < \nabla_1 x < \Delta_1 x < \Delta x$ entraîne :

$$\Delta_2(\nabla_1 x) = \Delta x$$

$$\nabla_2(\Delta_1 x) = \nabla x.$$

b) Si $\nabla_1 x$ et $\Delta_1 x$ sont tous deux éléments de T , $\Delta_2(\nabla_1 x) = \nabla x$, $\nabla_2(\Delta_1 x) = \Delta x$.

c) Si $\nabla_1 x = \nabla x$ et $\Delta_1 x < \Delta x$, l'égalité $\Delta_2(\nabla_1 x) = \nabla_2(\Delta_1 x)$ est vérifiée.

On peut remarquer qu'en particulier, si S se déduit de T par adjonction des demi-sommes de deux éléments consécutifs de T , $\Delta_2(\nabla_1 x)$ et $\nabla_2(\Delta_1 x)$ prennent, pour tout élément x de E , la valeur de T la plus proche de x (ou éventuellement les 2 valeurs optimales).

II - L'ARITHMETIQUE "APPROCHEE"

A un arrondi défini de E dans S , il est possible d'associer une arithmétique sur l'ensemble S .

(Dans les définitions ci-dessous, $y \neq 0$ si l'opération $*$ représente la division)

Définition I.3.

Un arrondi \square étant défini de E dans S , si $*$ désigne l'une des opérations $+$, $-$, \times , $/$ définies sur \mathbb{R} , on appellera "opération induite de $*$ sur S par \square ", la loi \otimes_{\square} de composition interne sur S définie pour $(x,y) \in S^2$ par :

$$x \otimes_{\square} y = \square(x * y)$$

De façon plus générale, lorsqu'à l'opération $*$ appartenant à l'ensemble $\{ + , - , \times , / \}$ des opérations élémentaires définies sur \mathbb{R} , on associe une loi de composition interne sur S , notée \otimes et qui sera dite l'opération correspondant à $*$ sur S , on donnera les définitions ci-dessous, énoncées par T.J. DEKKER [5] :

Définition I.4.

L'opération \otimes : $S \times S \rightarrow S$ sera dite "correcte" si :
 $\forall (x,y) \in S^2 \quad x \otimes y = \text{Max} \{z | z \in S \wedge z \leq x * y\} \vee \text{Min} \{z | z \in S \wedge z \geq x * y\}$

Définition I.5.

L'opération \otimes : $S \times S \rightarrow S$ sera dite "optimale" si :
 $\forall (x,y) \in S^2 \quad |x * y - x \otimes y| = \text{Min}_{z \in S} |x * y - z|$

Enfin, nous donnerons la définition suivante :

Définition I.6.

L'opération \otimes : $S \times S \rightarrow S$ sera dite "I - dirigée" (respectivement "E - dirigée") pour le couple $(x,y) \in S^2$ si :

$$|x \otimes y| \leq |x * y|$$

(resp. $|x \otimes y| > |x * y|$) .

Nous utiliserons aussi, relativement à une opération \otimes : $S \times S \rightarrow S$, la terminologie d'opération "G - dirigée" (resp. "D - dirigée") dont la définition est semblable à la définition I.2 énoncée pour un arrondi. Cette notion sera aussi considérée relativement à un couple donné de S^2 .

Avec les notations précédentes, la définition d'une opération \oplus correcte s'écrit :

$$x \oplus y = \nabla (x * y) \vee \Delta (x * y)$$

Il vient immédiatement les énoncés suivants, relatifs aux opérations \oplus_{\square} , induites d'une opération $*$ par un arrondi \square :

L'opération $\oplus_{\square} : S \times S \rightarrow S$, induite de $*$ par un arrondi \square , est correcte.

Toute opération \oplus_{\square} , induite, par un arrondi, d'une opération $*$ commutative, est aussi commutative.

Notation :

Les opérations \oplus ne sont pas généralement associatives. Nous choisirons, pour toute la suite, la convention suivante : L'ordre des opérations \oplus sera défini de la gauche à la droite.

N éléments x_i ($1 \leq i \leq N$) de S étant donnés, l'écriture

$$x_1 \oplus x_2 \oplus x_3 \dots \oplus x_N$$

représente donc, par cette convention, le résultat :

$$(\dots ((x_1 \oplus x_2) \oplus x_3) \oplus \dots) \oplus x_N$$

Nous noterons encore, de façon plus concise, ce nombre par :

$$\begin{aligned} \oplus \sum_{i=1}^N x_i & \quad \text{si } \oplus = \oplus \\ \otimes \prod_{i=1}^N x_i & \quad \text{si } \oplus = \otimes \end{aligned}$$

III - QUELQUES TYPES DE SOUS-ENSEMBLES DE \mathbb{R}

S étant un sous-ensemble fini, non vide, de \mathbb{R} , S^+ (resp S^-) désignera l'ensemble des éléments positifs (resp. négatifs) de S .

Nous utiliserons la terminologie suivante :

Définition I.7.

S étant un sous-ensemble fini de \mathbb{R} et x un élément de S , on appelle prédécesseur de x dans S (resp. successeur de x dans S), et on note $(x)'$ (resp. $(x)''$) le nombre :

$$\begin{aligned} (x)' &= \max \{ y \mid y \in S \wedge y < x \} \\ (\text{resp. } (x)'' &= \min \{ y \mid y \in S \wedge y > x \}) \end{aligned}$$

La définition de $(x)'$ (resp. $(x)''$) suppose x différent du plus petit élément (resp. du plus grand élément) de S . Nous désignerons par R_S l'ensemble des réels compris entre ces deux éléments, et :

Définition I.8.

S étant un sous-ensemble fini de \mathbb{R} contenant 0, et x un élément de R_S , on appellera intervalle de précision de x sur S , et on notera $\varepsilon(x)$, le nombre strictement positif suivant :

$$\begin{aligned} \varepsilon(x) &= \Delta x - \nabla x && \text{si } x \notin S \\ \varepsilon(x) &= (x)'' - x && \text{si } x \in S^+ \\ \varepsilon(x) &= x - (x)' && \text{si } x \in S^- \end{aligned}$$

Pour l'élément 0, on précisera éventuellement si l'intervalle de précision est considéré à droite ou à gauche.

III.1. SOUS-ENSEMBLES DE $[-1, 1]$ A INTERVALLES CONSTANTS.

Un sous-ensemble S de $[-1, 1]$ à intervalles constants, est l'ensemble de $2P + 1$ éléments de la forme :

$$x = m \varepsilon \quad \text{avec } \varepsilon = \frac{1}{P} \text{ et } m \text{ entier vérifiant } -P \leq m \leq P$$

Les éléments 1 ou -1 peuvent être exclus de S .

La représentation de \mathbb{R} par un sous-ensemble de $[-1, 1]$ à intervalles constants, implique un cadrage des operandes avant exécution d'une opération \otimes .

L'addition ou la soustraction sur S , de 2 éléments de cet ensemble, n'entraîne aucune erreur s'il ne se produit pas de dépassement de capacité. Le produit sur \mathbb{R} de 2 éléments de S appartient à l'intervalle $[-1, 1]$ et nous énoncerons une propriété relative au produit, sur S , d'éléments de cet ensemble, par des opérateurs corrects de multiplication :

Une multiplication correcte sur S de 2 éléments x et y de S donne pour résultat l'un des deux arrondis $\nabla(xy)$ ou $\Delta(xy)$.

Etant donnés N éléments x_i ($1 \leq i \leq N$) de S , considérons l'ensemble des résultats obtenus par la multiplication dans l'ordre donné de ces éléments sur S , lorsque, à chaque pas i , la multiplication du produit partiel des i premiers éléments par x_{i+1} est réalisée par l'une ou l'autre des opérations \otimes_{∇} ou \otimes_{Δ} ; on peut énoncer :

Proposition I.1.

Si S est un sous-ensemble de $[-1, 1]$ à intervalles constants, étant donnés N éléments x_i ($1 \leq i \leq N$) de S , l'un au moins des 2^{N-1} produits :

$$x_1 \otimes_{\alpha_1} x_2 \otimes_{\alpha_2} x_3 \otimes_{\alpha_3} \dots \otimes_{\alpha_{N-1}} x_N = \otimes_{\alpha_i} \prod_{i=1}^N x_i$$

où $\alpha_i = \nabla \vee \Delta$,

est égal à $\nabla \left(\prod_{i=1}^N x_i \right)$, respectivement $\Delta \left(\prod_{i=1}^N x_i \right)$.

Il suffit de raisonner avec $x_i > 0 \quad \forall_i$

Par récurrence :

. La propriété est vraie au rang 2 puisque, par définition :

$$x_1 \otimes_{\alpha_1} x_2 = \nabla (x_1 x_2) \vee \Delta (x_1 x_2)$$

. Supposons la vérifiée au rang k , c'est-à-dire supposons que dans l'ensemble des éléments $\otimes_{\alpha_i} \prod_{i=1}^k x_i$, il en existe au moins deux, éventuellement de même valeur, que nous noterons $\tilde{\pi}_k$ et $\tilde{\pi}'_k$ vérifiant :

$$\tilde{\pi}_k = \nabla \left(\prod_{i=1}^k x_i \right) \quad , \quad \tilde{\pi}'_k = \Delta \left(\prod_{i=1}^k x_i \right) .$$

De :

$$\tilde{\pi}'_k x_{k+1} - \tilde{\pi}_k x_{k+1} = x_{k+1} \varepsilon < \varepsilon \quad \text{si } x_{k+1} \neq 1$$

on déduit qu'il existe au plus un élément de S sur l'intervalle fermé $[\tilde{\pi}_k x_{k+1}, \tilde{\pi}'_k x_{k+1}]$.

L'égalité $\otimes_{\alpha_i} \prod_{i=1}^{k+1} x_i = \nabla \left(\prod_{i=1}^{k+1} x_i \right)$ est obtenue par exemple avec les

α_i ($i = 1, \dots, k-1$) ayant défini $\tilde{\pi}_k$, et le choix :

$\alpha_k = \nabla$ s'il n'existe pas d'éléments de S sur $[\tilde{\pi}_k x_{k+1}, \tilde{\pi}'_k x_{k+1}]$

ou si cet élément est supérieur strictement à $\prod_{i=1}^{k+1} x_i$,

$\alpha_k = \Delta$ sinon.

La propriété n'est plus vraie quels que soient les x_i de S , si l'on remplace l'une des opérations de multiplication par la division.

III.2. SOUS-ENSEMBLES DE \mathbb{R} A INTERVALLES NON DECROISSANTS

Définition I.9.

On dira que le sous-ensemble fini S de \mathbb{R} est à intervalles non décroissants, si :

$$|x| \leq |y| \implies \varepsilon(x) \leq \varepsilon(y) \quad \forall (x,y) \in S^2 \wedge x y > 0$$

Les opérateurs corrects d'addition arithmétique vérifient, sur un tel sous-ensemble de \mathbb{R} , la propriété suivante :

Proposition I.1'.

Si S est un sous-ensemble de \mathbb{R} à intervalles non décroissants, étant donnés N éléments x_i ($1 \leq i \leq N$) de S , de même signe, l'une au moins des 2^{N-1} sommes séquentielles :

$$x_1 \oplus_{\square_1} x_2 \oplus_{\square_2} x_3 \oplus_{\square_3} \dots \oplus_{\square_{N-1}} x_N = \oplus_{\square_i} \sum_{i=1}^N x_i \quad \text{où } \square_i = \nabla \vee \Delta, \\ \text{égale } \nabla \left(\sum_{i=1}^N x_i \right), \text{ respectivement } \Delta \left(\sum_{i=1}^N x_i \right).$$

Remarque :

Nous donnerons un contre-exemple grossier montrant que la propriété précédente n'est évidemment pas vérifiée pour l'addition arithmétique, sur un sous-ensemble quelconque de \mathbb{R} .

Soit le sous-ensemble S de $\mathbb{R} = \{0, 3, 6, 15, 21, 23, 25, 29, 31\}$; les 4 possibilités de réaliser sur S la somme $3 + 6 + 15 = 24$, par des opérations du type $3 \oplus_{\square_1} 6 \oplus_{\square_2} 15$ sont respectivement :

$$(3 \oplus_{\nabla} 6) \oplus_{\nabla} 15 = 21 = (3 \oplus_{\nabla} 6) \oplus_{\Delta} 15 \\ (3 \oplus_{\Delta} 6) \oplus_{\nabla} 15 = 29 \text{ et } (3 \oplus_{\Delta} 6) \oplus_{\Delta} 15 = 31$$

III.3. SOUS-ENSEMBLES DE \mathbb{R} DE TYPE FLOTTANT

Définition I.10.

Le sous-ensemble fini S de \mathbb{R} sera dit de type flottant si les deux axiomes suivants sont vérifiés :

- 1) $0 \in S$ et $x \in S \implies -x \in S$
- 2) pour tout x de S^+ tel que $\varepsilon(x)$ et $\varepsilon((x)')$ existent et vérifient $\varepsilon(x) \neq \varepsilon((x)')$, les rapports $\frac{x}{\varepsilon(x)}$, et $\frac{\varepsilon(x)}{\varepsilon((x)')}$ si $(x)' \neq 0$, sont entiers.

Un système de nombres à virgule flottante, de base b et de nombre de chiffres de mantisse s , est bien, en particulier, de ce type : Tout x positif du système vérifiant $\varepsilon(x) \neq \varepsilon((x)')$ satisfait à $\frac{x}{\varepsilon(x)} = b^{s-1}$; $\frac{\varepsilon(x)}{\varepsilon((x)'}) = b$ si $(x)' \neq 0$.

Un sous-ensemble de \mathbb{R} de type flottant est à intervalles non décroissants. La proposition I.1' est valable pour un tel système.

Tout élément x de S peut s'écrire sous la forme :

$x = \ell \varepsilon(x)$ où ℓ est un entier relatif.

L'addition correcte de 2 éléments de S vérifie le :

Lemme I.2.

Si S est un sous-ensemble de \mathbb{R} de type flottant, l'erreur entachant l'addition correcte de deux éléments x et y de S tels que $\varepsilon(y) \leq \varepsilon(x)$, $x + y \in \mathbb{R}_S$ vérifie :

$$|x + y - (x \oplus y)| \leq \varepsilon(x + y) - \varepsilon(y) \quad \text{avec}$$

$$= 0 \quad \text{si } \varepsilon(x + y) < \varepsilon(y)$$

Notation :

Etant donné un réel a , nous adopterons les notations suivantes :

$$[a] = \max \{n \mid n \in \mathbb{Z} \wedge n \leq a\}$$

$$[a] = \min \{n \mid n \in \mathbb{Z} \wedge n \geq a\}$$

Démonstration du lemme I.2. :

Par définition, si $u \in \mathbb{R}_S$, l'hypothèse que S est de type flottant entraîne :

$$\nabla u = \left\lfloor \frac{u}{\varepsilon(u)} \right\rfloor \varepsilon(u) \quad , \quad \Delta u = \left\lceil \frac{u}{\varepsilon(u)} \right\rceil \varepsilon(u)$$

Le lemme I.2 s'établit en considérant les 2 valeurs $\nabla(x + y)$ et $\Delta(x + y)$ que peut prendre l'addition correcte de x et y . Avec les notations $x = \ell \varepsilon(x)$, $y = m \varepsilon(y)$, $\ell, m \in \mathbb{Z}$

$$\nabla(x + y) = \left\lfloor \frac{\ell \varepsilon(x) + m \varepsilon(y)}{\varepsilon(x + y)} \right\rfloor \varepsilon(x + y)$$

Soit q l'entier relatif défini par :

$$\ell \varepsilon(x) + m \varepsilon(y) = q \varepsilon(x + y) + e \quad 0 \leq e < \varepsilon(x + y) .$$

L'hypothèse $\varepsilon(y) \leq \varepsilon(x)$ entraîne que $\ell \varepsilon(x) + m \varepsilon(y)$ est multiple de $\varepsilon(y)$.

Il en est de même de $\varepsilon(x + y)$ dès qu'est vérifié $\varepsilon(x + y) \geq \varepsilon(y)$ - voir note-
($x \oplus y = x + y$ sinon), et e satisfait à :

$$e = \lambda \varepsilon(y) \quad \text{avec } \lambda \in \mathbb{Z} \quad \text{et} \quad \left| \lambda \right| \leq \frac{\varepsilon(x + y)}{\varepsilon(y)} - 1 \quad \text{si } \varepsilon(x + y) \geq \varepsilon(y)$$

$$\lambda = 0 \quad \text{sinon}$$

D'où le lemme .

Définition.

Lorsque l'addition de deux éléments x et y de S tels que $\varepsilon(y) < \varepsilon(x)$, vérifie : $\varepsilon(x + y) < \varepsilon(y)$, il est dit se produire le phénomène de cancellation.

Note: Dans le cas où $0 < |x+y| < \varepsilon(0)$, la propriété $\varepsilon(x+y)$ multiple de $\varepsilon(y)$ peut n'être pas vérifiée. La majoration donnée de l'erreur doit être alors : $|x + y - (x \oplus y)| < \varepsilon(x + y) - \varepsilon(y)$.

-Toutefois la propriété ci-dessus est toujours vérifiée pour les systèmes de nombres à virgule flottante-

III.4. SYSTEMES DE NOMBRES A VIRGULE FLOTTANTE.

Rappel

Un système T_b^S de nombres à virgule flottante, de base b , de nombre de chiffres de mantisse s , est l'ensemble :

$$T_b^S = \{x | x = \varepsilon \left(\sum_{k=1}^s a_k b^{-k} \right) b^p \wedge \varepsilon = \bar{1} \wedge a_k \in \{0, 1, \dots, b-1\} \wedge p \in \mathbb{Z}, m_2 \leq p \leq m_1\}$$

où m_1 et m_2 sont des entiers relatifs donnés.

Nous considérons, afin de ne pas alourdir inutilement les démonstrations, uniquement une représentation normalisée ($a_1 \neq 0$ pour $x \neq 0$).

Tout élément de $T_b^S - \{0\}$ a alors une écriture unique.

Toutefois, si la représentation normalisée offre le maximum d'avantages, les résultats que nous obtiendrons par la suite, seraient semblables pour une écriture non normalisée des nombres.

Zéros de la représentation : $x = 0$ dès que $a_k = 0$, $\forall k$, mais l'exposant du zéro n'est pas défini. Il peut être nécessaire de lever cette ambiguïté.

Remarquant qu'on enregistre, dans la plupart des calculateurs, non pas l'exposant p du nombre à virgule flottante, mais le nombre $c = p +$ constante, appelé la "caractéristique" du nombre et permettant de ne pas introduire 2 positions "signe" différentes dans un même mot, nous désignerons dans ce cas par "zéro vrai" la représentation de ce nombre définie par $\varepsilon = +1$, $a_k = 0 \forall k$, $c = 0$; cette représentation se comporte dans tous les cas comme le nombre réel 0 , contrairement aux autres représentations, dites "zéros anormaux". Lorsqu'une ambiguïté pourrait se présenter, nous supposerons toujours qu'a été écrit un zéro vrai.

Notations :

1) La plus petite valeur absolue non nulle des éléments de T_b^s est $b^{m_2^{-1}}$: nous noterons ξ ce nombre.

2) Si $x \in T_b^s - \{0\}$, on peut écrire ce nombre :
 $x = \varepsilon \ell b^{p-s}$ où l'entier ℓ vérifie : $b^{s-1} \leq \ell < b^s$.

Nous noterons : $e(x) = p$ l'exposant de x .

La notation $e(u) = p$ sera étendue à l'entier relatif p associé au réel $u \in \mathbb{R}_{T_b^s} -]-\xi, \xi[$ par : $b^{p-1} \leq u < b^p$.

L'exposant d'un réel u tel que $|u| < \xi$ sera pris inférieur à tout autre exposant.

3) Lorsqu'aucune ambiguïté n'apparaîtra, nous noterons simplement T le système T_b^s de nombres à virgule flottante de base b , et de nombre de chiffres de mantisse s .

Enfin, on peut énoncer pour un système de nombres à virgule flottante en base 2 :

Proposition I.1''.

Si T_2^s est un système de nombres à virgule flottante en base 2, étant donnés N éléments x_i ($1 \leq i \leq N$) de T_2^s de même signe et $N - 1$ opérations \star_i d'addition, de multiplication ou de division, l'un au moins des éléments :

$$x_1 \overset{\star_1}{\square_1} x_2 \overset{\star_2}{\square_2} x_3 \dots \overset{\star_{N-1}}{\square_{N-1}} x_N$$

où $\square_i = \nabla$ ou Δ ,

égale $\nabla (x_1 \star_1 x_2 \dots \star_{N-1} x_N)$, respectivement $\Delta (x_1 \star_1 x_2 \dots \star_{N-1} x_N)$.

Reprenons le raisonnement par récurrence de la proposition I.1, en supposant au rang k l'existence de deux suites d'arrondis

\square_i et \square'_i ($i = 1, \dots, k-1$) permettant le calcul sur

T_2^s de $\tilde{r}_k = \nabla r_k$ et $\tilde{r}'_k = \Delta r_k$ où $r_k = x_1 \star_1 x_2 \dots \star_{k-1} x_k$.

Notant $e(\nabla r_k) = \pi_k$ et $x_{k+1} = \ell b^{p_{k+1}-s}$, il vient selon la nature

de l'opération \star_k :

1) $\star_k = \times$

$$\begin{aligned} \tilde{r}'_k x_{k+1} - \tilde{r}_k x_{k+1} &= \ell b^{-s} \varepsilon(r_{k+1}) < \varepsilon(r_{k+1}) \text{ si } e(r_{k+1}) = \pi_k + p_{k+1} \\ &= \ell b^{-s+1} \varepsilon(r_{k+1}) < 2 \varepsilon(r_{k+1}) \text{ si } e(r_{k+1}) = \pi_k + p_{k+1}^{-1} \wedge b=2 \end{aligned}$$

Il existe au plus 2 éléments de T_2^S sur l'intervalle fermé $[\tilde{r}_k x_{k+1}, \tilde{r}'_k x_{k+1}]$ et il est toujours possible de définir un choix de \square_k qui, à partir de \tilde{r}_k ou de \tilde{r}'_k , calcule ∇r_{k+1} par exemple.

2) $\star_k = /$

$$\begin{aligned} \tilde{r}'_k / x_{k+1} - \tilde{r}_k / x_{k+1} &= (b^s / \ell) \varepsilon(r_{k+1}) < 2 \varepsilon(r_{k+1}) \text{ si } e(r_{k+1}) = \pi_k - p_{k+1} \wedge b=2 \\ &= (b^{s-1} / \ell) \varepsilon(r_{k+1}) \leq \varepsilon(r_{k+1}) \text{ si } e(r_{k+1}) = \pi_k - p_{k+1} + 1 \end{aligned}$$

On remarquera que $\ell = b^{s-1}$ implique la seconde hypothèse pour $e(r_{k+1})$ - et $e(\tilde{r}'_k / x_{k+1})$ - et on conclura comme précédemment.

Remarque

La proposition I.1" est encore vérifiée si l'expression arithmétique définie à partir des éléments x_i et des opérations \star_i - d'addition arithmétique, de multiplication ou de division - comporte des produits de composition partiels.

Soit en effet :

$$\begin{aligned} r &= (x_1 \star_1 x_2 \dots \star_{p-1} x_p) \star_p (x_{p+1} \star_{p+1} \dots \star_{N-1} x_N) \\ &= u \star_p v \end{aligned}$$

D'après ce qui précède, il est possible de choisir des suites d'arrondis, pour $i = 1, \dots, p-1$, puis, pour $i = p+1, \dots, N-1$, qui permettent le calcul sur T_2^S des éléments suivants :

$$\begin{aligned} \tilde{u} &= \nabla u & , & & \tilde{u}' &= \Delta u \\ \tilde{v} &= \nabla v & , & & \tilde{v}' &= \Delta v \end{aligned}$$

En considérant l'intervalle :

$$I = [\min(\tilde{u} \star_p \tilde{v}, \tilde{u}' \star_p \tilde{v}, \tilde{u} \star_p \tilde{v}', \tilde{u}' \star_p \tilde{v}'), \max(\tilde{u} \star_p \tilde{v}, \tilde{u}' \star_p \tilde{v}, \tilde{u} \star_p \tilde{v}', \tilde{u}' \star_p \tilde{v}')]]$$

comme une réunion d'intervalles du type précédemment étudié - c'est-à-dire, par exemple, $[\tilde{u} \star_p \tilde{v}, \tilde{u} \star_p \tilde{v}']$ - il suffit alors de remarquer que tous les éléments de T_2^S contenus dans I , sont obtenus au moins une fois pour l'ensemble des choix d'arrondis.

CHAPITRE II

CALCUL AUTOMATIQUE DE L'ERREUR D'ARRONDI ELEMENTAIRE
EN ARITHMETIQUE A VIRGULE FLOTTANTE

INTRODUCTION

Soient un système $T_b^s = T$ de nombres à virgule flottante, et les opérations d'addition, soustraction, multiplication, division, définies sur T , que nous noterons respectivement \oplus , \ominus , \otimes , \oslash .

Le but de ce chapitre est de déterminer, sous forme de produit de composition sur T par des opérations $\oplus, \ominus, \otimes, \oslash$, des expressions de l'erreur d'arrondi élémentaire:

$$X * Y = X \otimes Y \quad X \in T, Y \in T, * \in \{+, -, \times, /\}$$

L'opération de soustraction n'est pas considérée différemment de l'addition : nous supposerons que cette opération vérifie la propriété :

$$X \ominus Y = X \oplus (-Y) \quad \forall (X, Y) \in T^2$$

Afin de ne pas compliquer inutilement les résultats, nous n'avons pas donné, dans les cas d'overflows où l'erreur d'arrondi reste élément de T , l'écriture du produit de composition représentant l'erreur. Les formules écrites dans ce chapitre, ne sont donc valables que s'il n'y a pas dépassement de capacité, ce dernier phénomène entraînant d'ailleurs généralement, l'arrêt des calculs.

I - L'ERREUR D'ARRONDI DANS L'ADDITION

I. 1. THEOREME PRELIMINAIRE

Proposition II. 1.

Si l'opération $\oplus : T \times T \rightarrow T$ est correcte et possède la propriété suivante : $e(X) - e(Y) > s \implies X \oplus Y = X$ (C)
 alors, $\forall (X, Y) \in T^2$, l'erreur $E = X + Y - (X \oplus Y)$ est élément de $T \cup]-\xi, \xi[$.

Soient, en effet, X et Y deux éléments de T .

Notons p et q leurs exposants respectifs et supposons $p \geq q$.

Lorsque $X + Y \in R_T$, par le chapitre I , on peut affirmer que E vérifie :

$$E = \lambda \varepsilon(Y) = \lambda b^{q-s} \quad \text{où } \lambda \text{ est un entier relatif qui satisfait à :}$$

$$|\lambda| \leq \frac{\varepsilon(X+Y)}{\varepsilon(Y)} - 1 \quad \text{si } \varepsilon(X+Y) \geq \varepsilon(Y)$$

$$= 0 \quad \text{sinon}$$

. Si $p - q \leq s$, il vient $|\lambda| < b^s$

($e(X+Y) \leq p+1$ si $p-q < s$, $e(X+Y) \leq p$ si $p-q = s$).

Lorsque $\lambda \neq 0$, il existe un entier non négatif unique v tel que :
 $b^{s-v-1} \leq \lambda < b^{s-v}$

$$E = (\lambda b^v) \cdot b^{q-v-s} \quad \text{avec } b^{s-1} \leq \lambda b^v < b^s$$

$$q - v < m_2 \implies |E| < \xi \wedge E \notin T$$

$$q - v \geq m_2 \implies E \in T .$$

. Si $p - q > s$, on a $E = Y \in T$ en raison de la propriété (C) .

Si un overflow se produit, notant $|X| = l b^{p-s}$, $|Y| = m b^{q-s}$,
 il vient $|E| = (m - (b^s - 1)l) b^{q-s}$ où $p = m_1$, et $E \in T$ si
 $|E| \geq \xi$ d'où la proposition.

Etant donnés X et Y éléments de T , il sera donc possible si \oplus
 est correcte et vérifie (C) , de calculer $E = X + Y - (X \oplus Y)$ par
 des opérations internes sur T .

Remarque

Les résultats obtenus, concernant l'erreur d'addition élémentaire en arithmétique à virgule flottante, seront utiles par la suite :

X et Y étant deux éléments de T tels que $e(X) \geq e(Y)$, $X + Y \in \mathbb{R}_T$,

$E = X + Y - (X \oplus Y)$ vérifie :

$$E = \lambda b^{e(Y) - s} \quad \text{avec } \lambda \in \mathbb{Z} \quad \text{et :}$$

. si $e(Y) \geq m_2 + s - 1$

$$|\lambda| \leq b^{e(X+Y) - e(Y) - 1}$$

$$\leq b^{e(X \oplus Y) - e(Y) - 1} \quad \text{si } e(X+Y) \geq e(Y) \vee e(X \oplus Y) \geq e(Y)$$

$\lambda = 0$

si $e(X+Y) \leq e(Y) \vee e(X \oplus Y) \leq e(Y)$

. si $e(Y) < m_2 + s - 1$

Les mêmes formules sont valables, sauf en cas d'apparition d'un underflow ($0 < |X+Y| < \xi$), où l'on peut seulement affirmer $|E| < \xi$.

I. 2. EXPRESSION, PAR UN PRODUIT DE COMPOSITION SUR T, DE L'ERREUR D'ARRONDI

I. 2. 1. Première formule

Soient $X = \varepsilon l b^{p-s}$ et $Y = \varepsilon' m b^{q-s}$ deux éléments de T, et soit $E = X + Y - (X \oplus Y)$.

On peut énoncer :

Proposition II. 2.

Si l'opération $\oplus : T \times T \rightarrow T$ est :

correcte

possède la propriété (C)

I - dirigée pour tout couple $(X, Y) \in T^2 \wedge X Y > 0$,

Alors, si $|E| \notin]0, \xi[$, E vérifie :

$$E = -(X \oplus Y) \ominus X \ominus Y \quad (\text{II} \cdot 1) \quad \text{où } e(X) \geq e(Y)$$

Démonstration :

Posons :

$$S = X + Y$$

$$\tilde{S} = X \oplus Y$$

$$S = \tilde{S} + E$$

Supposons $p \geq q$.

Tout d'abord, remarquons que $p - q > s$ entraîne $\tilde{S} = X$ par (C), et (II. 1) conduit bien, dans ce cas, à $E = Y$.

La démonstration dans le cas général s'effectuera en prouvant que $\tilde{S} - X$ est élément de T ; en effet, l'opération \oplus étant correcte, il viendra :

$$\tilde{S} \ominus X = \tilde{S} - X$$

Alors, $(\tilde{S} \ominus X) - Y = -E$ étant élément de T d'après la proposition II. 1 on aura :

$$(\tilde{S} \ominus X) \ominus Y = (\tilde{S} \ominus X) - Y = -E$$

Ecrivaint :

$$S = (\varepsilon 1 + \varepsilon' m b^{q-p}) b^{p-s}$$

envisageons les différentes possibilités pouvant se présenter lors de l'opération $\tilde{S} \ominus X$:

$$1) e(S) = e(X)$$

Avec la seule hypothèse que \oplus est correcte :

$$\tilde{S} = \lfloor \varepsilon 1 + \varepsilon' m b^{q-p} \rfloor b^{p-s} \vee \lceil \varepsilon 1 + \varepsilon' m b^{q-p} \rceil b^{p-s},$$

$$\tilde{S} - X = \varepsilon' (\lfloor m b^{q-p} \rfloor \vee \lceil m b^{q-p} \rceil) b^{p-s} \text{ est élément de } T;$$

en effet, on a : $e(\tilde{S} - X) = q \vee q + 1$, ou éventuellement $\tilde{S} - X = 0$.

Il suffit alors de remarquer que $e(\tilde{S} - X) = q + 1$ ne pourrait provoquer d'overflow que pour $p = q = m_1$, mais alors dans ce cas $\tilde{S} = S$.

$$2) e(S) < e(X)$$

X et Y sont de signes différents.

Mais remarquons alors que :

$$q = p \vee p - 1 \implies \tilde{S} = S \vee 0 \text{ et dans ce cas } \tilde{S} - X \text{ est élément de } T.$$

$$q \leq p - 2 \implies e(S) = p - 1, \text{ et :}$$

$$\tilde{S} = \varepsilon \lfloor b1 - mb^{q-p+1} \rfloor b^{p-1-s} \vee \varepsilon \lceil b1 - mb^{q-p+1} \rceil b^{p-1-s}$$

Par suite, $\tilde{S} - X = -\varepsilon \left(\lfloor mb^{q-p+1} \rfloor \vee \lceil mb^{q-p+1} \rceil \right) b^{p-1-s}$ est élément de T .

$$3) \quad e(S) = e(X) + 1$$

(+) étant I- dirigée pour tout couple (X, Y) tel que $X Y > 0$:

$$\tilde{S} = \varepsilon \left\lfloor \frac{1 + m b^{q-p}}{b} \right\rfloor b^{p+1-s}$$

$$\tilde{S} - X = \varepsilon \left(b \left\lfloor \frac{1 + m b^{q-p}}{b} \right\rfloor - 1 \right) b^{p-s}$$

$e(\tilde{S} - X) = q \vee q - 1$, et $\tilde{S} - X$ est élément de T à condition que $|\tilde{S} - X| \geq \xi$.

En remarquant que l'éventualité $|\tilde{S} - X| < \xi$ ne peut se produire que si les conditions nécessaires $q = m_2$, $p - q \leq s - 2$ sont réalisées, et que dans ce cas l'erreur E entachant \tilde{S} est, en module, inférieure à ξ , on a encore, dans les conditions du théorème : $\tilde{S} - X \in T$.

Une autre proposition est donnée par :

Proposition II. 3.

Si l'opération (+) : $T \times T \longrightarrow T$ est :
correcte

possède la propriété (C)

si $b = 2$

Alors, si $|E| \notin]0, \xi[$, E vérifie (II. 1) où $e(X) \geq e(Y)$

Pour établir cette proposition, remarquons que l'hypothèse de I. direction de l'opération (+) pour tout couple $(X, Y) \in T^2_{\wedge} X Y > 0$, n'intervient dans la démonstration précédente que lorsque $e(S) = e(X) + 1$, pour éliminer le cas où \tilde{S} est donné par :

$$\tilde{S} = \varepsilon \left\lceil \frac{1 + m b^{q-p}}{b} \right\rceil b^{p+1-s}$$

$$\text{Dans ce cas, } \tilde{S} - X = \varepsilon \left(b \left\lceil \frac{1 + m b^{q-p}}{b} \right\rceil - 1 \right) b^{p-s}$$

$e(\tilde{S} - X) = q \vee q + 1$, et $\tilde{S} - X$ peut ne pas être élément de T lorsque $q = p$.

Or, si $q = p$, $E = -\varepsilon \varphi b^{p-s}$ où $\varphi \in \{0, 1, \dots, b-1\}$

On vérifie immédiatement que si $b = 2$, $|E| = b^{p-s}$ et $|Y| \leq (b^s - 1) b^{p-s}$ entraînent $\tilde{S} - X = Y - E \in T$.

Remarque

Les soustractions apparaissant dans la formule (II.1) s'effectuant "sans erreur" dans les hypothèses des propositions II. 2 et II. 3, cette formule peut être remplacée par :

$$E = Y \ominus (\tilde{S} \ominus X) \quad (\text{II.1'})$$

T. J. DEKKER [5] a montré que la formule (II.1') est valable dans le cas où $b = 3$, à la condition que l'opération \oplus soit optimale. Cet auteur a aussi établi cette formule pour de nombreuses autres arithmétiques à virgule flottante.

I.2.2. Généralisation d'un théorème de D. E. KNUTH

D. E. KNUTH [15] suggère, afin d'éliminer la condition non symétrique $e(X) \geq e(Y)$, la formule suivante qu'il établit dans l'hypothèse où \oplus est optimale :

$$E = (Y \ominus (\tilde{S} \ominus X)) \oplus (X \ominus (\tilde{S} \ominus (\tilde{S} \ominus X))).$$

Gardant tout d'abord cette condition $e(X) > e(Y)$, on établit la proposition suivante :

Proposition II. 4.

Si l'opération $\oplus : T \times T \rightarrow T$ est :
correcte
possède la propriété (C)

Alors, si $|E| \notin]0, \xi[$, E vérifie :

$$E = (Y \ominus (\tilde{S} \ominus X)) \oplus (X \ominus (\tilde{S} \ominus (\tilde{S} \ominus X))) \quad (\text{II.2}) \text{ où } e(X) > e(Y)$$

Démonstration :

La formule (II. 2) vient immédiatement en remarquant que :

* Dans les conditions de la proposition II. 2. , on a toujours :
 $\tilde{S} \ominus X = \tilde{S} - X$, ce qui entraîne, l'opération \oplus étant correcte :
 $X \ominus (\tilde{S} \ominus (\tilde{S} \ominus X)) = 0$

* Si $e(S) = e(X) + 1$, et que \oplus est E - dirigée pour le couple $(X, Y) \in T^2 \wedge XY > 0$, avec de plus $q = p$, rappelons que :

$$\tilde{S} - X = \varepsilon (b \lceil \frac{1+m}{b} \rceil - 1) b^{p-s}$$

Le cas litigieux $e(\tilde{S} - X) = p + 1$ conduit à l'une des deux valeurs :

$$\tilde{S} \ominus X = \varepsilon (\lceil \frac{1+m}{b} \rceil - \lfloor \frac{1}{b} \rfloor) b^{p+1-s} \quad \vee \quad \varepsilon (\lceil \frac{1+m}{b} \rceil - \lfloor \frac{1}{b} \rfloor) b^{p+1-s}$$

Posons :

$$\hat{S} - X = \hat{S} \oplus X + E'$$

$$E' = \varepsilon' \psi b^{p-s} \quad \text{où } \psi \in \{0, 1, \dots, b-1\}$$

D'où :

$$Y - (\hat{S} \oplus X) = E + E' = (-\varepsilon \varphi + \varepsilon' \psi) b^{p-s}$$

ce qui montre que ce nombre est élément de T dès que $E \in T$,
et

$$Y \oplus (\hat{S} \oplus X) = E + E'$$

$$\hat{S} - (\hat{S} \oplus X) = \varepsilon \left(\left\lceil \frac{1}{b} \right\rceil \vee \left\lfloor \frac{1}{b} \right\rfloor \right) b^{p+1-s} \text{ est élément de } T,$$

donc :

$$\hat{S} \oplus (\hat{S} \oplus X) = X + E'$$

et E' étant élément de T lorsque $E \in T$, on en déduit (II. 2).

Rappelant la définition suivante, due à T. J. DEKKER :

Définition

L'opération \oplus : $T \times T \rightarrow T$ est dite "proprement tronquante" si elle est commutative, et si pour tout couple (X, Y) de T^2 satisfaisant $|X| \geq |Y|$, on a :

$$\begin{aligned} X \oplus Y &= \text{Sup} \{Z \mid Z \in T \wedge Z \leq X + Y\} && \text{si } Y \geq 0 \\ X \oplus Y &= \text{Inf} \{Z \mid Z \in T \wedge Z \geq X + Y\} && \text{si } Y < 0 \end{aligned}$$

On peut énoncer :

Proposition II. 5

Si l'opération \oplus : $T \times T \rightarrow T$ est proprement tronquante, Alors, si aucun underflow n'apparaît lors de la succession d'opérations définie par (II. 2), E vérifie (II. 2), quelles que soient les valeurs relatives de X et Y .

Démonstration :

Remarquons tout d'abord que l'hypothèse : \oplus proprement tronquante, entraîne que (C) est vérifiée.

Si donc $e(X) \geq e(Y)$, (II. 2) est vérifiée dès que $E \in T$ d'après la proposition précédente.

Démontrons la proposition dans le cas où $e(X) < e(Y)$, ou plutôt, gardant l'hypothèse $e(X) > e(Y)$ pour la commodité des notations, établissons la formule suivante :

$$E = (X \ominus (\tilde{S} \ominus Y)) \oplus (Y \ominus (\tilde{S} \ominus (\tilde{S} \ominus Y))) \quad (\text{II.3})$$

Posons :

$$\tilde{S} - Y = \tilde{S} \ominus Y + E'$$

Il vient :

$$X - (\tilde{S} \ominus Y) = E + E'$$

$$\tilde{S} - (\tilde{S} \ominus Y) = Y + E'$$

$$Y - (\tilde{S} - (\tilde{S} \ominus Y)) = - E'$$

La proposition II. 1 entraîne que E' est élément de T à la condition

$$|E'| \geq \xi \vee E' = 0$$

Pour prouver (II. 3) , il suffit, l'opération \oplus étant correcte, d'établir que $E + E'$ et $Y + E'$ sont éléments de T .

D'après la démonstration de la proposition II. 1 , E et E' s'écrivent :

$$E = M b^{q-s}$$

$$E' = N b^{q-s}$$

où M et N sont des entiers relatifs.

(en effet, pour $\tilde{S} \neq S$, on a toujours $e(\tilde{S}) \geq e(Y)$).

De plus, \oplus étant proprement tronquante, la troncature de la somme $X + Y$ ($|X| \geq |Y|$) lorsque $X + Y \notin T$, se fait dans la direction de $- Y$.

E a donc le signe de Y .

E' , lorsque $E' \neq 0$, a le signe de $- Y$.

En effet, l'inégalité $|\tilde{S}| < |Y|$ implique $\tilde{S} = S \vee 0$ et $E' = 0'$.

De l'écriture de E et E' , il vient immédiatement que $E + E'$ et $Y + E'$ sont éléments de T à condition que leur formation ne donne pas lieu à l'apparition d'un underflow.

Nous verrons d'ailleurs par la suite qu'une telle éventualité ne peut pas se produire lorsque $E \in T$.

Il est, par contre, possible que E' vérifie $0 < |E'| < \xi$ bien que $|E| > \xi$.

Exemple :

$$\begin{aligned} b = 10, \quad s = 6, \quad X &= 0,100000 \cdot 10^{m_2+6}, \quad Y = 0,234986 \cdot 10^{m_2+3} \\ \tilde{S} &= 0,100234 \cdot 10^{m_2+6}, \quad E = 0,986000 \cdot 10^{m_2} \\ \tilde{S} \ominus Y &= 0,999991 \cdot 10^{m_2+5}, \quad E' = -0,086000 \cdot 10^{m_2} \end{aligned}$$

Proposition II. 6.

Si l'opération $\oplus : T \times T \rightarrow T$ est :
correcte
possède la propriété (C)

Alors, si aucun underflow n'apparait lors de la succession d'opérations définie par (II. 2), E vérifie (II. 2), quelles que soient les valeurs relatives de X et Y.

Démonstration :

Démontrons (II. 3) avec les seules hypothèses que \oplus est correcte et possède la propriété (C).

E' étant défini comme précédemment, nous établirons la proposition en montrant que $E + E'$ et $Y + E'$ sont éléments de T.

La démonstration est une longue suite de cas à envisager.

Tout d'abord, remarquons : $p - q > s \vee E = 0 \implies$ (II. 3) vérifiée.

Dans le cas général, écrivons l'égalité :

$$\tilde{S} - Y = X - E$$

et étudions les différentes possibilités en utilisant les résultats du § I. 2. 1 de ce chapitre.

1) $e(S) = p$

$$|E| < b^{p-s} \xrightarrow[\text{si } l \neq b^{s-1}]{\implies} \tilde{S} \ominus Y = X \vee X - \varepsilon'' b^{p-s} \text{ où } \varepsilon'' \text{ est le signe de } E.$$

Donc :

$$E' = -E \vee \varepsilon'' b^{p-s} - E$$

ce qui établit que $E + E' \in T$ lorsque $|E| \geq \xi$.

Il vient ensuite :

$$Y + E' = \varepsilon' (\lfloor m b^{q-p} \rfloor \vee \lceil m b^{q-p} \rceil) b^{p-s}$$

et $Y + E'$ est bien élément de T.

Dans le cas particulier $l = b^{s-1}$ et E du signe de X, on vérifie les 2 égalités suivantes :

$$E + E' = \lambda b^{p-1-s} \text{ où } \lambda \text{ est un entier relatif avec } 0 \leq |\lambda| \leq b$$

$E + E' \in T$ lorsque E est élément de T.

$$Y + E' = \varepsilon' (\lfloor m b^{q-p+1} \rfloor \vee \lceil m b^{q-p+1} \rceil) b^{p-1-s}$$

2) $e(S) = p-1$

Dans ce cas :

$$|E + E'| = 0 \vee b^{p-s} \vee b^{p-1-s} \quad (1 = b^{s-1})$$

et

$E + E' \in T$ lorsque $E \in T$.

On a :

$$Y + E' = -\varepsilon (\lfloor m b^{q-p+1} \rfloor + b \theta) b^{p-1-s} \vee -\varepsilon (\lceil m b^{q-p+1} \rceil + b \varphi) b^{p-1-s}$$

avec

$$\theta = 0 \text{ ou } 1, \varphi = 0 \text{ ou } -1.$$

On vérifie bien que $Y + E'$ est élément de T lorsque $|E| \geq \xi$:
 en effet, les conditions $\varphi = -1$ et $q = m_2$ qui, dans certains cas, peuvent entraîner un underflow, impliquerait $|E| < \xi$.

3) $e(S) = p + 1$

On vérifie d'abord, en utilisant l'inégalité $|E| < b^{p+1-s}$ et en considérant le cas particulier $b = 2$, que l'on a toujours $e(X - E) \geq p$.

* si $e(X - E) = p$, il vient :

$$E + E' = \lambda b^{p-s}$$

où λ est un entier relatif avec $0 \leq |\lambda| \leq b$,

et $E + E' \in T$ si $E \in T$.

D'autre part, les résultats vus au début de ce chapitre conduisent à :

$$Y + E' = \varepsilon (\lfloor m b^{q-p} \rfloor \vee \lceil m b^{q-p} \rceil) b^{p-s}$$

et $Y + E'$ est élément de T .

* si $e(X - E) = p + 1$, il vient :

$$E + E' = \lambda b^{p-s}$$

où λ est un entier relatif avec $0 < |\lambda| \leq 2b - 1$

et

$$Y + E' = \varepsilon (\lfloor m b^{q-p-1} \rfloor \vee \lceil m b^{q-p-1} \rceil) b^{p+1-s}$$

La proposition est donc démontrée.

I. 3. QUELQUES EXEMPLES DE REGLES D'ADDITION EN ARITHMETIQUE A VIRGULE FLOTTANTE

Nous rappelons, dans ce paragraphe, différentes règles permettant de définir l'addition, en arithmétique à virgule flottante normalisée, de deux éléments $X = \varepsilon l b^{p-s}$ et $Y = \varepsilon' m b^{q-s}$ ($e(X) \geq e(Y)$) de T .

I. 3. 1. La troncature sans chiffre de garde

Rappel.

La règle définissant l'addition sur T est une règle de troncature sans chiffre de garde, si l'opération d'addition s'effectue suivant le schéma suivant :

- 1) L'exposant provisoire de la somme $X \oplus Y$ est pris égal à $e(X)$. Si $e(Y) < e(X)$, Y est cadré, ce cadrage s'effectuant en déplaçant de $e(X) - e(Y)$ positions vers la droite la mantisse de ce nombre, et en tronquant les $e(X) - e(Y)$ derniers chiffres de cette mantisse.
- 2) Les mantisses, qui sont à présent alignées, sont ajoutées algébriquement, puis le résultat obtenu est normalisé, et sa mantisse est tronquée à s chiffres lorsque nécessaire.

I. 3. 1. a. Les nombres négatifs sont représentés par leur signe et leur valeur absolue.

Le cadrage de Y s'est effectué en remplaçant Y par $Y^t = \epsilon' [m b^{q-p}] b^{p-s}$

$$X \oplus Y = (\epsilon 1 + \epsilon' [m b^{q-p}]) b^{p-s} \quad \vee \quad 0 \quad \text{si } |X + Y^t| < \xi \quad \text{lorsque } e(S) \leq p$$

$$X \oplus Y = \epsilon \left[\frac{1 + m b^{q-p}}{b} \right] b^{p+1-s} \quad \text{lorsque } e(S) = p+1$$

Donc :

Propriété.

1) L'opération \oplus pour la troncature sans chiffre de garde, les nombres négatifs étant représentés avec signe et valeur absolue, n'est pas correcte, mais l'erreur E vérifie :

$$\begin{aligned} |E| < b^{p-s} & \text{ où } P = \max\{e(X \oplus Y), e(X), e(Y)\} \text{ si } X \oplus Y \neq 0 \\ |E| < \xi & \text{ si } X \oplus Y = 0 \end{aligned} \quad (\text{II. 4})$$

2) E est élément de T dès que $|E| \notin]0, \xi[$

3) \oplus est :

I- dirigée pour tout couple $(X, Y) \in T^2 \wedge X Y > 0$

E- dirigée pour tout couple $(X, Y) \in T^2 \wedge X Y < 0 \wedge |X + Y^t| > \xi$

4) La symétrie $(X \rightarrow -X)$ est un automorphisme de T .

(II. 4) vient immédiatement.

On a la seconde partie de la propriété en remarquant que $|E|$ est du type λb^{q-s} où λ est entier, et :

$$p - q \leq s - 1 \implies \lambda < b^s$$

$$p - q \geq s \implies E = Y$$

3) et 4) sont évidents.

Rappelons que si l'addition de deux nombres de même signe est correcte, la soustraction arithmétique est de faible précision, puisqu'elle peut engendrer, lorsque se produit le phénomène de cancellation, des erreurs du même ordre de grandeur que le résultat obtenu.

Ainsi :

$$\text{si } b = 10, \quad s = 6, \quad X = 0,100000, \quad Y = 0,999999 \cdot 10^{-1}$$

on a :

$$X \ominus Y = 10^{-6} \quad \text{et} \quad E = -9 \cdot 10^{-7}.$$

Toutefois $p - q \geq 2$ assure $e(X \ominus Y) \geq p - 1$, et au plus, deux chiffres du résultat sont entachés d'erreur.

Les théorèmes démontrés au § I. 2. de ce chapitre, ne s'appliquent pas ici, mais on peut énoncer :

Proposition II. 7.

Si la règle définissant l'opération \oplus est une troncature sans chiffre de garde, les nombres négatifs étant représentés avec signe et valeur absolue, l'erreur E entachant $\tilde{S} = X \oplus Y$ est donnée,

si $|E| \notin [0, \xi[$, par la suite d'opérations :

$$U = \tilde{S} \ominus X$$

$$FX = U \ominus \tilde{S} \oplus X$$

$$E = - (U \ominus b \cdot FX \oplus (b \cdot FX \ominus FX) \ominus Y) \quad (\text{II. 5})$$

avec la condition $|X| \geq |Y|$.

Démonstration :

Notation - lors de l'addition $A \oplus B$ de 2 éléments A et B de T , ($|A| \geq |B|$), nous sommes conduits à examiner :

- l'erreur due au cadrage de B , et nous noterons :

$$(A + B)_i = A + B^t$$

- l'erreur due à la troncature après normalisation de ce résultat intermédiaire, c'est-à-dire l'apparition éventuelle d'une retenue.

\tilde{S} , si $\tilde{S} \neq 0$, a le signe de X .

Distinguons 2 cas :

1) $e(\tilde{S}) \leq p$

Si $\tilde{S} = 0$, l'erreur est, en module, inférieure à ξ , et (II. 5) conduit à $E = 0$, valeur convenable uniquement si $\tilde{S} = S$.

Si $\tilde{S} \neq 0$, la mantisse de \tilde{S} se cadrant sans erreur pour l'exposant p , et l'opération $\tilde{S} - X$ étant une soustraction arithmétique :

$$U = \tilde{S} \ominus X = (\tilde{S} - X)_i = \tilde{S} - X = \varepsilon' \lfloor m b^{q-p} \rfloor b^{p-s}$$

Ceci met en évidence $e(U) = q$. (sauf si $p-q > s$, mais alors par (II. 5), $E = Y$).

L'opération $U \ominus \tilde{S}$:

les mantisses de U et de \tilde{S} pouvant se cadrer sans erreur pour l'exposant p ($> \hat{a} e(\tilde{S})$ et $e(U)$),

$$(U - \tilde{S})_i = U - \tilde{S} = -\varepsilon' b^{p-s}$$

$$U \ominus \tilde{S} = -X$$

$$FX = 0 \text{ et } E = -(U \ominus Y)$$

$e(U) = e(Y)$ donc $(U - Y)_i = U - Y$, et $U \ominus Y$ s'effectuant sans erreur dès que $E \in T$, la proposition est vérifiée.

2) $e(\tilde{S}) = p + 1$

$$U = \tilde{S} \ominus X = \tilde{S} - X^t \text{ où } X^t = \varepsilon \lfloor \frac{1}{b} \rfloor b^{p+1-s}$$

Remarquons que si $m_2 = \lfloor m b^{q-p-1} \rfloor$, $U = \varepsilon (m_2 + \theta) b^{p+1-s}$ où $\theta = 0$ ou 1 , la deuxième alternative ayant toujours lieu lorsque $p - q = s - 1$, donc : $e(U) = q \vee q + 1$.

$$U \ominus \tilde{S} = U - \tilde{S} = -X^t$$

$$FX = X \ominus X^t, \text{ donc :}$$

FX égale le dernier chiffre x_s de la mantisse de X , $FX = X - X^t$, à la condition que $p - s + 1 \geq m_2$, sinon $FX = 0$ même si $X - X^t \neq 0$ mais dans ce cas, $|E| < \xi$. Si ce cas particulier \mathcal{P} n'a pas lieu :

$$V = U \ominus b \cdot FX = U - b \cdot FX, \text{ d'après l'écriture de } U.$$

Il vient dans tous les cas

$$e(V) \leq q \text{ (} e(U) = q + 1 \text{ n'ayant lieu que si } x_s \neq 0 \text{)}.$$

Les mantisses de V et de $W = b \cdot FX \ominus FX = b \cdot FX - FX$ se cadrant sans erreur pour l'exposant p :

$$(V + W)_i = V + W = \tilde{S} - X$$

$$|\tilde{S} - X| \leq |Y| \implies V \oplus W = \tilde{S} - X$$

et

$$E = -((\tilde{S} - X) \ominus Y) = X + Y - \tilde{S}$$

dès que $E \in T$.

D'où la proposition.

Si \mathcal{P} a lieu, la formule (II. 5) prend la forme : $E = Y \ominus U$

$$\star e(U) = q \implies (Y - U)_i = [X + Y - \hat{S}] + [X^t - X]$$

les 2 quantités entre crochets étant de signes contraires et de module inférieur à ξ , $Y \ominus U = 0$.

$$\star e(U) = q + 1 \implies (Y - U)_i = Y^\zeta - U \quad \text{où } Y^\zeta = \epsilon \lfloor \frac{m}{b} \rfloor b^{q+1-s}$$

$$(Y - U)_i = [X + Y - \hat{S}] + [X^t - X + Y^\zeta - Y]$$

$$q \neq p \implies Y \ominus U = 0$$

mais les conditions $p = q = m_2 + s - 2 \wedge \lfloor \frac{m}{b} \rfloor = b^{s-1} \wedge x_s + y_s \geq b$

où y_s est le dernier chiffre de la mantisse de Y , entraînent :

$$Y \ominus U = -\epsilon \xi$$

alors que l'erreur est de module inférieur à ξ .

En conclusion, si \hat{E} est la valeur calculée par (II. 5) lorsque $|E| < \xi$, on peut seulement affirmer :

$$|E - \hat{E}| \leq 2\xi \left(\frac{b-1}{b}\right)$$

Remarque :

L'écriture $b \cdot FX \oplus (b \cdot FX \ominus FX)$ est imposée par le cas très particulier où $p = q \wedge \lfloor \frac{m}{b} \rfloor = b^{s-1} - 1 \wedge x_s + y_s \geq b$.

Ainsi :

$b = 10$, $s = 3$, $X = 0,999$, $Y = 0,991$ entraîneraient, avec les notations précédentes :

$U \ominus FX \ominus Y = 0,01$ introduisant l'erreur b^{p+1-s} , alors que $\hat{S} = S$

I. 3. 1. b. Les nombres négatifs sont représentés en notation de complément à 2 (base puissance de 2)

Dans la définition donnée, de l'addition pour la règle de troncature sans chiffre de garde, on doit remplacer :

- la mantisse des nombres par la représentation de ces mantisses en complément à 2.
- l'addition algébrique des mantisses par la somme arithmétique, modulo 2, de leurs représentations, sauf dans le cas où se produit un dépassement de capacité de mantisse auquel cas on tronque le résultat du chiffre de poids le plus faible.

Donc :

Lorsque Y est positif, $X \oplus Y$ prend les valeurs écrites au § précédent
 Lorsque Y est négatif, le cadrage de Y s'est effectué en remplaçant

$$Y \text{ par } Y^a = - \lceil m b^{q-p} \rceil b^{p-s} .$$

Il vient immédiatement :

$$X \oplus Y = (\varepsilon 1 - \lceil m b^{q-p} \rceil) b^{p-s} \quad \vee \quad 0 \quad \text{lorsque } e(S) \leq p$$

$$X \oplus Y = - \lceil \frac{1 + mb^{q-p}}{b} \rceil b^{p+1-s} \quad \text{lorsque } e(S) = p + 1$$

Propriété.

Lorsque l'opération \oplus est définie par la règle de troncature sans chiffre de garde, les nombres négatifs étant représentés en notation de complément à 2, on peut énoncer :

- 1) L'erreur E vérifie (II. 4) .
- 2) Si l'on complète la définition de \oplus par :
 $p - q \geq s \Rightarrow X \oplus Y = X$, quels que soient les signes de X et de Y ,

Alors E est élément de T dès que $|E| \notin]0, \xi[$

- 3) L'opération \oplus est G - dirigée $\forall (X, Y) \in T^2 \wedge X + Y \notin]-\xi, 0[$
- 4) La symétrie ($X \rightarrow -X$) n'est plus un automorphisme de T .

Expression de E par un produit de composition sur T :

La succession d'opérations de la formule (II. 5) ne définit plus E dans tous les cas. Il suffit, en effet, de considérer par exemple, le cas suivant :

$$X > 0, Y < 0, \text{ donc si } Y^a = - \lceil m b^{q-p} \rceil b^{p-s}, \quad \tilde{S} = X + Y^a ;$$

$$\text{Il vient } U = Y^a, \text{ mais le cas particulier où } \lceil m b^{q-p} \rceil = b^{s-p+q}$$

entraînant $e(U) = q + 1$, le chiffre de poids le plus faible de la mantisse de Y sera perdu lors du calcul $U \ominus Y$.

O. MOLLER [28] a montré que dans le cas où $b = 2$, les nombres négatifs étant représentés en complément à 2, l'erreur E est donnée par :

$$E = (Y \ominus (\tilde{S} \ominus X)) \oplus (X \ominus (\tilde{S} \ominus (\tilde{S} \ominus X))) \oplus (Y \ominus ((\tilde{S} \ominus X) \oplus (Y \ominus (\tilde{S} \ominus X))))$$

avec la condition $|X| \geq |Y|$,

la partie principale de l'erreur E étant constituée par les deux premiers termes $(Y \ominus (\tilde{S} \ominus X)) \oplus (X \ominus (\tilde{S} \ominus (\tilde{S} \ominus X)))$.

O. MOLLER a établi en outre que cette dernière écriture donne, sauf dans de rares cas particuliers, une bonne approximation de l'erreur commise même si $|X| < |Y|$.

I. 3. 2. La troncature avec chiffre de garde [21].

Définition :

La règle définissant l'addition sur T est une règle de troncature avec chiffre de garde, si l'opération d'addition de deux éléments X et Y de T ($e(X) \geq e(Y)$), s'effectue suivant le schéma suivant :

1) L'exposant provisoire de la somme $X \oplus Y$ est pris égal à $e(X)$.

X est écrit comme un nombre normalisé de $(s + 1)$ chiffres de mantisse, par addition d'un zéro après le chiffre de poids le plus faible de sa mantisse.

Si $e(Y) = e(X)$, il en est de même de Y , sinon Y est cadré, ce cadrage s'effectuant en déplaçant de $e(X) - e(Y)$ positions vers la droite la mantisse de ce nombre, et en tronquant les $e(X) - e(Y) - 1$ derniers chiffres de cette mantisse.

2) Les mantisses qui sont à présent alignées, sont ajoutées algébriquement, puis le résultat obtenu est normalisé, et enfin sa mantisse est tronquée à s chiffres.

I. 3. 2. a. Les nombres négatifs sont représentés par leur signe et leur valeur absolue .

X est exprimé sous la forme $\varepsilon (bl) b^{p-s-1}$; Y est "approché" sous la forme $\varepsilon' \lfloor m b^{q-p+1} \rfloor b^{p-s-1}$ (cette valeur étant toujours exacte si $q - p + 1 \geq 0$) .

Les valeurs prises par $\tilde{S} = X \oplus Y$ sont :

Si $\varepsilon \varepsilon' > 0$:

$$\tilde{S} = \varepsilon (1 + \lfloor m b^{q-p} \rfloor) b^{p-s} \quad \text{lorsque } e(\tilde{S}) = p$$

$$\tilde{S} = \varepsilon \lfloor \frac{1 + m b^{q-p}}{b} \rfloor b^{p+1-s} \quad \text{lorsque } e(\tilde{S}) = p + 1$$

Si $\varepsilon \varepsilon' < 0$:

$$\tilde{S} = \varepsilon \lfloor \frac{bl - \lfloor mb^{q-p+1} \rfloor}{b} \rfloor b^{p-s} \quad \text{lorsque } e(\tilde{S}) = p$$

$$\tilde{S} = \varepsilon (bl - \lfloor mb^{q-p+1} \rfloor) b^{p-s-1} \quad (\text{ou } 0) \quad \text{lorsque } e(\tilde{S}) < p$$

(On a exprimé dans ce dernier cas, la valeur de \tilde{S} , et non l'écriture après normalisation mettant en évidence l'exposant réel de ce nombre).

Proposition II. 8.

L'opération \oplus pour la troncature avec chiffre de garde, les nombres négatifs étant représentés avec signe et valeur absolue, est :

- 1) correcte
- 2) I - dirigée pour tout couple $(X, Y) \in T^2 \wedge XY > 0$, mais n'a pas de direction définie pour l'ensemble des couples $(X, Y) \in T^2 \wedge XY < 0$.
- 3) La symétrie est un automorphisme de T .

L'addition de deux nombres de même signe est correcte.

Pour mettre en évidence ce résultat pour l'addition de deux nombres de signes différents, écrivons, si $q \neq p$, les divisions arithmétiques suivantes :

$$m = b^{p-q-1} \quad m_1 + r_1 \wedge r_1 < b^{p-q-1}$$

$$bl - m_1 = bq_1 + r_2 \wedge r_2 < b$$

On a alors une expression commode de $E = S - \tilde{S}$ lorsque $q \neq p$:

($E = 0$, s'il ne se produit pas d'underflow, pour $q = p$)

$$\cdot \text{ Si } e(\tilde{S}) = p \quad E = \varepsilon(r_2 b^{p-q-1} r_1) b^{q-s} \quad \text{et} \quad -b^{p-1-s} < \varepsilon E < b^{p-s}$$

$$\cdot \text{ Si } e(\tilde{S}) < p \quad E = -\varepsilon r_1 b^{q-s} \quad \text{et} \quad |E| < b^{p-1-s}$$

Dans le cas où $e(\tilde{S}) < p$ remarquons :

$$q = p \vee q = p - 1 \implies E = 0 \quad (\text{ou } |E| < \xi)$$

$$q < p - 1 \implies e(\tilde{S}) = p - 1$$

La majoration de E est donc, dans tous les cas :

$$|E| < b^{P-s} \quad \text{où } P = e(\tilde{S}) \quad \text{si } \tilde{S} \neq 0$$

$$|E| < \xi \quad \text{si } \tilde{S} = 0$$

Ceci établit que \oplus est correcte, à la condition de s'assurer que la propriété est bien vérifiée dans le cas particulier $\tilde{S} = \varepsilon b^{p-1}$, ce qui est immédiat par $\varepsilon E > -b^{p-1-s}$.

L'opération \oplus pour la troncature avec chiffre de garde, les nombres négatifs étant représentés avec signe et valeur absolue, vérifie les hypothèses de la proposition II. 2.

La formule (II. 1) permet donc, pour cette arithmétique, le calcul de E .

(Il est toutefois nécessaire de remarquer que, si E est la valeur théorique de l'erreur, et \tilde{E} la valeur calculée par (II. 1) :

$$\tilde{E} = E \quad \text{si } E \in T$$

$$|\tilde{E} - E| < \xi \quad \text{si } |E| < \xi$$

($\tilde{E} = 0$ en général lorsque $|E| < \xi$, sauf dans le cas très particulier où les conditions nécessaires $e(\tilde{S}) = p + 1$, $q = m_2$, $p - q \leq s - 2$ peuvent entraîner $\tilde{E} = Y$)).

I. 3. 2. b. Les nombres négatifs sont représentés en notation de complément à 2

Désignons par $\overset{\sim}{\sigma}_*$ la représentation en complément à 2 de la mantisse $\overset{\sim}{\sigma}$ de $\tilde{S} = X \oplus Y$.

σ étant défini par : $S = X + Y = \sigma \cdot b^{e(S)}$, posons :

$$\sigma_* = \sigma \quad \text{si } \frac{1}{b} \leq \sigma < 1 \quad ; \quad \sigma_* = 2 - |\sigma| \quad \text{si } -1 < \sigma \leq -\frac{1}{b} .$$

Il vient :

$$\overset{\sim}{\sigma}_* = (\sigma_*)_s$$

la notation $(\quad)_s$ définissant la troncature à s chiffres après le bit de signe.

La propriété est, en effet, évidente si $e(S) = p \vee p + 1 \vee p - 1$, et d'autre part $e(S) < p - 1$ ne peut avoir lieu que pour $q = p \vee p - 1$, et dans ce cas $\overset{\sim}{\sigma} = \sigma$ si un underflow n'apparaît pas.

Propriété.

L'opération \oplus pour la troncature avec chiffre de garde, les nombres négatifs étant représentés en notation de complément à 2, est :

- 1) correcte
- 2) G - dirigée $\forall (X, Y) \in T^2 \wedge X + Y \notin] - \xi, 0 [$.

Si l'on complète la définition de \oplus par :

$p - q > s \implies X \oplus Y = X$ quels que soient les signes de X et de Y, afin que (C) ait lieu,

E est élément de T si la formation de \hat{S} n'a pas donné lieu à un underflow.

Le calcul de E se fera par la formule (II . 1) si $b = 2$, par la formule (II . 2) sinon.

I. 3. 3. L'arrondi sur chiffre de garde

Définition :

La définition de $X \oplus Y$ s'effectue comme en I. 3. 2., mais en 2), si la normalisation en fin d'opération n'a pas nécessité un décalage à gauche, le $(s + 1)^{i\text{ème}}$ chiffre de la mantisse du résultat est non pas tronqué, mais "arrondi" à l'aide de la règle suivante :

a) Si la valeur du $(s + 1)^{i\text{ème}}$ chiffre de mantisse est inférieure à $b/2$, le résultat est tronqué à s chiffres.

b) Si la valeur de ce $(s + 1)^{i\text{ème}}$ chiffre est supérieure à $b/2$, on ajoute une unité au $s^{i\text{ème}}$ chiffre de mantisse, on renormalise éventuellement et on tronque la mantisse à s chiffres.

c) Si la valeur du $(s + 1)^{i\text{ème}}$ chiffre de mantisse est égal à $b/2$, selon les calculateurs :

c - 1 . Il est opéré comme indiqué en b)

c - 2 . Il est défini un processus d'addition ou de non-addition d'une unité au $s^{i\text{ème}}$ chiffre de mantisse, évitant le biais introduit par c - 1 .

La définition donnée de l'arrondi sur chiffre de garde s'appliquera à une représentation des nombres par signe et valeur absolue.

(Signalons que certains calculateurs utilisant la notation de complément à 2 pour représenter la mantisse des nombres négatifs traduisent en fin d'opération mais avant arrondi, le résultat sous forme valeur absolue et signe.)

Le mode $c - 1$ d'arrondi est supposé choisi ci-dessous.

On peut énoncer

Propriété.

L'opération \oplus pour l'arrondi sur chiffre de garde, est :

- 1) correcte
- 2) optimale $\forall (X, Y) \in T^2 \wedge XY > 0$ si b est pair (voir restriction R_1)
- 3) vérifie la propriété (C) (voir restriction R_1 et R_2)
- 4) pour $b = 2$ et $b = 3$, l'erreur E vérifie, si $E \in T$, la formule (II. 1)

Les propriétés 1) 2) 3) énoncées viennent immédiatement.

Il convient de noter, pour ces résultats, que :

1. L'addition arithmétique pour l'arrondi sur chiffre de garde n'est pas optimale si b est impair.

On remarquera [15], en effet, que lorsque b est impair, un seul chiffre ne suffit pas pour contrôler l'arrondi, sauf si la notation "alternée" est utilisée, auquel cas l'arrondi se ramène simplement à la troncature.

Nous excluons cette possibilité dans ce paragraphe.

2. Pour une base puissance de 2, lorsque les nombres négatifs sont représentés en notation de complément à 2, les propriétés 2) et 3) ne sont plus vérifiées, si précaution n'est pas prise dans les règles définissant représentation et normalisation des mantisses, dans les cas particuliers suivants :

$$\mathcal{R}_1 : p - q > s \text{ (ou =)} ; X = - (b^s - 1) b^{p-s} ; Y < 0 \text{ pour 2) et 3)}$$

$$\mathcal{R}_2 : p - q > s ; X = b^{p-1} ; Y < 0 \text{ pour 3)}$$

qui peuvent conduire à $\tilde{S} = (X)'$.

La formule (II. 1) calculerait, dans ces cas particuliers, E lorsque $E \in T$, la partie principale de E sinon.

Enfin, pour la propriété 4), on utilisera, d'une part la proposition II. 3, et d'autre part, on remarquera que pour $b = 3$, l'addition arithmétique pour l'arrondi sur chiffre de garde, de deux nombres de même exposant $q = p$ implique :

$|E| = \varphi \cdot b^{p-s}$ où $\varphi \in \{0, 1\}$, et par suite $\tilde{S} - X = Y - E$ est toujours élément de T.

I. 3. 4. Opérations \oplus optimales

Rappelons seulement qu'il n'est pas nécessaire d'effectuer l'addition sur un accumulateur à double longueur, pour réaliser une opération \oplus optimale, lorsque b est pair.

Lorsque les nombres sont représentés avec signe et valeur absolue D. E. KNUTH [15] souligne en effet, que pour un schéma optimal de l'addition de X et Y, il suffit, lors du cadrage de Y, de tronquer la mantisse de ce nombre de la façon suivante :

Si X et Y ont le même signe, Y sera "approché" sous la forme $\varepsilon' \lfloor m b^{q-p+2} \rfloor b^{p-s-2}$, sinon sous la forme $\varepsilon' \lceil m b^{q-p+2} \rceil b^{p-s-2}$ (L'arrondi sur le $(s+1)$ ^{ième} chiffre de la mantisse du résultat est alors effectué comme § précédent).

Pour une base puissance de 2, lorsque les nombres négatifs sont représentés en notation de complément à 2, il suffira de tronquer, lors du cadrage de Y, la mantisse de ce nombre avec deux chiffres de garde. Toutefois, \oplus ne sera pas optimale dans le cas particulier \mathcal{R}_1 , lorsque la configuration de mantisse 1 0 0 est transformée en 1 (b - 1) 0—0, avec augmentation d'une unité pour l'exposant correspondant.

I. 3. 5. Une remarque pratique

Nous signalerons le travail de M. A. MALCOLM [25] qui a établi deux algorithmes permettant de déterminer automatiquement, sur les calculateurs disposant d'une arithmétique à virgule flottante correcte, les caractéristiques suivantes de cette arithmétique :

- base de la représentation
- nombre de chiffres de mantisse du système de nombres flottants
- propriété satisfaite par l'opération d'addition :
 - \oplus correcte ou \oplus optimale.

Ces algorithmes sont présentés sous forme de procédures FORTRAN.

Dans le cas de calculateurs disposant de registres flottants d'une longueur supérieure à la longueur de la représentation du mot-machine-et effectuant dans ces registres les calculs avec une précision supérieure à la précision avec laquelle les résultats peuvent être représentés en mémoire-il conviendra de se reporter aux modifications des procédures précédentes établies par W.MORVEN GENTLEMAN et S.B.MAROVICH [10] .

Nous mentionnerons enfin la récente étude de S.LINNAINMAA [21]: Cet auteur postérieurement à [33], mais par des démonstrations différentes, établit des conditions nécessaires et suffisantes de validité des formules (II.1) et (II.2). L'arithmétique "paritaire" est aussi étudiée et des résultats expérimentaux sont donnés par l'auteur.

II - L'ERREUR D'ARRONDI DANS LA MULTIPLICATION

II . 1. EXPRESSION, PAR UN PRODUIT DE COMPOSITION SUR T , DE L'ERREUR D'ARRONDI

La valeur sur \mathbb{R} du produit de deux nombres à virgule flottante étant un mot de longueur $2s$, le calcul de l'erreur d'arrondi dans la multiplication reviendra à réaliser une quasi double-précision.

On peut énoncer :

Proposition II.9.

Si s est pair et si les opérations \otimes et \ominus sont correctes, l'erreur $E = X Y - (X \otimes Y)$, entachant la multiplication en arithmétique à virgule flottante de deux éléments X et Y de T , vérifie, si aucun underflow n'apparaît lors des opérations définies ci-dessous, la formule suivante :

$$E = - (X \otimes Y \ominus P_1 \ominus P_2 \ominus P_2' \ominus P_3) \quad (\text{II-6})$$

où

$$P_1 = X_1 \otimes Y_1 ; P_2 = X_1 \otimes Y_2 ; P_2' = X_2 \otimes Y_1 ; P_3 = X_2 \otimes Y_2$$

X_1 et X_2 (resp. Y_1 et Y_2) étant définis à partir de $X = \varepsilon 1 b^{p-s}$ (resp. $Y = \varepsilon' m b^{q-s}$) par :

$$X_1 = \varepsilon \lfloor 1 b^{-s/2} \rfloor b^{s/2} . b^{p-s}$$

$$X_2 = X - X_1 \quad (\text{II-7})$$

$$Y_1 = \varepsilon' \lfloor m b^{-s/2} \rfloor b^{s/2} . b^{q-s}$$

$$Y_2 = Y - Y_1$$

Tout d'abord, il vient immédiatement :

Si l'opération $\otimes : T \times T \rightarrow T$ est correcte, alors $\forall (X,Y) \in T^2$, l'erreur $E = X Y - (X \otimes Y)$ est élément de T dès que $|E| \notin]0, \xi[\wedge |X Y| < b^{m1}$.

En effet, si $\xi \leq |X Y| < b^{m1}$, $|E|$ est du type λb^{p+q-2s} où λ est entier et vérifie $0 \leq \lambda < b^s$.

Lorsque $X \otimes Y$ provoque un overflow, hypothèse que nous excluerons, E n'est pas en général élément de T .

(II-6) exprime E à l'aide d'opérations internes définies sur T .

Démonstration :

Posons :

$$P = X Y$$

$$\overset{\vee}{P} = X \otimes Y$$

$$E = P - \overset{\vee}{P}$$

Si la formation de $\overset{\vee}{P}$, P_1 , P_2 , P'_2 , P_3 et E ne provoque pas d'underflow, il suffit pour établir la proposition, \ominus étant correcte, de démontrer que les trois nombres $\overset{\vee}{P} - P_1$, $\overset{\vee}{P} - P_1 - P_2$, $\overset{\vee}{P} - P_1 - P_2 - P'_2$ sont éléments de T .

Soit $\pi = p + q - v$ ($v = 0$ ou 1) l'exposant de P .

Si $l_1 = \lfloor 1 b^{-\frac{s}{2}} \rfloor$, $m_1 = \lfloor m b^{-\frac{s}{2}} \rfloor$, $l_2 = l - l_1$, $m_2 = m - m_1$,

posons :

$$P_2^t = \varepsilon \varepsilon' \lfloor l_1 m_2 b^{-\frac{s}{2}+v} \rfloor b^{\pi-s} \quad \text{et} \quad R_2 = P_2 - P_2^t$$

$$P_2'^t = \varepsilon \varepsilon' \lfloor l_2 m_1 b^{-\frac{s}{2}+v} \rfloor b^{\pi-s} \quad R_2' = P_2' - P_2'^t$$

$$P_3^t = \varepsilon \varepsilon' \lfloor l_2 m_2 b^{-s+v} \rfloor b^{\pi-s} \quad R_3 = P_3 - P_3^t$$

En écrivant :

$$R_2 + R_2' + R_3 = \varepsilon \varepsilon' (\lambda b^{\pi-s} + R) \quad \text{où } \lambda \in \{0, 1, 2\} \text{ et } 0 \leq R < b^{\pi-s}$$

il vient :

$$\overset{\vee}{P} = P_1 + P_2^t + P_2'^t + P_3^t + \varepsilon \varepsilon' (\lambda + \theta) b^{\pi-s}$$

avec $\theta = 0$ ou 1 suivant que \otimes est I ou E -dirigée pour le couple (X, Y) .

De l'écriture .

$$\overset{\vee}{P} - P_1 = P_2^t + P_2'^t + P_3^t + \varepsilon \varepsilon' (\lambda + \theta) b^{\pi-s}$$

on tire immédiatement que $\overset{\vee}{P} - P_1$ est de la forme $\varepsilon \varepsilon' \mu b^{\pi-s}$, où μ est entier et vérifie :

$$0 \leq \mu \leq 2(b^{\frac{s}{2}} - 2) b^v + b^v - 1 + 3 < b^s \quad \text{dès que } s \geq 4.$$

$$\overset{\vee}{P} - P_1 \in T.$$

Puis :

$$\overset{\vee}{P} - P_1 - P_2 = P_2'^t + P_3^t + \varepsilon \varepsilon' (\lambda + \theta) b^{\pi-s} - R_2$$

. Si $v = 1$, $P_2'^t + P_3^t + \varepsilon \varepsilon' (\lambda + \theta) b^{\pi-s}$ est de la forme $\varepsilon \varepsilon' \mu_1 b^{\pi-s}$
où : μ_1 est entier et vérifie $0 \leq \mu_1 \leq b^{\frac{s}{2}+1} - b + 2 \leq b^{\frac{s}{2}+1}$

En tenant compte de l'écriture de R_2 :

$$\overset{\vee}{P} - P_1 - P_2 \in T.$$

. Si $v = 0$, alors $P_3^t = 0$, $\max |P_2'^t| = (b^{\frac{s}{2}} - 2) b^{\pi-s}$ cette valeur ne pouvant être obtenue que lorsque $l_2 = m_1 = b^{\frac{s}{2}} - 1$.

La condition $l_2 = m_1 = b^{\frac{s}{2}-1}$ impliquant $\lambda \neq 2$,
 $P_2^t + \varepsilon \varepsilon' (\lambda + \theta) b^{\pi-s}$ est de la forme $\varepsilon \varepsilon' \mu_2 b^{\pi-s}$
 où μ_2 est entier et vérifie $0 \leq \mu_2 \leq b^{\frac{s}{2}}$.

Par ce résultat et l'écriture de R_2 , il vient encore :
 $\overset{\vee}{P} - P_1 - P_2 \in T$.

Enfin :

$$\overset{\vee}{P} - P_1 - P_2 - P_2' = P_3^t + \varepsilon \varepsilon' (\lambda + \theta) b^{\pi-s} - R_2 - R_2'$$

implique que ce nombre est du type $N b^{p+q-\frac{3s}{2}}$

où N est entier et est composé de $\frac{s}{2} + 2$ chiffres au maximum.

Par suite

$$\overset{\vee}{P} - P_1 - P_2 - P_2' \in T \text{ dès que } s \geq 4.$$

Remarques

- 1) On ne change pas (II-6) en permutant P_2 et P_2' , sauf en cas d'apparition d'un underflow.
- 2) En raison des underflows, qui peuvent survenir lors de la formation de P_1, P_2, P_2', P_3 ou lors des soustractions définies en (II-6), on n'est assuré d'obtenir par (II-6) l'erreur E dès que $E \in T$, seulement si est réalisée la condition suffisante suivante :

$$p + q \geq m_2 + 2s - 1$$

impliquant $E \in T$.

En effet, si cette condition sur les exposants n'est pas vérifiée, on n'a pas les implications suivantes, où $\overset{\vee}{E}$ désigne la valeur calculée par (II-6) :

$$|E| < \xi \implies \overset{\vee}{E} = 0$$

et

$$E \in T \implies \overset{\vee}{E} = E$$

Dans le cas de produits trop "petits", il faudrait pour utiliser (II-6) procéder préalablement à un recadrage.

On peut aussi énoncer :

Proposition II.9'.

La proposition II.9 est encore vraie si l'on remplace X_1, X_2, Y_1, Y_2 par X'_1, X'_2, Y'_1, Y'_2 définis par

$$X'_1 = \varepsilon \left[1 b^{-\frac{s}{2}} \right] b^{\frac{s}{2}} \cdot b^{p-s}$$

$$X'_2 = X - X'_1 \quad (\text{II-7'})$$

$$Y'_1 = \varepsilon' \left[m b^{-\frac{s}{2}} \right] b^{\frac{s}{2}} \cdot b^{q-s}$$

$$Y'_2 = Y - Y'_1$$

La démonstration de cette proposition est analogue à la démonstration précédente.

Lorsque s est impair, un calcul exact du produit de 2 nombres flottants à l'aide d'une arithmétique à mot court, nécessite que l'on partage X (et Y) en utilisant un arrondi afin que les produits P_1, P_2, P'_2, P_3 soient tous éléments de T .

G.W. VELTKAMP [40] a étudié (II-6) dans de telles hypothèses.

T.J. DEKKER [5] a établi un algorithme différent du précédent pour exprimer par un couple de mots courts, la valeur exacte du produit $X Y$, avec les hypothèses suivantes :

s quelconque

⊕ optimale

⊗ correcte.

Le procédé de cet auteur, calculant directement poids fort et poids faible du produit, nécessite une multiplication de moins et une addition de plus que (II-6).

II.2. FORMULES PRATIQUES

Les nombres X_1, X_2, Y_1, Y_2 (resp. X'_1, X'_2, Y'_1, Y'_2) ont été définis à partir de X et Y par (II.7) ou (II.7') de façon théorique. Lorsqu'on cherche à exprimer ces nombres à l'aide d'opérations internes sur T , une réponse simple à ce problème est donnée dans le cas où $b = 2$ par :

Proposition II.10.

Si les hypothèses de la proposition II.9 sont vérifiées, et si, de plus, $b = 2$ et \oplus est I - dirigée pour tout couple $(U, V) \in T^2 \wedge U \vee V > 0$, alors l'égalité (II.6) est vérifiée pour x_1, x_2, y_1, y_2 définis à partir de X et Y par :

$$\begin{aligned} x_1 &= z_1 \oplus X \ominus z_1, & x_2 &= X - x_1 \\ y_1 &= z_2 \oplus Y \ominus z_2, & y_2 &= Y - y_1 \end{aligned} \quad (\text{II.8})$$

où

$$z_1 = B \otimes X, \quad z_2 = B \otimes Y \quad \text{et} \quad B = b^{\frac{s}{2}}.$$

Démonstration :

Montrons, en effet, que la formule définissant x_1 (respectivement y_1) implique $x_1 = X_1$, sauf dans certains cas particuliers :

Soit x_s le dernier chiffre de la mantisse de X .

. Si $1 \leq b^s - b^{\frac{s}{2}} \vee (1 > b^s - b^{\frac{s}{2}} \wedge x_s = 1)$, il vient :

$$z_1 \oplus X = z_1 + X_1 \quad \text{d'où} \quad x_1 = X_1.$$

. Si $1 > b^s - b^{\frac{s}{2}} \wedge x_s = 0$ (C_p), il vient :

$$x_1 = X_1 - (b - 1) b^{\frac{s}{2}} \cdot b^{p-s}$$

$$x_2 = X_2 + (b - 1) b^{\frac{s}{2}} \cdot b^{p-s}$$

Il suffit, pour établir la proposition II.10, de montrer que l'égalité (II.6) est encore valable lorsque X (ou Y ou X et Y) vérifie la condition particulière C_p :

On remarquera que lorsque l'une de ces éventualités a lieu, les produits $x_1 y_2$, $x_2 y_1$, $x_2 y_2$ sont, sauf en cas d'underflow, éléments de T , et en reprenant la démonstration de (II.6) pour ces 3 cas particuliers, on vérifiera bien la propriété dès que $s \geq 6$.

Les formules (II.8) ne conviennent qu'avec les hypothèses précédentes ; en effet, si $b = 2$ mais \oplus est E -dirigée pour des couples d'éléments de T de même signe, ou si $b \neq 2$, on n'est plus assuré dans tous les cas que $x_1 y_2$, $x_2 y_1$ et $x_2 y_2$ appartiennent à T . Un "découpage" du nombre flottant X est alors une fonction plus compliquée de T dans T .

Ainsi :

Proposition II.11.

Si s est pair, et si l'opération \oplus est définie par une règle de troncature avec chiffre de garde, les nombres négatifs étant représentés par leur signe et leur valeur absolue,

alors, $\forall X \in T \wedge X \neq \epsilon(b^{s-1} + b^{\frac{s}{2}-1}) b^{p-s}$, l'élément X_1 de T , déduit de X par (II.7) vérifie :

$$X_1 = Z \oplus X \ominus Z$$

$$\text{où } Z = b^{\frac{s}{2}} X \ominus X.$$

Démonstration :

On établit cette proposition en vérifiant que les hypothèses choisies entraînent l'égalité :

$$Z \oplus X = Z + X_1 \quad \forall X \in T \wedge X \neq \epsilon(b^{s-1} + b^{\frac{s}{2}-1}) b^{p-s}$$

En effet, si $X = \epsilon l b^{p-s}$, il vient :

$$* l > b^{s-1} + b^{\frac{s}{2}-1} \implies e(Z) = p + \frac{s}{2}, \text{ avec de plus, en désignant par } x_{\frac{s}{2}+1} \text{ le } \left(\frac{s}{2} + 1\right)^{\text{ième}} \text{ chiffre de la mantisse de } X :$$

$$Z = b^{\frac{s}{2}} X - X_1 - \varepsilon b^{p-\frac{s}{2}} \quad \text{si } x_{\frac{s}{2}+1} \neq 0$$

$$Z = b^{\frac{s}{2}} X - X_1 \quad \text{si } x_{\frac{s}{2}+1} = 0$$

Par suite, dans les 2 cas :

$$Z \oplus X = Z + X_1$$

$$\star 1 < b^{s-1} + b^{\frac{s}{2}-1} \implies e(Z) = p + \frac{s}{2} - 1, \text{ avec}$$

$$Z = b^{\frac{s}{2}} X - X_1 \quad \text{car } x_{\frac{s}{2}+1} = 0$$

$$\text{Il vient encore } Z \oplus X = Z + X_1$$

$$\star 1 = b^{s-1} + b^{\frac{s}{2}-1} \implies Z \oplus X \ominus Z = X$$

On conclura pour la proposition, l'opération \oplus étant correcte.

Remarquons que, par contre, si \oplus et \otimes sont correctes, et I - dirigées $\forall (U, V) \in T^2 \wedge UV > 0$ pour \oplus , et $\forall (U, V) \in T^2$ pour \otimes , la fonction :

$$\varphi(X) = (\mathcal{B} \otimes X) \oplus X \ominus (\mathcal{B} \otimes X) \quad \text{où } \mathcal{B} = b^{\frac{s}{2}-1}$$

n'implique pas $\varphi(X) = X_1$, $\forall X \in T \wedge X \neq \varepsilon(b^{s-1} + b^{\frac{s}{2}-1}) b^{p-s}$

Ainsi :

$$b = 10, s = 6, X = 0,100034 \implies \varphi(X) = 0,0991.$$

Remarque :

Les expressions permettant de calculer X_1 à l'aide d'opérations internes sur T étant souvent compliquées, il sera parfois préférable de calculer X_1 en utilisant des propriétés simples de représentation des nombres et de programmation. (Voir programme ci-dessous).

Pour une arithmétique satisfaisant les hypothèses de la proposition II.11, l'application successive, de la formule établie dans cette proposition et de (II.6), calcule exactement l'erreur élémentaire de multiplication sauf dans le cas de l'élevation au carré d'un nombre de mantisse $\varepsilon(b^{-4} + b^{-s/4-4})$ [$\pm 0,100100$ sur IBM-360].

II.3. ETUDE EXPERIMENTALE

A titre d'exemple, nous avons effectué le calcul des produits $X Y$ où $X = \frac{1}{10+N}$, $Y = \frac{1}{250+N}$, $0 \leq N \leq 999$.

Ces calculs ont été effectués sur une IBM 360-30 : Les hypothèses des propositions II.9 et II.11 sont vérifiées.

Trois programmes ont été comparés :

Etant donnés les deux éléments X et Y de T et leur produit sur T ,
 $P = X \otimes Y$:

les programmes 1 et 2 utilisaient respectivement, pour le calcul de $E = X Y - P$, les instructions FORTRAN suivantes :
($b=16, s=6, B=b^{s/2}$)

Programme 1. $B = 4096$
 $Z = B * X - X$
 $X1 = Z + X - Z$
 $X2 = X - X1$
 $Z = B * Y - Y$
 $Y1 = Z + Y - Z$
 $Y2 = Y - Y1$
 $E = - P + X1 * Y1 + X1 * Y2 + X2 * Y1 + X2 * Y2$

Programme 2. EQUIVALENCE (J,X1) , (K,Y1)
 $X1 = X$
 $Y1 = Y$
 $J = (J/4096) * 4096$
 $K = (K/4096) * 4096$
 $X2 = X - X1$
 $Y2 = Y - Y1$
 $E = - P + X1 * Y1 + X1 * Y2 + X2 * Y1 + X2 * Y2$

Le programme 3 effectuait en double précision le produit XY .

Les différents temps d'exécution ont été successivement :

Programme 1 1 ' 49
Programme 2 1 ' 50
Programme 3 1 ' 54

note : Le programme 2, ainsi écrit, s'applique à des nombres flottants X et Y de signes positifs. Soulignons en effet que la représentation des nombres négatifs, dans un système IBM 360, utilise :

- la complémentation à 2 pour les nombres entiers
- signe et valeur absolue pour la mantisse des nombres réels.

III - L'ERREUR D'ARRONDI DANS LA DIVISION

Si (X, Y) est un élément de $T \times T - \{0\}$, l'erreur $E = X/Y - X \oslash Y$ n'est pas, en général, élément de T. Cependant E est exprimable comme la somme d'une série d'éléments de T, calculables par opérations internes sur T.

III.1. THEOREME PRELIMINAIRE

Proposition II.12.

Si l'opération \oslash est correcte, alors, quel que soit le couple (X, Y) de $T \times T - \{0\}$ vérifiant $\left| \frac{X}{Y} \right| < b^{m_1}$, le nombre $C = E Y$ où $E = X/Y - (X \oslash Y)$, est élément de $T U] - \xi, \xi [$.

Démonstration :

Soient, en effet, $X = \epsilon l b^{p-s}$ et $Y = \epsilon' m b^{q-s}$ ($b^{s-1} \leq l, m < b^s$) deux éléments de T.

Si $X \oslash Y$ provoque un overflow, hypothèse que nous avons exclue, C n'est pas en général élément de T, et dans le cas d'underflow, la propriété est évidente.

Dans le cas général, il vient :

$$\epsilon \epsilon' \left(\left\lfloor \frac{l b^s}{m} \right\rfloor \vee \left\lceil \frac{l b^s}{m} \right\rceil \right) b^{p-q-s} \quad \text{si } l < m$$

$X \oslash Y =$

$$\epsilon \epsilon' \left(\left\lfloor \frac{l b^{s-1}}{m} \right\rfloor \vee \left\lceil \frac{l b^{s-1}}{m} \right\rceil \right) b^{p-q+1-s} \quad \text{si } l \geq m$$

En posant, dans chacun des cas, la division arithmétique :

$$l b^{s-\theta} = m d + r \wedge r < m \quad \theta = 0, 1$$

on a :

$$C = \varepsilon (r_{\vee} r - m) b^{p+\theta-2s}$$

d'où immédiatement la proposition.

III.2. CALCUL DE L'ERREUR

III.2.1. Calcul de C

Notons :

$$D = X/Y$$

$$\tilde{D} = X \oslash Y$$

On obtiendra C en calculant sans erreur le nombre $X - \tilde{D} Y$, c'est-à-dire :

- soit, si l'on travaille en simple précision mais que l'on dispose de la double précision pour la multiplication et l'addition, en réalisant les 2 opérations $\tilde{D} \cdot Y$ et $X - \tilde{D} Y$ en double précision - sous réserve que l'arithmétique définie pour la double précision soit correcte -
- soit si l'on calcule uniquement par une arithmétique à mot fixé, en utilisant, par exemple, sous réserve de sa validité dans l'arithmétique flottante employée, la formule suivante :

$$C = X \ominus P_1 \ominus P_2 \ominus P'_2 \ominus P_3 \quad (\text{II.9})$$

où

$$P_1 = \tilde{D}_1 \otimes Y_1 ; P_2 = \tilde{D}_1 \otimes Y_2 ; P'_2 = \tilde{D}_2 \otimes Y_1 ; P_3 = \tilde{D}_2 \otimes Y_2$$

$\tilde{D}_1, \tilde{D}_2, Y_1, Y_2$ étant définis à partir de \tilde{D} et Y par des formules analogues à (II.7).

Ainsi on établit :

Proposition II.13.

Si s est pair, si les opérations $\ominus, \otimes, \oslash$ sont correctes, et si \oslash est I - dirigée pour tout couple de $T \times T - \{0\}$, alors, (II.9) est vérifiée à la condition suffisante qu'aucun underflow n'apparaisse lors des opérations définies par cette formule.

Démonstration :

La démonstration est analogue à la démonstration de la proposition II.9.

Posant :

$$m = m_1 b^{\cancel{s/2}} + m_2 \wedge m_2 < b^{\cancel{s/2}}$$

$$d = d_1 b^{\cancel{s/2}} + d_2 \wedge d_2 < b^{\cancel{s/2}}$$

on a la relation en entiers :

$$1 b^{s-\theta} = m_1 d_1 b^s + (d_1 m_2 + d_2 m_1) b^{\cancel{s/2}} + d_2 m_2 + r \quad (\text{II.10})$$

traduction de :

$$X = P_1 + P_2 + P'_2 + P_3 + C$$

où tous les termes écrits sont du signe ϵ de X .

Ceci met en évidence :

$$X - P_1 = \epsilon (1 - m_1 d_1 b^\theta) b^{p-s} = \epsilon \mu_1 b^{p-s}$$

où μ_1 est entier et vérifie : $0 \leq \mu_1 < 1$, impliquant $X - P_1 \in T$.

(II.10) exprime que $d_2 m_2 + r$ est divisible par $b^{\cancel{s/2}}$, soit :

$$d_2 m_2 + r = \lambda b^{\cancel{s/2}} , \text{ et il vient : } \lambda \leq 2 b^{\cancel{s/2}} - 2 .$$

Par suite :

$$X - P_1 - P_2 = \epsilon (d_2 m_1 + \lambda) b^{p+\theta-\frac{3s}{2}} = \epsilon \mu_2 b^{p+\theta-\frac{3s}{2}}$$

où μ_2 est entier et vérifie $0 \leq \mu_2 \leq b^s - 1$, impliquant $X - P_1 - P_2 \in T$.

Enfin

$$X - P_1 - P_2 - P'_2 = \epsilon \lambda b^{p+\theta-\frac{3s}{2}} , \text{ donc } X - P_1 - P_2 - P'_2 \in T.$$

D'où la proposition , Θ étant correcte.

III.2.2. Algorithme d'obtention de E

Notons \ominus et \otimes les opérations de soustraction et de multiplication : en double précision par des opérations correctes

ou en pseudo double précision par des formules du type II.9 ,
c'est-à-dire vérifiant dans les 2 cas :

$$X \otimes Y = X Y \quad X \in T, Y \in T$$

$$X \ominus Y Z = X - Y Z \quad X \in T, Y \in T, Z \in T, X - Y Z \in T$$

Si l'opération \oslash est correcte, le quotient $D = X/Y$ s'obtient immédiatement comme la somme de la série $\{\tilde{D}_n\}$ définie, ainsi que la série auxiliaire $\{C_n\}$, par les formules suivantes :

$$\begin{aligned} \tilde{D}_1 &= X \oslash Y \\ C_1 &= X \ominus \tilde{D}_1 \otimes Y \end{aligned} \tag{II.11}$$

et pour $n > 1$:

$$\begin{aligned} \tilde{D}_n &= C_{n-1} \oslash Y \\ C_n &= C_{n-1} \ominus \tilde{D}_n \otimes Y \end{aligned}$$

III.2.3. Résultats numériques

Nous avons testé l'algorithme (II.11) sur une IBM 360-30 en simple précision, pour le calcul des trois premiers éléments, de la suite des nombres à virgule flottante, approchant $\frac{1}{X}$ ($X = 1, 2, \dots$)

L'écriture de chaque élément du triplet représentant $\frac{1}{X}$ a été réalisée en hexadécimal, afin d'éviter toute erreur de conversion d'écriture des nombres selon des bases différentes. On a ainsi obtenu exactement la représentation de $\frac{1}{X}$ en triple précision.

X=5	$\tilde{D}_1=40333333$	$\tilde{D}_2=3A333333$	$\tilde{D}_3=34333333$
X=6	$\tilde{D}_1=402AAAAA$	$\tilde{D}_2=3AAAAAAA$	$\tilde{D}_3=34AAAAAA$
X=7	$\tilde{D}_1=40249249$	$\tilde{D}_2=3A249249$	$\tilde{D}_3=34249249$

IV - CONCLUSION

Lorsque se pose la question de comparer les 2 méthodes suivantes d'augmentation de la précision d'un résultat :

- utilisation d'un mot long
 - ou expression du résultat d'opération élémentaire pour un couple de mots courts (partie principale et erreur) ,
- il est possible de faire plusieurs remarques :

1) Il est préférable, lorsque cela est possible, d'exprimer le résultat R d'une opération élémentaire sous la forme d'un mot long, plutôt que par 2 mots courts, dans les cas de multiplications, et dans les cas d'additions où la propriété (C) n'intervient pas. Cela résulte, en effet, de la forme de l'erreur qui vérifie dans les cas ci-dessus :

$$E = \lambda b^{e(R) - 2s} \quad \text{où } \lambda \text{ est entier et vérifie } 0 \leq |\lambda| < b^s .$$

Par contre, lorsque la condition $e(X) - e(Y) > s$ est réalisée lors de l'addition des 2 nombres X et Y , E n'est pas exprimable exactement dans l'écriture par mot long, et c'est l'écriture par couple de mots courts qui permet de garder toute la précision.

Cette forme est également la meilleure lorsqu'un quotient est nécessaire avec de nombreux chiffres de mantisse, l'algorithme II. 11 permettant d'approcher un quotient avec une précision arbitraire.

2) Globalement -et non plus au niveau d'opération élémentaire- un résultat calculé en double précision est meilleur qu'un résultat, calculé en simple précision et corrigé par une correction.

Il ne faut, en effet, pas oublier l'importance prépondérante, pour l'obtention d'une précision suffisante, des 2 facteurs suivants :

- * nombre suffisant des chiffres de mantisse
- * utilisation d'une arithmétique correcte.

Par contre, l'amélioration apportée par une arithmétique, non seulement, correcte, mais optimale, n'est pas décisive.

3) Cependant, l'intérêt des méthodes exposées dans ce travail, est, outre que la double précision n'est pas accessible sur tous les calculateurs, leur possibilité d'itération.

Ces méthodes sont applicables à toute arithmétique, et leur utilisation pour une arithmétique à mot long, assurera une très grande précision dans les cas où cela est nécessaire.

4) Nous soulignerons, pour terminer, que le temps nécessité par les calculs dans la méthode d'expression de l'erreur comme poids faible d'un couple de mots courts, est d'un ordre de grandeur analogue au temps nécessité par la double précision.

V - APPENDICE. SOUS-PROGRAMMES FORTRAN

Nous joignons à ce chapitre, les sous-programmes FORTRAN, relatifs à l'étude précédente, que nous avons utilisé sur le système I.B.M. 360.

Les sous-programmes :

FUNCTION ES (X, Y, S)

" EP (X, Y, P)

" CD (X, Y, D)

calculent respectivement, à partir des données X et Y et des valeurs $S = X \oplus Y$, $P = X \otimes Y$, $D = X \oslash Y$, les erreurs :

$$ES = X + Y - S$$

$$EP = XY - P$$

et le terme correctif $CD = (X/Y - D) Y$.

SIMPLE PRECISION

```
FUNCTION ES(X,Y,S)
  IF(ABS(X)-ABS(Y)) 1,1,2
1 ES=-S+Y+X
  GOTO 3
2 ES=-S+X+Y
3 CONTINUE
  RETURN
  END
```

```
FUNCTION EP(X,Y,P)
EQUIVALENCE (J,X1),(K,Y1)
L=4096
U=ABS(X)
V=ABS(Y)
Q=ABS(P)
X1=U
Y1=V
J=(J/L)*L
K=(K/L)*L
X2=U-X1
Y2=V-Y1
EP=-Q+X1*Y1+X1*Y2+X2*Y1+X2*Y2
IF(P.LT.0.0) EP=-EP
RETURN
END
```

```
FUNCTION CD(X,Y,D)
EQUIVALENCE (J,Y1),(K,D1)
L=4096
U=ABS(X)
V=ABS(Y)
Q=ABS(D)
Y1=V
D1=Q
J=(J/L)*L
K=(K/L)*L
Y2=V-Y1
D2=Q-D1
CD=U-D1*Y1-D1*Y2-D2*Y1-D2*Y2
IF(X.LT.0.0) CD=-CD
RETURN
END
```

DOUBLE PRECISION

```
REAL FUNCTION ES*8(X,Y,S)
DOUBLE PRECISION X,Y,S
IF(DABS(X)-DABS(Y)) 1,1,2
1 ES=-S+Y+X
  GOTO 3
2 ES=-S+X+Y
3 CONTINUE
  RETURN
  END
```

```
REAL FUNCTION EP*8(X,Y,P)
REAL*8 X,Y,P
REAL*8 Z,X1,X2,Y1,Y2
B=268435456.0D0
Z=B*X-X
X1=Z+X-Z
X2=X-X1
Z=B*Y-Y
Y1=Z+Y-Z
Y2=Y-Y1
EP=-P+X1*Y1+X1*Y2+X2*Y1+X2*Y2
  RETURN
  END
```

```
REAL FUNCTION CD*8(X,Y,D)
REAL*8 X,Y,D
REAL*8 Z,D1,D2,Y1,Y2
B=268435456.0D0
Z=B*Y-Y
Y1=Z+Y-Z
Y2=Y-Y1
Z=B*D-D
D1=Z+D-Z
D2=D-D1
CD=X-Y1*D1-Y1*D2-Y2*D1-D2*Y2
  RETURN
  END
```

CHAPITRE III

APPLICATIONS

Introduction

Ce chapitre est relatif aux applications du calcul de l'erreur élémentaire comme poids faible d'un couple de mots en arithmétique à virgule flottante.

T désignera encore dans tout le chapitre le système T_b^S de nombres flottants, et les opérations élémentaires définies sur cet ensemble seront toujours notées \oplus , \ominus , \otimes , \oslash .

Le paragraphe I est consacré à un algorithme de correction d'une somme algébrique -ou d'un produit scalaire- .

Les paragraphes II, III et IV examinent l'utilisation des techniques décrites au chapitre précédent dans les trois domaines suivants :

- L'analyse à posteriori des erreurs d'arrondi
- L'utilisation de méthodes de raffinement itératif
- La résolution numérique de problèmes mal conditionnés

Ces trois aspects différents sont traités sur des exemples.

I - UN ALGORITHME DE CORRECTION D'UNE SOMME ALGEBRIQUE

I . 1 . LE PROBLEME ETUDIE

Etant donné un système T de nombres à virgule flottante, muni de la loi d'addition \oplus , soient N nombres $X_i \in T$, ($1 \leq i \leq N$), et $S = \sum_{i=1}^N X_i$ la somme de ces nombres sur \mathbb{R} .

Nous nous proposons d'obtenir, par des opérations internes sur T , un élément de T dont la distance à S soit la plus petite possible.

Posons : $S^1 = \oplus \sum_{i=1}^N X_i$

Si e_i désigne l'erreur d'arrondi lors de la i ème addition dans le calcul de S^1 , on a immédiatement :

$$S = S^1 + \sum_{i=1}^{N-1} e_i .$$

Lorsque les erreurs e_i sont éléments de T , un procédé de correction de S^1 sera le suivant : on ajoutera, à cette valeur trouvée S^1 , la somme des e_i calculée sur T ; puis, cette dernière somme étant entachée d'erreur, on pourra effectuer une deuxième correction, et ainsi de suite.

Définition de l'algorithme α

Supposons que la fonction définie sur T^2 par :

$$(X,Y) \longrightarrow e = X + Y - (X \oplus Y)$$

soit une application de T^2 dans $T \cup]-\xi, \xi[$.

Définissons les suites $\{S^n\}$, $\{e_i^n\}$ ($1 \leq i \leq N-1$), $\{E_i^n\}$ ($1 \leq i \leq N-1$) et $\{C^n\}$ de la façon suivante :

Les premiers termes de chaque suite sont définis par :

$$S^1 = \oplus \sum_{i=1}^N X_i$$

$$e_i^1 = \left(\oplus \sum_{j=1}^i X_j \right) + X_{i+1} - \left(\oplus \sum_{j=1}^{i+1} X_j \right) \quad 1 \leq i \leq N-1$$

e_i^1 est "l'erreur" commise lors de la i ème addition dans le calcul de S^1 .

La suite $\{E_i^n\}$ se déduit de $\{e_i^n\}$ par :

$$E_i^n = \begin{cases} e_i^n & \text{si } e_i^n \in T \\ 0 & \text{sinon} \end{cases}$$

Enfin :

$$C^1 = \bigoplus_{i=1}^{N-1} E_i^1 .$$

Puis :

$$S^n = S^{n-1} \bigoplus C^{n-1}$$

Soient e_i^n ($1 \leq i \leq N-2$) définies par :

$$e_i^n = \left(\bigoplus_{j=1}^i E_j^{n-1} \right) + E_{i+1}^{n-1} - \left(\bigoplus_{j=1}^{i+1} E_j^{n-1} \right)$$

et soit

$$e_{N-1}^n = S^{n-1} + C^{n-1} - S^n$$

(e_i^n est "l'erreur" commise lors de la i ème addition dans le calcul de C^{n-1}).

Enfin :

$$C^n = \bigoplus_{i=1}^{N-1} E_i^n$$

Des définitions, il vient :

$$S = S^1 + \sum_{i=1}^{N-1} e_i^1$$

$$S^{n-1} + \sum_{i=1}^{N-1} E_i^{n-1} = S^n + \sum_{i=1}^{N-1} e_i^n \quad \text{pour } n > 1.$$

Donc :

$$S = S^n + \sum_{i=1}^{N-1} e_i^n + \sum_{j=1}^{n-1} \sum_{i=1}^{N-1} (e_i^j - E_i^j)$$

ou plus simplement :

$$\text{si } I_n = \{i ; e_i^n \neq 0\}$$

et $\Lambda_n = \{(i,j) ; j < n \wedge 0 < |e_i^j| < \xi\}$

vérifiant : $\text{Card } I_n + \text{Card } \Lambda_n \leq N-1$

on a :

$$S = S^n + \sum_{i \in I_n} e_i^n + \sum_{(i,j) \in \Lambda_n} e_i^j \quad (\text{III.1})$$

I . 2 . CONVERGENCE DE L'ALGORITHME

Avant d'aborder l'étude du comportement de l'algorithme \mathcal{A} , examinons le point suivant :

Lorsque n augmente, les erreurs e_i^n vont tendre à diminuer en valeur absolue. Sera-t-il alors possible de conserver, dans les éléments E_i^n toute l'information sur les termes d'erreur, ou au contraire, l'existence d'inégalités $0 < |e_i^n| < \xi$ entraînera-t-elle, par la mise à 0 de E_i^n , une perte de précision ?

Il vient :

Si l'opération \oplus est correcte et si les N éléments X_i de T vérifient l'inégalité suivante :

$$e(X_i) \geq m_2 + s-1 \quad 1 \leq i \leq N ,$$

alors, les erreurs e_i^n sont, quel que soit n , éléments de T .

On remarquera pour ce résultat, les deux propriétés suivantes :

1) Si X et Y sont deux éléments de T tels que $e(X) \geq e(Y) \wedge X + Y \in R_T$

$$E = X + Y - (X \oplus Y) = \lambda b^{e(Y)-s} \quad \text{avec } \lambda \in \mathbb{Z}$$

entraîne $E \in T$ à la condition suffisante $e(Y) \geq m_2 + s-1$.

(Remarque du § I.1 chapitre II)

2) La somme $\tilde{X} = X \oplus Y$ de deux éléments de T tels que $e(X) \geq e(Y)$ vérifie l'un des deux points a ou b suivants :

- a) $e(\overset{\vee}{S}) \geq e(Y)$
- b) Si $e(\overset{\vee}{S}) < e(Y)$ et si t est l'entier positif défini par $t = e(Y) - e(\overset{\vee}{S})$, alors les t derniers chiffres de la mantisse de $\overset{\vee}{S}$ sont nuls. L'erreur entachant la somme de $\overset{\vee}{S}$ et d'un élément de T d'exposant supérieur ou égal à $e(\overset{\vee}{S}) + t$ satisfait encore à $E = \lambda b^{e(Y)-s_\lambda} \in \mathbf{Z}$.

- La plus petite erreur e_i^n , non nulle, possible est

$$\min_b \{e(X_i) - s\}$$

Il s'en suit que, dans le cas où les éléments X_i sont d'exposants trop proches de m_2 , il conviendra de les multiplier par un facteur de cadrage, afin de ne perdre aucune précision.

Nous ne considérerons donc pas dans la suite l'éventualité d'underflows, mais les modifications de démonstrations dues à ce phénomène sont dans [32] ; $E_i^n = e_i^n$, $\bigwedge_n = \emptyset$.

I. 2. 1. Deux théorèmes relatifs à une opération d'addition correcte

On peut énoncer :

Proposition III . 1 .

Si l'opération \oplus est :

correcte

I - dirigée pour tout couple $(X, Y) \in T^2 \wedge X Y > 0$,

Alors, si $2N \leq b^{s-1}$, il existe un entier $k > 0$ tel que :

$$|S - S^k| < 2 b^{P^k - s} \quad (\text{III.2})$$

où $P^k = e(S^k)$.

Démonstration :

La démonstration ci-dessous suppose exclue toute possibilité d'overflow.

Examinons le passage du pas $n - 1$ au pas n .

Soit Q^n le plus petit entier tel que :

$$\sum_{i=1}^{N-1} |E_i^{n-1}| < b^{Q^n} \quad (\text{III.3})$$

et montrons que nous avons deux possibilités :

- * soit $Q^{n+1} < Q^n$
- * soit l'algorithme est terminé.

Remarquons que, si $Z_i \in T$, ($1 \leq i \leq P$), l'opération \oplus étant I - dirigée pour tout couple $(X, Y) \in T^2 \wedge X Y > 0$, et de plus correcte, on a :

$$| \oplus_{i=1}^1 Z_i | \leq \sum_{i=1}^P |Z_i| \quad (\text{III.4}) \quad \text{pour } \forall 1 \leq P$$

(en effet, \oplus étant correcte, pour tout couple $(X, Y) \in T^2 \wedge X Y < 0$ on a :

$$|X \oplus Y| \leq \text{Max} (|X|, |Y|) \leq |X| + |Y|).$$

\oplus étant correcte et aucun overflow n'étant possible, (III.3) et (III.4) entraînent :

$$|e_i^n| < b^{Q^n - s} \quad 1 \leq i \leq N-2$$

Notant $e(S^n) = P^n$:

$$\sum_{i \in I_n} |e_i^n| < (N-2) b^{Q^n - s} + b^{P^n - s} < \frac{1}{2} b^{Q^n - 1} + b^{P^n - s}$$

* Si $\frac{1}{2} b^{Q^n - 1} \leq b^{P^n - s}$, on a $|S - S^n| < 2 b^{P^n - s}$, l'algorithme est terminé et seul le dernier bit de la mantisse de S^n peut être entaché d'erreur.

* Sinon, il vient $Q^{n+1} < Q^n$, et si l'inégalité (III.2) n'est pas vérifiée au pas n , on forme le pas suivant, en remarquant que la suite des Q^n est finie.

D'où la proposition.

(Si $S^k = 0$, on fera $P^k = -\infty$ dans l'inégalité (III.2) qui sera prise au sens large).

Notation. Lorsque (III.2) est vérifiée, l'algorithme \mathcal{A} sera dit convergent.

Une autre proposition est donnée par :

Proposition III.2.

Si l'opération \oplus est :

correcte

isotone

une condition suffisante de convergence de l'algorithme \mathcal{A} est que soit réalisée $2N \ll b^{s-2}$.

Démonstration :

En effet, (III.4) n'est plus valable, mais sera remplacée par l'inégalité

$$\left| \oplus \sum_{i=1}^1 Z_i \right| \ll \oplus \sum_{i=1}^1 |Z_i| \quad (\text{III.4}')$$

lorsque la symétrie est un automorphisme de T ,

ou par :

$$\left| \oplus \sum_{i=1}^1 Z_i \right| \ll \text{Max} \left\{ \oplus \sum_{i=1}^1 |Z_i|, \left| \oplus \sum_{i=1}^1 -|Z_i| \right| \right\} \quad (\text{III.4}''')$$

sinon.

On a aussi la double inégalité moins large :

$$\oplus \sum_{i \in J} Z_i \ll \oplus \sum_{i=1}^1 Z_i \ll \oplus \sum_{i \in I} Z_i$$

où $I = \{i ; Z_i > 0\}$, $J = \{i ; Z_i < 0\}$

valable dès que \oplus est isotone.

On utilisera ensuite le lemme suivant :

Lemme III.1.

Si l'opération \oplus est correcte, et si $2l \ll b^{s-1}$, la somme en arithmétique à virgule flottante de l nombres de même signe, $Z_i \in T$ ($1 \leq i \leq l$), vérifie :

$$e \left(\oplus \sum_{i=1}^l Z_i \right) \ll e \left(\sum_{i=1}^l Z_i \right) + 1 \quad (\text{III.5})$$

La démonstration du lemme est immédiate :

$$\text{Soient } S = \sum_{i=1}^1 Z_i, \quad \tilde{S} = \oplus \sum_{i=1}^1 Z_i, \quad p_i = e(\oplus \sum_{j=1}^{i+1} Z_j).$$

Il vient :

$$|S - \tilde{S}| \leq \sum_{i=1}^{l-1} b^{p_i - s} \quad \text{et} \quad b^{p_i - s} \leq b^{1-s} \left| \oplus \sum_{j=1}^{i+1} Z_j \right| \quad \text{entraîne :}$$

$$|S - \tilde{S}| \leq b^{1-s} | \tilde{S} | \quad (\text{III.6})$$

d'où le résultat, que l'on déduit aussi des formules (25.5) à (25.10) de [42].

I . 2 . 2 . Un théorème de convergence pour une addition définie par une troncature sans chiffre de garde.

Afin de ne pas compliquer inutilement l'écriture, nous démontrerons la proposition ci-dessous avec l'hypothèse :

$$e(X_i) \geq m_2 + s - 1 \quad 1 \leq i \leq N$$

pour encore exclure toute possibilité d'underflow.

Si l'hypothèse précédente n'est pas vérifiée, il convient de modifier (III.1) en remarquant que si, par exemple, les nombres flottants sont représentés par leur signe et leur valeur absolue et les erreurs d'arrondi élémentaires calculées par la formule (II.5), la suite $\{E_i^n\}$ se déduit de $\{e_i^n\}$ par :

$$E_i^n = \begin{cases} e_i^n & \text{si } e_i^n \in T - \{0\} \\ 0 \text{ ou } \pm \xi & \text{sinon} \end{cases}$$

Proposition III.3.

Si la règle définissant l'addition est une troncature sans chiffre de garde, une condition suffisante de convergence de l'algorithme \mathcal{A} est que soit réalisée $2N \leq b^{s-1}$.

Modifions, avec les hypothèses présentes, la démonstration de la proposition III.1.

Soit encore $Q^n = e(\sum_{i=1}^{N-1} |E_i^{n-1}|)$, et examinons le passage du pas $n-1$ au pas n .

L'inégalité (III.4) reste valable ; en effet :

\oplus est I - dirigée pour tout couple $(X, Y) \in T_{\wedge}^2$ $X Y > 0$
 et $|X \oplus Y| \leq \max (|X| , |Y|) \quad \forall (X, Y) \in T_{\wedge}^2$ $X Y < 0$.

Le module de l'erreur e_i^n est majoré à l'aide de l'inégalité (II.4) :

$$|e_i^n| < b^{p_i^n - s}$$

où

$$p_i^n = \max \{ e (\oplus_{j=1}^{i+1} E_j^{n-1}) , e (\oplus_{j=1}^i E_j^{n-1}) , e (E_{i+1}^{n-1}) \}$$

En utilisant (III.4) , il vient immédiatement :

$$p_i^n \leq Q^n$$

et

$$|e_i^n| < b^{Q^n - s} \quad 1 \leq i \leq N-2$$

Posant $\pi^n = \max \{ e(S^n) , e(S^{n-1}) , e(C^{n-1}) \}$,

$$|e_{N-1}^n| < b^{\pi^n - s}$$

Distinguons 3 cas :

* Si $\pi^n = e(S^n)$, la démonstration de la proposition s'effectue comme celle de la proposition III.1.

* Si $\pi^n = e(C^{n-1})$, il vient immédiatement $\pi^n \leq Q^n$ et la proposition se déduit par : $\sum_{i=1}^{N-1} |e_i^n| < N b^{Q^n - s} \implies Q^{n+1} < Q^n$.

* Si $\pi^n = e(S^{n-1})$, l'opération $S^n = S^{n-1} \oplus C^{n-1}$ est une soustraction arithmétique ; rappelons que $e(S^{n-1}) - e(C^{n-1}) \geq 2$ entraîne $e(S^n) \geq e(S^{n-1}) - 1$.

Posons : $\pi^n = Q^n + L$

Si $L \leq s-2$, on a :

$$\sum_{i=1}^{N-1} |e_i^n| \leq (N-2) b^{Q^n - s} + b^{Q^n + L - s} < (\frac{1}{2} + \frac{1}{b}) b^{Q^n - 1} \text{ donc } Q^{n+1} < Q^n$$

Si $L = s - 1$, l'algorithme est terminé, soit au pas n , soit au pas $n+1$.

I . 3 . REMARQUES

Remarque 1. Précision apportée par la première correction :

$$\text{Si } P = \max_{2 \leq i \leq N} \left\{ e \left(\oplus \sum_{j=1}^i X_j \right) \right\}$$

on a :

$$|S - S^2| < (N - 2) b^{P+\gamma-2s} + b^{P^2-s}$$

où γ est le plus petit entier tel que $N \leq b^\gamma$.

Remarque 2. Correction d'un produit scalaire :

Si $X = \{X_i\}$ et $Y = \{Y_i\}$ sont deux vecteurs de T^N , lorsque l'on désigne par e_i les erreurs de multiplication et d'addition suivantes :

$$e_i = X_i Y_i - (X_i \otimes Y_i) \quad i = 1, \dots, N$$

$$e_{N+i} = \left\{ \left(\oplus \sum_{j=1}^i X_j \otimes Y_j \right) + X_{i+1} \otimes Y_{i+1} \right\} - \oplus \sum_{j=1}^{i+1} X_j \otimes Y_j$$

$i = 1, \dots, N$

il vient immédiatement :

$$X Y = X \otimes Y + \sum_{i=1}^{2N-1} e_i$$

La correction d'un produit scalaire se ramène, dès que les e_i sont calculés, à la correction d'une somme.

I . 4 . TEST D'ARRET DE L'ALGORITHME -a

Il est particulièrement simple lorsque l'opération \oplus est G - dirigée (resp. D - dirigée), puisque l'égalité

$$S^n = S^{n-1}$$

permet d'affirmer la convergence .

S'il n'en est pas ainsi, on adjoindra au calcul de $C^{n-1} = \bigoplus_{i=1}^{N-1} E_i^{n-1}$,

le calcul de :

$$\text{TEST} = \bigoplus_{i=1}^{N-1} |E_i^{n-1}|$$

et par le lemme III.1, la convergence sera assurée par la condition suffisante :

$$\frac{3}{2} \text{TEST} < \theta \cdot b^{P^n - s} \quad (\text{III.7})$$

où $\theta = 1$ ou 2 suivant la précision voulue.

On peut aussi adjoindre au calcul de C^{n-1} , le calcul suivant :

$$\text{Posant } F_i = \bigoplus_{j=1}^{i+1} E_j^{n-1}$$

on calcule

$$\text{TEST} = \bigoplus_{i=1}^{N-2} |F_i|$$

et on utilise la majoration suivante obtenue par (III.6)

$$|S - S^n| < b^{1-s} (1 + b^{1-s}(N-2)) \text{TEST} + |E_{N-1}^n|$$

I . 5 ETUDE NUMERIQUE

Nous avons considéré 2 sommes arithmétiques :

$$\sigma_1 = \sum_{i=1}^N \frac{1}{i} \quad \text{et} \quad \sigma_2 = \sum_{i=1}^N \frac{1}{i^2}$$

et 2 sommes algébriques :

$$\sigma_3 = \sum_{i=1}^N \frac{(-1)^{i-1}}{i} \quad \text{et} \quad \sigma_4 = \sum_{i=1}^N \frac{(-1)^{i-1}}{i^2}$$

Nous avons utilisé l'arithmétique à virgule flottante définie sur une IBM 360-30 : $b = 16$, $s = 6$; l'opération \bigoplus est correcte, I - dirigée pour tout couple $(X, Y) \in T^2$ \wedge $XY > 0$.

Il vient immédiatement, par la remarque 1 du § I.3, que pour $n < b^3$, une seule correction sera nécessaire pour rendre significatifs tous les chiffres du résultat.

Dans chacun des 4 cas envisagés, soient :

S^1 la valeur de la somme calculée de façon habituelle, en arithmétique à virgule flottante

S^2 la valeur de cette somme corrigée par une seule correction

S D la valeur de la somme effectuée en double précision.

Les calculs ont été faits pour N variant de 100 à 500, par pas de 100. La lecture des résultats montre que tous les chiffres de la somme corrigée S^2 sont significatifs, une éventuelle différence entre les derniers chiffres de S^2 et S D provenant de la conversion des résultats en décimal, comme nous l'avons vérifié en effectuant la sortie S^1 , S^2 , S D en hexadécimal.

	N	100	200	300	400	500
σ_1	S^1	5,187340	5,877946	6,282538	6,569756	6,792601
	S^2	5,187377	5,878030	6,282662	6,569928	6,792822
	S D	5,187377	5,878030	6,282663	6,569929	6,792822
σ_2	S^1	1,634939	1,639858	1,641470	1,642253	1,642706
	S^2	1,634983	1,639946	1,641605	1,642437	1,642936
	S D	1,634983	1,639946	1,641606	1,642437	1,642936
σ_3	S^1	0,6881702	0,6906486	0,6914758	0,6918888	0,6921354
	S^2	0,6881722	0,6906534	0,6914833	0,6918988	0,6921481
	S D	0,6881722	0,6906534	0,6914833	0,6918987	0,6921482
σ_4	S^1	0,8224151	0,8224491	0,8224530	0,8224527	0,8224512
	S^2	0,8224175	0,8224546	0,8224615	0,8224639	0,8224650
	S D	0,8224175	0,8224545	0,8224614	0,8224639	0,8224650

Enfin, signalons que les temps nécessités pour le calcul des valeurs S^2 et S D, relatives à la somme théorique σ_1 lorsque $N = 2000$, ont été identiques, le test de comparaison $|X| \geq |Y|$ relatif au calcul de l'erreur élémentaire d'addition par (II.1) étant omis (le test est inutile ici puisque l'on connaît les grandeurs respectives des opérandes).

I . 6 . CONCLUSION

Rappelons que le problème de l'amélioration de la précision d'une somme algébrique, calculée par une arithmétique à mot fini, apparaît dans deux circonstances, éventuellement simultanées :

1) Le nombre de termes de la somme $S = \sum_{i=1}^N X_i$ est grand.

Dans de nombreux cas, les X_i sont approximativement du même ordre de grandeur (intégration numérique, résolution d'équations différentielles avec un grand nombre de pas, calcul de somme de séries) mais, N étant grand, les arrondis ou troncatures successives deviennent excessives.

2) Les exposants des termes de la somme sont d'ordres de grandeur très différents, et peuvent ainsi introduire une erreur substantielle, même si N n'est pas très grand : c'est le phénomène de "cancellation", apparaissant en particulier quand une somme intermédiaire est beaucoup plus grande en valeur absolue que la somme finale. De nombreux chiffres significatifs sont alors perdus, et il est nécessaire de corriger les résultats obtenus.

L'obtention de la quasi double-précision d'une somme, à l'aide d'une arithmétique à simple précision, avait déjà été envisagée dès 1965 par W.KAHAN [14] et O.MOLLER [28] qui établissent un schéma de correction au 1er degré (c'est-à-dire en négligeant les erreurs du second ordre), d'une somme $S = \sum_{i=1}^N X_i$, différant de celui ci-dessus : au pas i , on ajoute préalablement à X_{i+1} la valeur ou la partie principale de l'erreur commise au pas précédent. O.MOLLER a établi, à cette occasion, une formule exacte d'obtention de l'erreur élémentaire d'addition.

Une autre approche du problème de l'amélioration de la précision d'une somme, a aussi été réalisée par JM.WOLFE [44] qui propose l'utilisation d'accumulateurs en cascade, chaque accumulateur servant à stocker des sommes partielles dont les exposants sont compris entre des limites définies. M.A.MALCOLM [24] établit une modification de l'algorithme de WOLFE, permettant d'assurer une précision fixée, le coût de la méthode étant augmenté par un découpage préalable, en un nombre q fixé de parties, de chaque terme de la somme.

Enfin D.R.ROSS [36] avait modifié l'algorithme de WOLFE pour permettre son application, non plus seulement à une somme arithmétique, mais à une somme algébrique dont le nombre de chiffres significatifs est fixé à l'avance.

II - OBTENTION DU PROBLEME PERTURBE DANS L'ANALYSE A POSTERIORI DES ERREURS D'ARRONDI

II.1. INTRODUCTION

Soit un problème (\mathcal{P}) , défini par des paramètres a_1, a_2, \dots, a_n , dont la solution x s'exprime comme le résultat d'une succession d'opérations élémentaires sur \mathbb{R} .

Notons $x = f(a_1, a_2, \dots, a_n)$. (\mathcal{P})

Lorsque la résolution est effectuée à l'aide d'une arithmétique flottante, le résultat obtenu \bar{x} diffère de x en raison des erreurs d'arrondi ; dans de nombreux cas, il est commode, réalisant une analyse à postériori des erreurs d'arrondi, d'exprimer que \bar{x} est la solution exacte (c'est-à-dire calculée sur \mathbb{R}) d'un problème (\mathcal{P}') obtenu par perturbation du problème initial, soit

$$\bar{x} = f(a_1 + e_1, a_2 + e_2, \dots, a_n + e_n) \quad (\mathcal{P}')$$

écrivons alors le problème suivant :

Obtenir, par des calculs en arithmétique à virgule flottante, la solution exacte du problème perturbé (\mathcal{P}') .

Nous soulignerons plusieurs points :

1) Pour de nombreux algorithmes, dans ce type d'analyse, les erreurs interviennent linéairement et l'on est ramené à des problèmes de sommation d'erreurs.

2) Rappelons que l'écriture du problème perturbé n'est pas, en général, unique, mais que l'intérêt de l'analyse à postériori des erreurs d'arrondi, est certain pour les problèmes expérimentaux dont les paramètres sont entachés d'erreurs de mesure ;

Le problème des erreurs d'entrée des données, peut être aussi considéré de façon analogue.

Nous traiterons, à titre d'exemple, l'algorithme de Gauss.

II.2. UN EXEMPLE : L'ALGORITHME DE GAUSS.

II.2.1. Notations

$A \in \mathcal{M}_{n,n}(\mathbb{R})$, non singulière, et $b \in \mathbb{R}^n$ étant donnés, considérons la résolution, par l'algorithme de Gauss général, du système linéaire :

$$A x = b \quad (\text{III.8})$$

A et b sont multipliés à gauche successivement par les matrices J_k ($k = 1, \dots, n-1$) :

où $\rho_{ik}^k = \frac{a_{ik}}{a_{kk}^k}$ ($i = k+1, \dots, n$)

pour former les deux suites $\{A^k\}$ et $\{b^k\}$ ($k = 1, \dots, n$) :

$$A^1 = A, \quad b^1 = b$$

$$A^{k+1} = J_k A^k, \quad b^{k+1} = J_k b^k \quad (k = 1, \dots, n-1)$$

(III.8) est remplacé par le système triangulaire équivalent $Tx = d$, que l'on résoud, et où nous avons posé pour la commodité des notations :

$$T = A^n = J_{n-1} \dots J_1 A$$

$$d = b^n = J_{n-1} \dots J_1 b$$

II.2.2. Rappel des résultats de l'analyse à postériori des erreurs d'arrondi

Dans ce paragraphe, nous rappelons des résultats formulés et démontrés par N. GASTINEL [8] (Voir aussi J. H. WILKINSON [42]) :

Par une arithmétique à virgule flottante, il est obtenu une suite de matrices $\{\bar{A}^k\}$ et une suite de seconds membres $\{\bar{b}^k\}$ ($k = 1, \dots, n$, avec $\bar{A}^1 = A$, $\bar{b}^1 = b$) :

Définissant les erreurs absolues e_i^k , e_{ij}^k , e''_{ij}^k suivantes, commises au cours de la triangularisation pour :

$$k = 1, \dots, n-1$$

$$i = k+1, \dots, n$$

$$j = k+1, \dots, n+1 \text{ si l'on pose } \bar{a}_{i,n+1}^k = \bar{b}_i^k :$$

$$e_i^k = \bar{a}_{ik}^k / (\bar{a}_{kk}^k - \rho_i^k) \quad \text{avec} \quad \rho_i^k = \bar{a}_{ik}^k \ominus \bar{a}_{kk}^k$$

$$e_{ij}^k = \frac{\bar{a}_{ij}^k}{\rho_i^k} \bar{a}_{kj}^k - \rho_i^k \otimes \bar{a}_{kj}^k$$

$$e''_{ij}^k = \bar{a}_{ij}^k - \rho_i^k \otimes \bar{a}_{kj}^k - \bar{a}_{ij}^{k+1}$$

De plus, définissant les erreurs absolues ϵ'_{ij} , ϵ''_{ij} , ϵ_i suivantes, commises lors de la résolution du système triangulaire $\bar{T}x = \bar{d}$ pour :

$$i = n, \dots, 1$$

$$j = n, \dots, i+1 :$$

$$\epsilon'_{ij} = \bar{t}_{ij} \bar{x}_j - \bar{t}_{ij} \otimes \bar{x}_j$$

$$\text{et si } S_{i,n+1} = 0, \quad S_{ik} = \oplus \sum_{l=n}^k \bar{t}_{il} \otimes \bar{x}_l$$

$$\epsilon''_{ij} = S_{i,j+1} + \bar{t}_{ij} \otimes \bar{x}_j - S_{i,j} \quad j \neq n$$

$$\epsilon''_{in} = \bar{d}_i - S_{i,i+1} - (\bar{d}_i \ominus S_{i,i+1})$$

$$\epsilon_i = (\bar{d}_i \ominus S_{i,i+1}) / \bar{t}_{ii} - \bar{x}_i$$

et notant :

\bar{J}_k la matrice déduite de J_k en remplaçant dans cette matrice les

valeurs ρ_i^k par les valeurs calculées $\bar{\rho}_i^k$

et

$$\bar{J} = \bar{J}_{n-1} \dots \bar{J}_1$$

on a :

Théorème (N. GASTINEL [8] , ch V, théorème VII)

La solution approchée du système $A x = b$ et obtenue par un calcul en arithmétique à virgule flottante, peut être considérée comme la solution exacte du système :

$$(A + \delta A) x = b + \delta b + \bar{J}^{-1} \delta \bar{d} \quad (\text{III.9})$$

où

$$\delta A = \sum_{k=1}^{n-1} M^k, \quad \delta b = \sum_{k=1}^{n-1} z^k,$$

M^k et z^k étant respectivement la matrice et le vecteur définis par :

$$M^k \begin{cases} m_{ik}^k = -e_i^k \bar{a}_{kk}^k, & i > k \\ m_{ij}^k = e_{ij}^k - e_{ij}^{''k}, & i, j = k+1, \dots, n \\ m_{ij}^k = 0 & \text{ailleurs} \end{cases}$$

$$z^k \begin{cases} (z^k)_i = 0, & i \leq k \\ (z^k)_i = e_{i,n+1}^k - e_{i,n+1}^{''k}, & i > k \end{cases}$$

et $\delta \bar{d}$ est le vecteur dont la $i^{\text{ème}}$ composante s'écrit :

$$(\delta \bar{d})_i = -\epsilon_i \bar{t}_{ii} + \left(\sum_{j=i+1}^n \epsilon_{ij}^{\prime} \right) + \left(\sum_{j=i+1}^{n-1} \epsilon_{ij}^{\prime\prime} \right) - \epsilon_{in}^{\prime\prime}$$

II.2.3. Etude expérimentale

Nous avons écrit en FORTRAN le programme suivant, calculant au 1er degré (c'est-à-dire en négligeant les erreurs du second ordre), la matrice δA et les vecteurs δb et $\delta \bar{d}$ du problème perturbé (III.9) ; les paramètres DA, DB, JM1, DD de la subroutine PGUSS représentent respectivement δA , δB , \bar{J}^{-1} et $\delta \bar{d}$.

L'écriture ES(X,Y,S), EP(X,Y,P), CD(X,Y,D) désigne les sous-programmes FUNCTION auxiliaires, de calcul des erreurs élémentaires, définis au chapitre II.


```
SUBROUTINE PGAUSS(A,B,N,X,IER,DA,DB,JM1,DD)
REAL JM1
DIMENSION A(N,N),B(N),X(N),DA(N,N),DB(N),JM1(N,N),DD(N)
DO 1 I=1,N
DB(I)=.0
JM1(I,I)=1.0
DO 1 J=1,N
DA(I,J)=.0
1 CONTINUE
IER=0
C TRIANGULARISATION DE A
NM1=N-1
DO 2 K=1,NM1
IF(A(K,K))10,20,10
20 IER=-1
GOTO 3
10 KP1=K+1
DO 2 I=KP1,N
R=A(I,K)/A(K,K)
DA(I,K)=DA(I,K)-CD(A(I,K),A(K,K),R)
JM1(I,K)=R
V=R*B(K)
W=B(I)-V
Z=-V
DB(I)=DB(I)+(EP(R,B(K),V)-ES(B(I),Z,W))
R(I)=W
DO 2 J=KP1,N
V=R*A(K,J)
W=A(I,J)-V
Z=-V
DA(I,J)=DA(I,J)+(EP(R,A(K,J),V)-ES(A(I,J),Z,W))
A(I,J)=W
2 CONTINUE
C RETOUR ARRIERE
X(N)=B(N)/A(N,N)
DD(N)=-CD(B(N),A(N,N),X(N))
DO 3 L=2,N
I=N+1-L
IP1=I+1
S=.0
Z=.0
DO 4 J=IP1,N
V=A(I,J)*X(J)
W=S+V
Z=Z+EP(A(I,J),X(J),V)+ES(S,V,W)
S=W
4 CONTINUE
S=-S
V=B(I)+S
X(I)=V/A(I,I)
DD(I)=Z-ES(B(I),S,V)-CD(V,A(I,I),X(I))
3 CONTINUE
RETURN
END
```

Nous donnerons à titre d'exemple, le résultat suivant, obtenu par l'activation du sous-programme PGAUSS, pour la résolution de $Ax = b$ avec :

$$A = \begin{vmatrix} 3 & 2 & -1 & 4 & 1 & 1 & 3 & 2 \\ 1 & -3 & 2 & 1 & 0 & -2 & 1 & 1 \\ 0,5 & 1 & 3 & 0,25 & -4 & -1 & 2 & 0 \\ 3 & 2 & -3 & 4 & 1 & 2 & -1 & 1 \\ 2 & -1 & 0,125 & 0,25 & -3 & 2 & 1,25 & -2 \\ -1 & 1 & 0 & 0,5 & 2 & 3 & 4 & 1 \\ -2 & -2 & 1 & 0 & 1 & 3 & 2 & 3 \\ 1 & -1 & 1 & 6 & -3 & 2 & 4 & 3 \end{vmatrix} \quad b = \begin{vmatrix} 19 \\ 18 \\ 16,25 \\ -2 \\ 7 \\ 8,5 \\ 1 \\ 31 \end{vmatrix}$$

$$\delta A = 10^{-8} \begin{vmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -6 & 60 & -30 & 24 & 0 & 30 & 0 & 0 \\ -12 & -17 & -95 & 0 & 65 & 96 & -13 & 5 \\ 0 & 0 & 19 & 1 & 25 & 4 & 30 & 30 \\ -12 & 36 & 40 & 56 & -165 & -184 & -1324 & 16 \\ 6 & -74 & -66 & -65 & -101 & -78 & 1860 & -126 \\ 12 & -4 & 3 & -80 & 1561 & -447 & 1470 & -121 \\ -6 & 74 & 13 & -212 & 2379 & -402 & -671 & 704 \end{vmatrix}$$

$$\delta b = 10^{-8} \begin{vmatrix} 0 \\ 89 \\ 42 \\ -75 \\ -1464 \\ -1462 \\ 3016 \\ -8165 \end{vmatrix} \quad \bar{J}^{-1} \delta \bar{a} = 10^{-8} \begin{vmatrix} 489 \\ 322 \\ 47 \\ 447 \\ 1294 \\ 713 \\ 958 \\ 36 \end{vmatrix}$$

Solution théorique x et solution calculée \bar{x} :

$$x = \begin{array}{c} 1 \\ -1 \\ 2 \\ 3 \\ 0 \\ -2 \\ 4 \\ -1 \end{array} \quad \bar{x} = \begin{array}{c} 1,000028 \\ -1 \\ 1,999955 \\ 2,999960 \\ -0,0000275 \\ -2,000015 \\ 4,000010 \\ -0,9999652 \end{array}$$

les résidus $r_1 = b - A\bar{x}$ et $r_2 = b + \delta b + \bar{J}^{-1} \delta \bar{d} - (A + \delta A) \bar{x}$
ont été évalués en double précision :

$$r_1 = 10^{-7} \begin{array}{c} -49 \\ -52 \\ -52 \\ -25 \\ -457 \\ 822 \\ 270 \\ 485 \end{array} \quad r_2 = 10^{-14} \begin{array}{c} 0 \\ -25 \\ -340 \\ -0,3 \\ -21 \\ 110 \\ -2000 \\ -2800 \end{array}$$

III - UTILISATION DE METHODES DE RAFFINEMENT ITERATIF

III . 1 . INTRODUCTION

Les méthodes de raffinement itératif, de la solution d'un système linéaire ou de l'inverse calculée d'une matrice, ne sont guère utilisables lorsque le calcul numérique du résidu est effectué avec la même précision que celle utilisée pour la résolution du problème initial.

J.H. WILKINSON [42] , [43] a étudié ces méthodes et donné les conditions de convergence du processus itératif.

Cet auteur souligne l'importance de l'accumulation sur 2 s chiffres des produits scalaires dans le calcul du résidu, afin d'éviter que les erreurs entachant ces résidus n'atteignent l'ordre de grandeur de ceux-ci. Les techniques du calcul des erreurs d'arrondi élémentaires permettent, en arithmétique flottante de simple ou de double précision, le calcul presque exact des résidus.

Nous avons considéré la méthode la plus simple de raffinement itératif de l'inverse d'une matrice.

III . 2 . UN EXEMPLE : CORRECTION DE L'INVERSE D'UNE MATRICE

III.2.1. Formule théorique

Soit A une matrice inversible, et B_0 une valeur approchée de l'inverse A^{-1} de A . Posant :

$$R = I - AB_0$$

il vient :

$$A^{-1} = B_0 (I - R)^{-1}$$

La suite $\{ B_k = B_0 \sum_{j=0}^k R^j \}$ que l'on calcule par la récurrence :

$$B_k = B_{k-1} + B_0 R^k \quad (\text{III. 10})$$

vérifie :

$$B_k = A^{-1} (I - R^{k+1})$$

et converge vers A^{-1} quand $k \rightarrow \infty$, lorsque le rayon spectral de R est inférieur à 1.

L'erreur au pas k , $A^{-1} - B_k = E_k$ est alors :

$$E_k = A^{-1} R^{k+1} \quad (\text{III. 11})$$

III. 2. 2. Calculs en arithmétique à virgule flottante

R définie ci-dessus, est calculée, donc entachée d'erreur, c'est-à-dire que si \bar{R} est la valeur "calculée" :

$$R = \bar{R} + \Omega$$

Notons $B \oplus C$ et $B \otimes C$ la somme et le produit en arithmétique flottante, de 2 matrices B et C définies sur T.

Posant $\bar{B}_0 = B_0$

$$\bar{C}_0 = B_0, \bar{C}_k = \bar{C}_{k-1} \otimes \bar{R} \quad k = 1, \dots$$

on obtient par l'arithmétique à virgule flottante, la suite de matrices $\{\bar{B}_k\}$:

$$\bar{B}_k = \bar{B}_{k-1} \oplus \bar{C}_k$$

et soient les suites $\{M_k\}$ et $\{M'_k\}$ de matrices d'erreur suivantes :

$$M_k = \bar{C}_{k-1} \bar{R} - \bar{C}_k \quad k = 1, 2, \dots$$

$$M'_k = (\bar{B}_{k-1} + \bar{C}_k) - \bar{B}_k$$

Par la sommation des égalités précédentes :

$$\bar{C}_k = B_0 \bar{R}^k - \sum_{j=1}^k M_j \bar{R}^{k-j}$$

$$\bar{B}_k = B_0 + \sum_{j=1}^k \bar{C}_j - \sum_{j=1}^k M'_j$$

et l'erreur due aux calculs $E'_k = B_k - \bar{B}_k$ est :

$$E'_k = B_0 \sum_{j=1}^k (R^j - \bar{R}^j) + \sum_{j=1}^k \sum_{i=1}^j M_i \bar{R}^{j-i} + \sum_{j=1}^k M'_j \quad (\text{III. 12})$$

L'erreur totale commise sera :

$$A^{-1} - \bar{B}_k = E_k + E'_k$$

En utilisant par exemple la norme de matrice suivante :

$$\|A\|_{\infty} = \max_i \left(\sum_{j=1}^N |a_{ij}| \right) \quad A \in \mathcal{X}_{N,N}(\mathbb{R})$$

il est possible de donner une majoration de l'erreur due aux calculs :

Nous rappellerons au préalable les résultats suivants de l'analyse d'erreur (J.H. WILKINSON [42] pp. 18-19 et 83-84) :

$X = (x_1, \dots, x_N)$ et $Y = (y_1, \dots, y_N)$ étant 2 vecteurs de $(T_b^s)^N$,

l'erreur entachant le produit scalaire de X et Y , calculé en arithmétique flottante par des opérations \oplus et \otimes correctes, et noté $X \otimes Y$, vérifie :

$$|XY - X \otimes Y| < b^{1-s_1} (N|x_1y_1| + N|x_2y_2| + (N-1)|x_3y_3| + \dots + 2|x_Ny_N|)$$

(III.13)

cette écriture supposant $Nb^{1-s} < 0, 1$, et s_1 étant défini pour une base b par :

$$s_1 = s - (0,08406 / \log_2 b)$$

b^{1-s_1} est à remplacer dans (III. 13) par :

$\frac{1}{2} b^{1-s_1}$ si \oplus et \otimes sont optimales

b^{1-s} si \oplus et \otimes sont I-dirigées pour respectivement, tout couple de T^2 de composantes de même signe, et tout couple de T^2 ,

et la conséquence :

A et B étant 2 matrices à éléments dans T , d'ordre N , l'erreur entachant le produit $A \otimes B$ des 2 matrices, calculé sur T par des opérations correctes, satisfait à :

$$|AB - A \otimes B| < b^{1-s_1} |A| |B| \quad (\text{III. 14})$$

où $D = (d_i)$ est la matrice diagonale d'éléments $d_1 = N$, $d_i = N-i+2$ $i > 1$
 et à fortiori :

$$||AB - A \otimes B||_{\infty} \leq N b^{1-s_1} ||A||_{\infty} ||B||_{\infty} \quad (\text{III. 14}')$$

Relativement à la somme $A \oplus B$, de 2 matrices A et B , calculée sur T par une addition correcte, on a :

$$||A + B - (A \oplus B)||_{\infty} \leq b^{1-s} ||A + B||_{\infty} \quad (\text{III. 15})$$

Posant alors :

$$||\bar{R}||_{\infty} = q$$

$$||\Omega||_{\infty} = \lambda ||\bar{R}||_{\infty} = \lambda q$$

majorons les normes des différents termes de (III. 12). Afin de ne pas compliquer inutilement l'écriture, nous donnons les calculs relatifs à des opérations correctes, \oplus et \otimes étant I -dirigées pour respectivement, tout couple de T^2 de composantes de même signe, et tout couple de T^2 .
 Pour une telle arithmétique, on aura en particulier :

$$||A \otimes B||_{\infty} \leq ||A||_{\infty} ||B||_{\infty} \quad (\text{III. 16})$$

(III. 16) appliquée aux matrices \bar{C}_i entraîne :

$$||\bar{C}_i||_{\infty} \leq ||B_0||_{\infty} q^i$$

et :

$$||\sum_{j=1}^k \sum_{i=1}^j M_i \bar{R}^{j-i}||_{\infty} \leq N b^{1-s} ||B_0||_{\infty} \sum_{j=1}^k j q^j$$

En prenant les majorations suivantes :

$$||\bar{B}_j||_{\infty} \leq \sum_{i=0}^j ||\bar{C}_i||_{\infty} \leq ||B_0||_{\infty} (1 + q + \dots + q^j)$$

et

$$||R^j - \bar{R}^j||_{\infty} = ||(\bar{R} + \Omega)^j - \bar{R}^j||_{\infty} \leq q^j (C_j^1 \lambda + \dots + \lambda^j) = q^j [(1 + \lambda)^j - 1]$$

et en remarquant que :

$$B_0 = A^{-1} (I - R) \text{ donc :}$$

$$\|B_0\|_\infty \leq \|A^{-1}\|_\infty (1 + (1+\lambda) q)$$

il vient finalement :

$$\frac{\|E'_k\|_\infty}{\|A^{-1}\|_\infty} < (1 + (1+\lambda) q) \left[\sum_{j=0}^k \{((1+\lambda)^j - 1 + b^{1-s} (Nj + k - j + 1))\} q^j \right]$$

En particulier, l'utilisation de l'algorithme \mathcal{A} pour le calcul de \bar{R} permet de choisir $\lambda = 2b^{1-s}$, et la formule précédente s'écrit :

$$\frac{\|E'_k\|_\infty}{\|A^{-1}\|_\infty} \leq (1 + q) b^{1-s} \sum_{j=0}^k (Nj + 2j + k + 1) q^j$$

Cette majoration donne alors, en fonction de B_0 et q , une indication sur l'erreur due aux calculs.

Soulignons que si \bar{R} est calculée sur T sans correction, $|\Omega|$ est de l'ordre de $b^{1-s} \|A\| \|D\| \|B_0\|$ et dans le cas où $\|A\| \|B_0\| \gg \|A B_0\|$ λ a une valeur non négligeable.

III. 2. 3. Exemple numérique

Nous avons effectué les calculs pour la correction de l'inverse de la matrice A d'ordre 6, d'éléments $a_{ij} = \frac{M}{i+j-1}$ ($M=27720$)

Nous avons écrit ci-dessous, en chaque position i, j :

- 1) La valeur exacte de l'élément de A^{-1}
- 2) La valeur calculée par la méthode de Gauss
- 3) La valeur obtenue après une correction
- 4) La valeur obtenue après deux corrections

(dans les 3 derniers cas, nous n'avons pas réécrit tous les chiffres obtenus, mais seulement ceux qui diffèrent de la valeur exacte).

q = 0,825

0,001298701 3864 677					
-0,02272727 57182 639	0,5303030 265818 2819 29				
0,1212121 01084 056 0	-3,181818 55308 664 6	20,36363 21120 275 2			
-0,2727272 697844 098 1	7,636363 565578 5947 1	-50,90909 50089 671 6	130,9090 29 9878 38 1		
0,2727272 694311 076 1	-7,954545 875210 076 2	54,54545 08720 277 3	-143,1818 2,1455 759 .	159,0909 8,0802 852 .	
-0,1000000 9868890 9999210 9999985	3,000000 2,968432 2,999812 2,999998	-21,00000 0,81747 0,99892 0,99997	56,00000 5,58675 5,99762 5,99997	-63,00000 2,59642 2,99771 2,99997	25,20000 05740 19920 19998

Rappelons que l'on utilise, de préférence à la méthode précédente, la méthode, de raffinement itératif de l'inverse d'une matrice, à convergence quadratique :

$$B_{k+1} = B_k (2 I - A B_k)$$

Les produits $A B_k$ pourront être calculés avec précision par les techniques du chapitre II .

IV - RESOLUTION NUMERIQUE DE PROBLEMES MAL CONDITIONNES

IV . 1 . INTRODUCTION

L'impossibilité de résoudre numériquement un problème par arithmétique flottante, est subordonnée à la longueur de la représentation des nombres.

Lorsque la possibilité de se tourner vers la double précision a été épuisée, les techniques de calcul à chaque pas de l'erreur élémentaire, permettront la résolution du problème, la limite imposée à la précision devenant fonction de l'augmentation du coût et de l'encombrement.

Comme lors de la correction d'une somme, on calcule conjointement aux valeurs à déterminer, les erreurs les entachant.

L'accumulation des erreurs ne se fait généralement plus linéairement, comme au paragraphe I.

Nous illustrons la méthode sur deux exemples simples.

IV . 2 . EVALUATION NUMERIQUE D'UN POLYNOME

Si l'on calcule la valeur $P(x)$ d'un polynôme P en un point x voisin d'un zéro mal conditionné, le danger est grand d'obtenir une valeur numérique présentant une perte importante de chiffres significatifs ou même de signe erroné. Ce phénomène, dû à des cancellations successives, est particulièrement gênant dans les processus itératifs lorsqu'il génère une décision incorrecte dans le déroulement de l'algorithme (J.H.WILKINSON [42])

Soit P un polynôme défini sur T par ses coefficients a_i ($i = 0, \dots, N$) et x un élément de T .

Le calcul de la valeur de P en x :

$$P(x) = a_0 x^N + \dots + a_N$$

par le schéma de Horner, sur respectivement \mathbb{R} et T , définit les suites b_i , \bar{b}_i et l'erreur $\epsilon_i = b_i - \bar{b}_i$ ($i = 0, \dots, N$) par :

$$\begin{cases} b_0 = a_0 \\ b_i = b_{i-1} x + a_i \end{cases} \quad i = 1, \dots, N$$

$$\begin{cases} \bar{b}_0 = a_0 \\ \bar{b}_i = \bar{b}_{i-1} \otimes x \oplus a_i \end{cases} \quad i = 1, \dots, N$$

Désignant par e_i et e'_i les erreurs absolues commises au pas i :

$$e_i = \bar{b}_{i-1} x - (\bar{b}_{i-1} \otimes x) \quad i = 1, \dots, N$$

$$e'_i = (\bar{b}_{i-1} \otimes x + a_i) - \bar{b}_i \quad i = 1, \dots, N$$

il vient immédiatement que la suite ε_i vérifie la relation de récurrence :

$$\begin{cases} \varepsilon_0 = 0 \\ \varepsilon_i = \varepsilon_{i-1} x + e_i + e'_i \end{cases} \quad i = 1, \dots, N$$

L'adjonction au calcul de \bar{b}_i , de l'algorithme :

$$\begin{cases} \bar{\varepsilon}_0 = 0 \\ \bar{\varepsilon}_i = \bar{\varepsilon}_{i-1} \otimes x \oplus e_i \oplus e'_i \end{cases} \quad i = 1, \dots, N$$

permet, par la correction $\bar{P}(x) = \bar{b}_N \oplus \bar{\varepsilon}_N$, de regagner de nombreux chiffres significatifs.

EXEMPLE : Evaluation au point $x = N - 0,01$ du polynôme $P_N(x) = \prod_{j=1}^N (x-j)$

Les valeurs obtenues pour l'arithmétique IBM-360 en double précision ont été :

N	VALEUR MACHINE	VALEUR CORRIGEE	VALEUR EXACTE
5	-0.23503490009996 <u>5</u> D 00	-0.235034900100000D 00	-0.235034900100000D 00
6	-0.117282415150 <u>6</u> 80D 01	-0.117282415149900D 01	-0.117282415149900D 01
7	-0.70252166671 <u>5</u> 256D 01	-0.702521666747902D 01	-0.702521666747902D 01
8	-0.4910626451 <u>4</u> 5694D 02	-0.491062645056783D 02	-0.491062645056783D 02
9	-0.3923590534 <u>2</u> 2144D 03	-0.392359053400370D 03	-0.392359053400370D 03
10	-0.35273078791 <u>8</u> 950D 04	-0.352730789006932D 04	-0.352730789006932D 04
11	-0.352378058 <u>3</u> 79144D 05	-0.352378058217926D 05	-0.352378058217925D 05
12	-0.38726348 <u>2</u> 490301D 06	-0.387263485981500D 06	-0.387263485981500D 06
13	-0.4643289 <u>3</u> 9717770D 07	-0.464328919691819D 07	-0.464328919691819D 07
14	-0.603162 <u>3</u> 82884369D 08	-0.603163266679672D 08	-0.603163266679672D 08
15	-0.843825491 <u>9</u> 82666D 09	-0.843825410084862D 09	-0.843825410084861D 09
16	-0.126490201741 <u>1</u> 72D 11	-0.126489428971718D 11	-0.126489428971718D 11
17	-0.2022596617919 <u>3</u> 8D 12	-0.202256596925812D 12	-0.202256596925812D 12
18	-0.343649541515100D 13	-0.343633958176954D 13	-0.343633958176954D 13
19	-0.619193936971520D 14	-0.618197490760341D 14	-0.618197490760340D 14

Nous avons souligné les chiffres inexacts.

On remarque la concordance entre valeurs exacte et corrigée, et la perte progressive de chiffres significatifs dans l'évaluation simple de $P(x)$.

IV . 3 . L'ELIMINATION DE GAUSS

De la même façon qu'en IV . 2 , dans la résolution par l'algorithme de Gauss, du système linéaire :

$$A x = b$$

où $A = (a_{ij})$ et $b = (b_i)$ sont définis sur T ,

on associe aux calculs les algorithmes auxiliaires suivants :

- les notations sont celles du § II -

A l'algorithme de triangularisation :

$$\bar{\rho}_i^{-k} = a_{ik}^{-k} \oslash \bar{a}_{kk}^{-k}$$

$$\bar{a}_{ij}^{-k+1} = a_{ij}^{-k} \ominus \bar{\rho}_i^{-k} \otimes \bar{a}_{kj}^{-k}$$

on adjoint l'algorithme suivant, calculant la partie principale de l'erreur $\theta_{ij}^k = a_{ij}^k - \bar{a}_{ij}^k$:

$$(\bar{b}_i^{-k} = \bar{a}_{i,n+1}^{-k})$$

$$\Delta \bar{\rho}_i^{-k} = (e_i^k \cdot \bar{a}_{kk}^{-k} \oplus \bar{\theta}_{ik}^{-k} \ominus \bar{\rho}_i^{-k} \otimes \bar{\theta}_{kk}^{-k}) \oslash \bar{a}_{kk}^{-k}$$

$$\bar{\theta}_{ij}^{-k+1} = \bar{\theta}_{ij}^{-k} \oplus e_{ij}^{nk} \ominus e_{ij}^{k} \ominus \Delta \bar{\rho}_i^{-k} \otimes \bar{a}_{kj}^{-k} \ominus \bar{\theta}_{kj}^{-k} \otimes \bar{\rho}_i^{-k}$$

A l'algorithme de retour-arrière :

$$\bar{x}_i = \left[\bar{d}_i \ominus \left\{ \oplus \sum_{j=i+1}^n \bar{t}_{ij} \otimes \bar{x}_j \right\} \right] \oslash \bar{t}_{ii}$$

on adjoint le calcul de la partie principale de l'erreur $\xi_i = x_i - \bar{x}_i$:

E C et E P ci-dessous représentent respectivement, l'erreur "générée" dans le calcul de \bar{x}_i , et l'erreur "propagée" par les calculs antécédents.

$$E P = \bar{\theta}_{i,n+1}^n \ominus \left[\sum_{j=i+1}^n \bar{\theta}_{ij}^n \otimes \bar{x}_j \oplus \bar{t}_{ij} \otimes \bar{\xi}_j \right] \ominus \bar{x}_i \otimes \bar{\theta}_{ii}^n$$

$$E C = \ominus \sum_{j=i+1}^n \varepsilon'_{ij} \ominus \sum_{j=i+1}^{n-1} \varepsilon''_{ij} \oplus \varepsilon''_{in} \oplus \varepsilon_i \bar{t}_{ii}$$

$$\bar{\xi}_i = (E P \oplus E C) \oslash \bar{t}_{ii}$$

TECHNIQUE DE CORRECTION : Remarquons que la correction des calculs peut s'effectuer selon les procédés suivants :

1) Correction en fin de résolution uniquement : Au vecteur \bar{x} calculé par l'algorithme de Gauss, on ajoute, à la fin des calculs, le vecteur correction $\bar{\xi}$.

2) La correction 1 est précédée, lorsque la triangularisation de la matrice A est terminée et avant le retour-arrière, de la correction de la matrice triangulaire T et du second membre d obtenus.

3) Correction, à chaque pas, des valeurs calculées par l'algorithme de Gauss :

- par addition de $\bar{\theta}_{ij}^{k+1}$ à \bar{a}_{ij}^{k+1} $\forall k, i, j$

- par addition de $\bar{\xi}_i$ à \bar{x}_i $\forall i$

l'erreur résultant de l'addition sur T des quantités ci-dessus étant remplacée dans respectivement $\bar{\theta}_{ij}^{k+1}$ ou $\bar{\xi}_i$.

Seule la troisième méthode assure, pour le coût, une précision suffisante.

EXEMPLE : Résolution de $H_N x = u$, où H_N est la matrice de Hilbert d'ordre N , u le 1er vecteur de la base naturelle de \mathbb{R}^N .

Les résultats suivants ont été obtenus :

	Chiffres significatifs	
Solution calculée par l'algorithme de Gauss :	5	4
Solution corrigée par le procédé 1 :	11	7 ou 8
Solution corrigée par le procédé 2 :	12	9 ou 10
Solution corrigée par le procédé 3 :	15	15
	N = 10	N = 11

Ainsi :

N = 11	VALEUR MACHINE	VALEUR CORRIGEE
	<u>0.120994980151904D</u> 03	0.121000000000000D 03
	<u>-0.725946209190183D</u> 04	-0.726000000000000D 04
	<u>0.141555777288999D</u> 06	0.141570000000000D 06
	<u>-0.132115844472851D</u> 07	-0.132132000000000D 07
	<u>0.693595457285204D</u> 07	0.693693000000000D 07
	<u>-0.221947070248717D</u> 08	-0.221981760000000D 08
	<u>0.449172516859284D</u> 08	0.449248800000000D 08
	<u>-0.577500700985906D</u> 08	-0.577605600000000D 08
	<u>0.457183299670648D</u> 08	0.457271100000000D 08
	<u>-0.203190701897057D</u> 08	-0.203231600000000D 08
	<u>0.387906321794531D</u> 07	0.387987600000000D 07

COUT DE LA METHODE : La correction à chaque pas des valeurs calculées par l'algorithmme de Gauss, est coûteuse.

Le temps nécessité par les calculs est environ multiplié par 7 lorsque les erreurs d'opérations élémentaires sont déterminées par les méthodes du chapitre II .

Ce temps de calcul est multiplié par 4 lorsque l'accès au poids faible du résultat d'une opération flottante élémentaire peut être effectué par une instruction arithmétique (exemple : CDC série 6000) .

PROGRAMME FORTRAN : Les sous-programmes FUNCTION auxiliaires utilisés : ES (X, Y, S) , EP (X, Y, P) , CD (X, Y, D) sont ceux du chapitre II.

```

SUBROUTINE CGAUSS(A,B,N,X,EA,EB,EX)
REAL*8 A(20,20),EA(20,20),B(20),EB(20),X(20),EX(20)
REAL*8 ES,EP,CD
DOUBLE PRECISION R,ER,V,S,Z,W
NM1=N-1
DO 7 I=1,N
EB(I)=0.000
DO 7 J=1,N
EA(I,J)=0.000
7 CONTINUE
DO 2 K=1,NM1
KP1=K+1
DO 2 I=KP1,N
R=A(I,K)/A(K,K)
ER=EA(I,K)-EA(K,K)*R+CD(A(I,K),A(K,K),R)
ER=ER/A(K,K)
S=R
R=S+ER
ER=ES(S,ER,R)
V=R*B(K)
S=B(I)-V
EB(I)=EB(I)-ER*B(K)-R*EB(K)-EP(R,B(K),V)+ES(B(I),-V,S)
B(I)=S+EB(I)
EB(I)=ES(S,EB(I),B(I))
DO 2 J=KP1,N
V=R*A(K,J)
S=A(I,J)-V
EA(I,J)=EA(I,J)-ER*A(K,J)-EA(K,J)*S-EP(R,A(K,J),V)+ES(A(I,J),-V,S)
A(I,J)=S+EA(I,J)
EA(I,J)=ES(S,EA(I,J),A(I,J))
2 CONTINUE
X(N)=B(N)/A(N,N)
EX(N)=EB(N)-EA(N,N)*X(N)+CD(B(N),A(N,N),X(N))
EX(N)=EX(N)/A(N,N)
S=X(N)
X(N)=S+EX(N)
EX(N)=ES(S,EX(N),X(N))
DO 3 L=2,N
I=N+1-L
IP1=I+1
S=0.000
Z=0.000
R=0.000
DO 4 J=IP1,N
V=A(I,J)*X(J)
W=S+V
Z=Z+EP(A(I,J),X(J),V)+ES(S,V,W)
S=W
R=R+EA(I,J)*X(J)+A(I,J)*EX(J)
4 CONTINUE
S=-S
V=B(I)+S
X(I)=V/A(I,I)
Z=-Z+ES(B(I),S,V)+CD(V,A(I,I),X(I))
R=-R+EB(I)-X(I)*EA(I,I)
EX(I)=(Z+R)/A(I,I)
S=X(I)
X(I)=S+EX(I)
EX(I)=ES(S,EX(I),X(I))
3 CONTINUE
RETURN
END

```


V - CRITERE D'APPLICATION POUR LES METHODES DES § III et IV

La question essentielle qui se pose, avant l'application à une solution calculée, d'une méthode de correction, est la question de la validité de la solution approchée.

Pour déterminer si une méthode de raffinement est nécessaire et possible, on appliquera, au préalable, au résultat approché, les méthodes de contrôle de résultats d'algorithmes numériques définies par

M. LA PORTE et J. VIGNES [19] :

La méthode générale de permutation-perturbation définie par ces auteurs, ou le critère des résidus normés pour la résolution de systèmes linéaires ou d'équations algébriques [20] , permettra d'évaluer la précision de la solution approchée.

Seules les solutions de précision "moyenne" devront et pourront être corrigées .

CHAPITRE IV

L'ERREUR D'ARRONDI DANS LES METHODES ITERATIVES

LINEAIRES

INTRODUCTION

Ce chapitre est relatif à la résolution numérique d'un système d'équations linéaires par itérations, en arithmétique à virgule flottante.

Deux problèmes sont abordés :

1) Un système d'équations ayant été résolu par la convergence, sur le système de nombres à virgule flottante, de la suite des itérés d'une méthode itérative linéaire, existe-t-il une possibilité d'amélioration de la précision ?

Une méthode est proposée, utilisant les techniques définies au chapitre II.

2) Possibilité de prévoir à priori, en fonction de la matrice de l'itération linéaire et du vecteur initial, le type de convergence satisfait par la suite des itérés calculés sur le système T de nombres à virgule flottante.

I - RAFFINEMENT DE LA SOLUTION NUMERIQUE D'UN SYSTEME LINEAIRE, CALCULEE PAR UNE METHODE ITERATIVE.

I . 1 . NOTATIONS

T désigne le système T_b^S de nombres à virgule flottante muni des opérations \oplus , \ominus , \otimes , \oslash .

$M = (m_{ij})$, $x = (x_i)$ et $y = (y_i)$ étant respectivement une matrice d'ordre N définie sur T et 2 vecteurs de T^N , $M \otimes x$ et $x \oplus y$ représenteront les vecteurs définis sur T^N par :

$$M \otimes x = \left(\oplus_{j=1}^N m_{ij} \otimes x_j \right)$$

$$x \oplus y = (x_i \oplus y_i)$$

$B = (b_{ij}) \in \mathcal{M}_{N,N}(T)$ et $c = (c_i) \in T^N$ étant donnés, considérons la résolution du système linéaire :

$$x = B x + c \quad (\text{IV.1})$$

$x^{(0)} = \bar{x}^{(0)}$ étant un vecteur initial sur T^N , soient $\{x^{(n)}\}$ et $\{\bar{x}^{(n)}\}$ les suites définies respectivement sur T^N et T^N par :

$$x^{(n+1)} = B x^{(n)} + c \quad (\text{IV.2})$$

$$\bar{x}^{(n+1)} = B \otimes \bar{x}^{(n)} \oplus c \quad (\text{IV.3})$$

La convergence de la suite $\{x^{(n)}\}$ vers la solution x de (IV.1), est assurée à la condition nécessaire et suffisante que le rayon spectral de B soit inférieur à 1, ce que nous supposerons.

La suite $\{\bar{x}^{(n)}\}$ est dite "numériquement" convergente lorsqu'est réalisée l'une des deux possibilités suivantes :

a) Il existe un entier K vérifiant :

$$\bar{x}^{(K+1)} = \bar{x}^{(K)} \quad (\text{convergence forte})$$

b) Il existe un entier K à partir duquel les itérés cyclent sur un nombre fini et petit de valeurs, c'est-à-dire qu'il existe K et y^1, \dots, y^p vérifiant :

$$\bar{x}^{(K+k)} = y^j \quad \forall k \geq 0 \wedge k \equiv j - 1 \pmod{p} \text{ (convergence faible)}$$

I . 2 . METHODE

Lorsque (IV.1) est résolu numériquement par la convergence sur T de la suite $\{\bar{x}^{(n)}\}$, le vecteur erreur - écart entre les solutions calculées sur R et T de (IV.1) - est solution d'un système linéaire de même matrice que (IV.1). En effet :

I . 2 . 1 $\{\bar{x}^{(n)}\}$ converge fortement

La limite $\bar{x} = (\bar{x}_i)$ de $\{\bar{x}^{(n)}\}$ vérifie :

$$\bar{x} = B \otimes \bar{x} \oplus c$$

Soit $e = (e_i)$ le vecteur défini par :

$$e = (B \bar{x} + c) - \bar{x} \quad (\text{IV.4})$$

(IV.1) et (IV.4) entraînent pour le vecteur erreur $\varepsilon = x - \bar{x}$:

$$\varepsilon = B \varepsilon + e \quad (\text{IV.5})$$

I . 2 . 2 $\{\bar{x}^{(n)}\}$ converge faiblement

On a les égalités :

$$\left\{ \begin{array}{l} y^2 = B \otimes y^1 \oplus c \\ y^3 = B \otimes y^2 \oplus c \\ \dots \\ y^1 = B \otimes y^p \oplus c \end{array} \right.$$

Particularisant y^1 par exemple, considérons ce vecteur comme la solution de (IV.1) calculée sur T .

Soit e' le vecteur défini par :

$$e' = B y^1 + c - y^2 \quad (\text{IV.6})$$

il vient pour le vecteur erreur $\varepsilon = x - y^1$:

$$\varepsilon = B \varepsilon + e \quad \text{avec} \quad e = e' + y^2 - y^1 \quad (\text{IV.7})$$

Un procédé de "correction" peut donc être défini par la résolution de (IV.5) - resp. (IV.7) - selon la méthode itérative employée pour la résolution de (IV.1) .

Les ordres de convergence asymptotique, de la méthode de résolution du système et de la méthode de correction, sont égaux, entraînant précisions et vitesses de convergence analogues. Le procédé peut être itéré.

I . 2 . 3 Calcul de e et e' définis en (IV.4) et (IV.6)

e(ou e') représente l'erreur d'arrondi générée par une itération sur T . Rappelons :

Notant u un vecteur de T^N et définissant V et \bar{V} par :

$$V = B u + c$$

$$\bar{V} = B \otimes u \oplus c$$

si e_{ij}^P et e_{ij}^S ($1 \leq j \leq N$) sont respectivement les erreurs absolues de multiplication et d'addition dans le calcul de la ième composante de \bar{V} :

$$\bar{V}_i = \left(\oplus \sum_{j=1}^N b_{ij} \otimes u_j \right) \oplus c_i \quad 1 \leq i \leq N$$

le vecteur $V - \bar{V}$ s'exprime en fonction des erreurs élémentaires définies, que l'on calculera par les méthodes décrites au chapitre II, suivant :

$$(V - \bar{V})_i = \sum_{j=1}^N e_{ij}^P + \sum_{j=1}^N e_{ij}^S \quad 1 \leq i \leq N$$

Remarque Lorsque le système linéaire à résoudre n'est pas explicitement donné sous la forme (IV.1) mais par :

$$A x = b \quad \text{où } A \in \mathcal{K}_{N,N}(T) \text{ et } b \in T^N$$

l'expression de e doit être légèrement modifiée suivant la méthode itérative employée pour résoudre le système.

Ainsi, pour les méthodes de Jacobi ou de Gauss-Seidel, la convergence, forte par exemple, de la suite numérique s'écrit :

$$\bar{x}_i = \left[\left\{ \oplus \sum_{\substack{j=1 \\ j \neq i}}^N (-a_{ij}) \otimes \bar{x}_j \right\} \oplus b_i \right] \oslash e_{ii} \quad 1 \leq i \leq N$$

Le vecteur erreur ϵ vérifie :

$$\epsilon_i = \frac{\sum_{j \neq i} (-a_{ij}) \epsilon_j + e_i}{a_{ii}} \quad 1 \leq i \leq N$$

avec :

$$e_i = \sum_{\substack{j=1 \\ \neq i}}^N e_{ij}^P + \sum_{j=1}^{N-1} e_{ij}^S + a_{ii} e_i^D$$

et les notations suivantes :

$$e_{ij}^P = (-a_{ij} \bar{x}_j) - (-a_{ij} \otimes \bar{x}_j) \quad j \neq i$$

e_{ij}^S = erreurs absolues commises lors des sommes partielles dans le calcul

de :

$$\{ \oplus \sum_{j \neq i} (-a_{ij}) \otimes \bar{x}_j \} \oplus b_i = s_i$$

$$e_i^D = s_i / a_{ii} - \bar{x}_i$$

(le nombre $a_{ii} e_i^D$ se calcule exactement sur T).

I . 3 . BORNES DE L'ERREUR

I.3.1. Une borne supérieure du vecteur erreur ϵ

Le vecteur ϵ vérifie :

$$\epsilon = (I-B)^{-1} e$$

Avec les notations des pages 73 - 74 où $D = (d_i)$ désigne la matrice diagonale d'ordre N , d'éléments $d_1 = N$, $d_i = N-i+2$ $i > 1$, il vient :

$$|e| \leq b^{1-s_1} (|B| D + I) |\bar{x}|$$

$$|e'| \leq b^{1-s_1} (|B| D |y^1| + |y^2|)$$

et en normes, pour l'une des trois normes de matrices $\| \cdot \|_1$, $\| \cdot \|_2$ ou $\| \cdot \|_\infty$, suivant le type de convergence de $\{ x^{(n)} \}$:

$$\| \varepsilon \| \leq b^{1-s_1} \| (I-B)^{-1} \| (N \| |B| \| + 1) \| \bar{x} \|$$

$$\| \varepsilon \| \leq \| (I-B)^{-1} \| [b^{1-s_1} (N \| |B| \| \| y^1 \| + \| y^2 \|) + \| y^2 - y^1 \|] \quad (IV.8)$$

où l'on peut aussi faire apparaître le conditionnement de $I - B$.

Cas particulier important où $\| B \| < 1$

Les formules (IV.8) deviennent, avec la majoration $\| (I-B)^{-1} \| \leq \frac{1}{1 - \| B \|}$

$$\| \varepsilon \| \leq b^{1-s_1} \frac{N \| |B| \| + 1}{1 - \| B \|} \| \bar{x} \|$$

$$\| \varepsilon \| \leq \frac{b^{1-s_1} (N \| |B| \| \| y^1 \| + \| y^2 \|) + \| y^2 - y^1 \|}{1 - \| B \|}$$

I.3.2. Amélioration de la précision après une correction

Numériquement, la résolution de :

$$\varepsilon = B \varepsilon + e$$

et la correction de la solution \bar{x} ou y^1 , est réalisée par les calculs suivants :

1) calcul sur T de e

L'élément \bar{e} de T est déterminé avec :

$$e = \bar{e} + \eta$$

2) résolution sur T de $\varepsilon = B \varepsilon + \bar{e}$

le vecteur $\bar{\varepsilon}$ de T^N est calculé .

Considérons par exemple le cas de convergence numérique forte des itérés sur T du système initial et du système de correction.

Le nouveau vecteur solution x_c est :

$$x_c = \bar{x} + \bar{\varepsilon}$$

vérifiant avec la notation $e^2 = (B \bar{\varepsilon} + \bar{e}) - \bar{e}$

$$x - x_c = (I-B)^{-1} (\eta + e^2) \quad (IV.9)$$

Si les opérations \otimes sont correctes, on a au moins les majorations suivantes

$$|n_i| \leq (2N-1) b^{1-s} (1+2N b^{1-s}) \left(\sum_{j=1}^N |e_{ij}^P| + \sum_{j=1}^N |e_{ij}^S| \right)$$

$$\|n\| \leq 2N b^{2-2s_1} (N \|B\| + 1) \|\bar{x}\|$$

et

$$\|e^2\| \leq b^{1-s_1} (N \|B\| + 1) \|\bar{\epsilon}\|$$

soit, en négligeant les erreurs du 2ème ordre, lorsque $\|B\| < 1$:

$$\|x-x_c\| \leq b^{2-2s_1} \frac{N \|B\| + 1}{1 - \|B\|} \left(2N + \frac{N \|B\| + 1}{1 - \|B\|} \right) \|\bar{x}\|$$

Les majorations précédentes sont souvent, malheureusement, beaucoup trop larges. Nous considérons ci-dessous, le problème de l'évaluation précise de l'incertitude sur la solution numérique à corriger.

I. 4. EVALUATION, AVANT CORRECTION, DE L'INCERTITUDE SUR LA SOLUTION NUMERIQUE A RAFFINER

Dans ce paragraphe, nous utilisons les résultats de M. LA PORTE et J. VIGNES [19] , [20] , relatifs au contrôle des erreurs d'arrondi :

Lorsque (IV. 1) a été résolu numériquement par la convergence sur T de la suite $\{\bar{x}^{(n)}\}$ vers \bar{x} , il est, comme pour une méthode directe de résolution, nécessaire, avant d'aborder la correction, de connaître le nombre de chiffres significatifs des composantes du vecteur solution numérique. En effet, rappelons :

i) De même que précédemment, si tous les chiffres de la solution ont été obtenus significatifs, une correction est évidemment inutile, mais si aucun ne l'a été, la méthode proposée sera insuffisante pour améliorer la précision.

ii) La connaissance du nombre de chiffres significatifs permettra de déterminer combien de fois doit être itérée la méthode de correction : ainsi, si les calculs s'effectuent par exemple sur un système de nombres flottants à 6 chiffres hexadécimaux de mantisse, l'obtention de 2 chiffres hexadécimaux corrects à la 1ère résolution du système, appellera 2 corrections successives.

La détermination du nombre de chiffres significatifs d'un résultat numérique par la méthode de permutation-perturbation est effectuée selon le schéma suivant [20] :

Désignant par :

$$Z^v = \{ z^1, z^2, \dots, z^v \}$$

une population de vecteurs solutions sur T de (IV. 1), on définit les vecteurs ci-dessous :

$$\bar{z} \text{ de composantes } \bar{z}_j = \frac{1}{v} \sum_{i=1}^v z_j^i$$

$$\delta^2 \text{ de composantes } \delta_j^2 = \frac{1}{v-1} \sum_{i=1}^v (z_j^i - \bar{z}_j)^2$$

$$\hat{\epsilon} \text{ de composantes } \hat{\epsilon}_j = \sqrt{(z_j^1 - \bar{z}_j)^2 + \delta_j^2}$$

$$C^v \text{ de composantes } C_j^v = \log_{10} \left| \frac{z_j^1}{\hat{\epsilon}_j} \right|$$

L'algorithme d'obtention du nombre de chiffres significatifs de chaque composante du résultat numérique est alors le suivant :

- 1) On résoud numériquement (IV. 1) sur T : la solution \bar{x} est calculée.
- 2) Choissant \bar{x} comme vecteur initial, on calcule par itérations, à l'aide d'une arithmétique flottante différente de celle précédemment utilisée, une seconde solution de (IV.1) sur le système de nombres à virgule flottante, notée z^1 .

La nouvelle arithmétique s'obtient en prenant aléatoirement \otimes_{∇} ou \otimes_{Δ} pour chacune des opérations définies sur T .

3) On calcule C^2 correspondant à $Z^2 = \{ \bar{x}, z^1 \}$. Deux cas sont alors possibles :

a) si $C_j^2 < 1 \quad \forall j$, la solution \bar{x} est à rejeter.

b) si l'inégalité précédente n'est pas vérifiée, on calcule de nouveaux éléments de la population. Chaque nouvel élément s'obtient par résolution itérative de (IV. 1), à l'aide d'opérations flottantes générées aléatoirement comme précédemment ; le vecteur initial est pris égal à \bar{x} ou à l'une des solutions calculées.

A chaque nouvel élément z^i on associe C^{i+1} correspondant à $Z^i = \{ \bar{x}, z^1, \dots, z^i \}$ et :

$C_j^{i+1} < 1 \quad \forall j$ entraîne la même conclusion que précédemment.

Sinon : ou C^{i+1} est stationnaire et chacune de ses composantes donne le nombre de chiffres significatifs de la composante correspondante de \bar{x} ou C^{i+1} n'a pas atteint sa valeur stationnaire et il faut poursuivre le procédé.

Soulignons que la méthode de permutation-perturbation pour la détermination du nombre de chiffres significatifs d'un résultat numérique calculé par méthode itérative, est peu coûteuse, un petit nombre d'itérations étant généralement nécessaire pour calculer, à partir d'un vecteur initial solution de (IV.1) sur le système de nombres à virgule flottante, une autre solution de (IV.1) sur ce système.

I. 5. EXEMPLE NUMERIQUE

Nous avons considéré la résolution, par les méthodes de Jacobi et de Gauss - Seidel, du système linéaire $Ax = b$ pour :

$$A = \begin{bmatrix} 3 & -1 & -0,5 & -0,25 & -0,5 & -0,125 & 0 & -0,5 \\ -2 & 4 & -0,25 & -1 & 0 & -0,25 & 0 & -0,25 \\ -4 & -1 & 8,25 & -0,5 & -0,5 & -1 & -0,5 & -0,25 \\ -1 & -0,5 & -1 & 4 & 0 & 0 & -1 & -0,25 \\ 0 & -1 & -1 & -1,5 & 4,5 & -0,5 & -0,25 & 0 \\ -5 & -1,5 & 0 & -1,5 & 0 & 9 & -0,5 & 0 \\ -0,5 & -2,5 & -0,5 & 0 & -0,5 & 0 & 5 & 0 \\ -3 & -0,5 & -0,5 & 0 & 0 & -0,25 & -0,25 & 5 \end{bmatrix} \quad B = \begin{bmatrix} 0,125 \\ 0,25 \\ 0,5 \\ 0,25 \\ 0,25 \\ 0,5 \\ 1 \\ 0,5 \end{bmatrix}$$

Les vecteurs x , \bar{x} et x_c désignent respectivement :

la solution exacte du système

la solution obtenue après convergence forte en K_1 itérations

la solution corrigée après K_2 itérations supplémentaires

$$x = \begin{bmatrix} 1 \\ " \\ " \\ " \\ " \\ " \\ " \\ " \end{bmatrix} \quad \bar{x} = \begin{bmatrix} 0,9999905 \\ " \quad 09 \\ " \quad 11 \\ " \quad 14 \\ " \quad 15 \\ " \quad 12 \\ " \quad 25 \\ " \quad 14 \end{bmatrix} \quad x_c = \begin{bmatrix} 0,9999999 \\ " \\ " \\ " \\ " \\ " \\ " \\ " \end{bmatrix}$$

avec :

$K_1 = 208$, $K_2 = 65$ pour la méthode de Jacobi

$K_1 = 106$, $K_2 = 33$ pour la méthode de Gauss - Seidel

Le vecteur initial $x^{(0)}$ ayant, dans tous les cas, été choisi égal au vecteur nul.

Une seule correction était nécessaire, puisque la méthode de permutation-perturbation conduisait à la stationnarité sur 5 du nombre de chiffres significatifs de chaque composante, après calcul des deux solutions supplémentaires z^1 et z^2 :

$$z^1 = \begin{array}{l} 0.9999946E 00 \\ 0.9999957E 00 \\ 0.9999949E 00 \\ 0.9999959E 00 \\ 0.9999953E 00 \\ 0.9999956E 00 \\ 0.9999959E 00 \\ 0.9999956E 00 \end{array} \quad c^2 = \begin{array}{l} 5,4 \\ 5,4 \\ 5,5 \\ 5,4 \\ 5,5 \\ 5,4 \\ 5,5 \\ 5,4 \end{array}$$

$$z^2 = \begin{array}{l} 0.9999997E 00 \\ 0.9999986E 00 \\ 0.9999998E 00 \\ 0.9999988E 00 \\ 0.9999996E 00 \\ 0.9999991E 00 \\ 0.9999995E 00 \\ 0.9999990E 00 \end{array} \quad c^3 = \begin{array}{l} 5 \\ 5 \\ 5 \\ 5 \\ 5 \\ 5 \\ 5 \\ 5 \end{array} \quad \text{(partie entière)}$$

II - ITERATIONS LINEAIRES DEFINIES PAR UNE MATRICE POSITIVE (RESP. NEGATIVE)

Des remarques particulières peuvent être faites, lorsque la matrice $B = (b_{ij})$ vérifie $b_{ij} \geq 0 \quad \forall i, j$ (resp. $b_{ij} \leq 0 \quad \forall i, j$).

Considérons le cas $B \geq 0$:

Si le vecteur $x^{(1)} - x^{(0)}$ a toutes ses composantes de même signe, la suite des itérés théoriques sera monotone par l'égalité :

$$x^{(n+1)} - x^{(n)} = B (x^{(n)} - x^{(n-1)}) = B^n (x^{(1)} - x^{(0)}) \quad (IV. 10)$$

D'autre part, par exemple de l'hypothèse :

$$x^{(1)} - x^{(0)} \geq 0$$

soit

$$B x^{(0)} + c - x^{(0)} \geq 0$$

$$B x^{(0)} + x - Bx - x^{(0)} \geq 0$$

$$(I - B) (x - x^{(0)}) \geq 0$$

il vient, en remarquant que $(I - B)^{-1}$ est une matrice positive ($\rho(B) < 1$) :

$$x \geq x^{(0)}$$

Dans l'hypothèse $x^{(1)} \geq x^{(0)}$ (resp. $x^{(1)} \leq x^{(0)}$), la suite $\{x^{(n)}\}$ est monotone croissante (resp. décroissante), bornée par x .

La propriété de monotonie est conservée par la suite $\{x^{(n)}\}$ des itérés calculés sur T , lorsque dans le calcul des itérés successifs, les opérateurs en positions identiques gardent, pour tous les calculs, la même définition (par exemple \otimes_{∇} ou \otimes_{Δ}) :

Deux tableaux de N lignes et N colonnes, définissant en effet les opérations de multiplication $\otimes = (\otimes_{ij})$ et d'addition $\oplus = (\oplus_{ij})$, dès que les opérations \otimes sont isotones, il vient :

L'hypothèse :

$$\bar{x}_j^{(n-1)} \leq \bar{x}_j^{(n)} \quad 1 \leq j \leq N$$

entraînera :

$$b_{ij} \bar{x}_j^{(n-1)} \leq b_{ij} \bar{x}_j^{(n)} \quad 1 \leq i \leq N, \quad 1 \leq j \leq N$$

$$b_{ij} \otimes_{ij} \bar{x}_j^{(n-1)} \leq b_{ij} \otimes_{ij} \bar{x}_j^{(n)} \quad " \quad "$$

et par l'isotonie des opérations d'addition :

$$\oplus_{ij} \left(\sum_{j=1}^N b_{ij} \otimes_{ij} \bar{x}_j^{(n-1)} \right) \oplus_{iN} c_i \leq \oplus_{ij} \left(\sum_{j=1}^N b_{ij} \otimes_{ij} \bar{x}_j^{(n)} \right) \oplus_{iN} c_i$$

d'où :

$$\bar{x}_i^{(n)} \leq \bar{x}_i^{(n+1)} \quad 1 \leq i \leq N$$

$$\bar{x}^{(0)} \leq \bar{x}^{(1)} \wedge B \geq 0 \implies \bar{x}^{(n)} \leq \bar{x}^{(n+1)} \quad \forall n$$

$$\bar{x}^{(0)} \geq \bar{x}^{(1)} \wedge " \implies \bar{x}^{(n)} \geq \bar{x}^{(n+1)} \quad "$$

La suite de vecteurs $\{ \bar{x}^{(n)} \}$ est non décroissante si $\bar{x}^{(0)} \leq \bar{x}^{(1)}$, non croissante si $\bar{x}^{(0)} \geq \bar{x}^{(1)}$:

Si la suite converge, elle converge fortement.

De plus, dans le cas $\bar{x}^{(0)} \leq \bar{x}^{(1)}$ par exemple, soit, s'il en existe, $\bar{x} \in T^N$ vérifiant :

$$\bar{x} = B \otimes \bar{x} \oplus c \wedge \bar{x}^{(0)} \leq \bar{x} \quad (\text{IV. 11})$$

(IV. 11) implique :

$$\bar{x}^{(1)} \leq \bar{x}$$

et par récurrence :

$$\bar{x}^{(n)} \leq \bar{x} \quad \forall n$$

d'où la :

Proposition IV. 1.

Si $b_{ij} \geq 0 \quad \forall i, j$

si les opérateurs \otimes_{ij} sont isotones et fixés $\forall i, j$,

si $\bar{x}^{(1)} \geq \bar{x}^{(0)}$ ou $\bar{x}^{(1)} \leq \bar{x}^{(0)}$

alors :

si la suite $\{\bar{x}^{(n)}\}$ converge, elle converge fortement vers :

$$\min \{ \bar{x} : \bar{x} = B \otimes \bar{x} \oplus c \wedge \bar{x} \geq \bar{x}^{(0)} \} \quad \text{si } \bar{x}^{(1)} \geq \bar{x}^{(0)}$$

$$\max \{ \bar{x} : \bar{x} = B \otimes \bar{x} \oplus c \wedge \bar{x} \leq \bar{x}^{(0)} \} \quad \text{si } \bar{x}^{(1)} \leq \bar{x}^{(0)}$$

On remarquera en effet que si les hypothèses de la proposition IV. 1 sont vérifiées avec $\bar{x}^{(0)} \leq \bar{x}^{(1)}$ par exemple, et si la suite $\{\bar{x}^{(n)}\}$ est convergente, donc fortement convergente, vers une limite \bar{x} , on peut affirmer :

$$\bar{x} = \inf \{ \bar{x} : \bar{x} = B \otimes \bar{x} \oplus c \wedge \bar{x} \geq \bar{x}^{(0)} \}$$

puisque :

Tout vecteur \bar{x} défini par :

$$\bar{x} = B \otimes \bar{x} \oplus c \wedge \bar{x} \geq \bar{x}^{(0)}$$

satisfaisant, relativement aux itérés $\bar{x}^{(n)}$:

$$\bar{x}^{(n)} \leq \bar{x} \quad \forall n$$

on a l'implication :

$$\bar{x}_i \leq \bar{x}_i \quad 1 \leq i \leq N$$

On peut d'ailleurs, étant donnés deux tableaux \otimes et \oplus définissant les opérations de multiplication et d'addition, caractériser l'ensemble des vecteurs \bar{x} de T^N satisfaisant $\bar{x} = B \otimes \bar{x} \oplus c$; par exemple :

Proposition IV. 2.

Si B est une matrice à éléments non négatifs, de rayon spectral inférieur à 1 ,

si les opérations \otimes_{ij} sont G-dirigées (resp. D-dirigées).

alors :

l'ensemble $E = \{\bar{x} : \bar{x} \in T^N \wedge \bar{x} = B \otimes \bar{x} \oplus c\}$ est un sup demi-treillis (resp. inf demi-treillis).

Exemple (grossier) : Soient $T = T_{10}^2$ et :

$$B = \begin{vmatrix} 0,99 & 0 \\ 0 & 0,99 \end{vmatrix} \quad c = \begin{vmatrix} 0,011 \\ 0,011 \end{vmatrix}$$

Si $\otimes_{ij} = \otimes_{\nabla}$, alors E contient les vecteurs \bar{x} , de composantes \bar{x}_1 et \bar{x}_2 définies par :

$$\bar{x}_1 \in T \wedge 0,11 \leq \bar{x}_1 \leq 1$$

$$\bar{x}_2 \in T \wedge 0,11 \leq \bar{x}_2 \leq 1$$

La démonstration de la proposition IV.2 s'établit par $(\otimes_{ij} = \otimes_{\nabla})$:

Si $\bar{x} = (\bar{x}_i)$ et $\bar{y} = (\bar{y}_i)$ sont 2 éléments de E ne vérifiant pas la relation d'ordre, soit $\bar{z}^{(0)} = (\bar{z}_i^{(0)})$ défini par :

$$\bar{z}_i^{(0)} = \max\{\bar{x}_i, \bar{y}_i\} \quad i = 1, \dots, N.$$

et soit pour $n \geq 1$: $\bar{z}^{(n)} = B \otimes_{\nabla} \bar{z}^{(n-1)} \oplus_{\nabla} c$.

$$\bar{z}^{(0)} \geq \bar{x} \wedge \bar{x} \in E \implies \bar{z}^{(1)} \geq \bar{x} \quad \text{et} \quad \bar{z}^{(1)} \geq \bar{z}^{(0)}$$

$$\bar{z}^{(0)} \geq \bar{y} \wedge \bar{y} \in E \implies \bar{z}^{(1)} \geq \bar{y}$$

La suite $\{\bar{z}^{(n)}\}$ est non décroissante.

Elle est, d'autre part, bornée par x car les opérations \otimes_{ij} étant G-dirigées, on a : $\bar{z}^{(1)} \leq B \bar{z}^{(0)} + c$ et :

$$\bar{z}^{(0)} \leq B \bar{z}^{(0)} + c \wedge \rho(B) < 1 \implies \bar{z}^{(0)} \leq x, \text{ ce qui implique } \bar{z}^{(n)} \leq x, \forall n$$

(En remarquant que $\bar{z}^{(n)} \leq x \implies \bar{z}^{(n+1)} \leq B \bar{z}^{(n)} + c \leq x$).

$\{\bar{z}^{(n)}\}$ converge fortement vers $\bar{t} = \min\{\bar{x} : \bar{x} \in E \wedge \bar{x} \geq \bar{z}^{(0)}\}$, qui est borne supérieure de la paire $\{\bar{x}, \bar{y}\}$.

CONSEQUENCES : Une conséquence des résultats précédents est :

Pour des opérations \otimes G-dirigées (resp. D-dirigées)
alors,

le choix du vecteur initial tel que :

$$\bar{x}^{(0)} \geq x \wedge \bar{x}^{(0)} \geq \bar{x}^{(1)} \quad (\text{resp. } \bar{x}^{(0)} \leq x \wedge \bar{x}^{(0)} \leq \bar{x}^{(1)})$$

assure, lorsque la convergence a lieu, la précision maximum.

En effet, dans l'hypothèse $\otimes_{ij} = \otimes_{\nabla} \forall i, j$ par exemple, tout élément \bar{x} de T^N vérifiant $\bar{x} = B \otimes_{\nabla} \bar{x} \oplus_{\nabla} c$, satisfait la relation d'ordre :

$$\bar{x} \leq x$$

Pour le choix énoncé du vecteur initial, la convergence des itérés $\bar{x}^{(n)}$ a lieu vers :

$$\max \{ \bar{x} : \bar{x} = B \otimes_{\nabla} \bar{x} \oplus_{\nabla} c \}$$

L'effet des erreurs d'arrondi, pour l'arithmétique considérée, est minimisé.

Nous énoncerons enfin le corollaire suivant de la proposition IV. 1 :

Corollaire IV. 1.

Si la matrice B du système linéaire $x = Bx + c$ est positive, et si le vecteur solution x est élément de T^N et vérifie $x = B \otimes_{\nabla} x \oplus_{\nabla} c$ (resp. $x = B \otimes_{\Delta} x \oplus_{\Delta} c$)

alors :

la suite $\{\bar{x}^{(n)}\}$ des itérés, calculés sur T par des opérations \otimes G-dirigées (resp. D-dirigées), à partir d'un vecteur initial $\bar{x}^{(0)}$ vérifiant $\bar{x}^{(0)} \geq x \wedge \bar{x}^{(0)} \geq \bar{x}^{(1)}$ (resp. $\bar{x}^{(0)} \leq x \wedge \bar{x}^{(0)} \leq \bar{x}^{(1)}$) converge fortement vers x .

EXEMPLES :

1) La résolution, par la méthode de Jacobi, du système linéaire $Ax = b$ du § I. 5, avec : $(*)_{ij} = (*)_{\nabla}$, $\bar{x}_i^{(0)} = 1,5 \quad \forall i$, conduit, par convergence forte, à la solution exacte du système : $x_i = 1 \quad \forall i$, en 160 itérations.

Soulignons l'effet d'accélération de convergence produit par les erreurs d'arrondi dès que $\bar{x}^{(1)} \leq \bar{x}^{(0)}$ $\vee (*)_{ij} = (*)_{\nabla}$, puisque la résolution en double précision du système, à partir du même vecteur initial, détermine avec le test d'arrêt suivant :

$$|\bar{x}_i^{(n+1)} - \bar{x}_i^{(n)}| \leq 10^{-7} \quad \forall i$$

le vecteur de composantes $(\bar{x}_{DP})_i = 1,000001$ en 190 itérations.

2) Soit le problème (\mathcal{P}) aux dérivées partielles :

$$(\mathcal{P}) \begin{cases} -\Delta u + u = 0 & \text{dans le carré (C) : } 0 < x < 1, 0 < y < 1 \\ u = ye^{-x} & \text{sur les bords du carré} \end{cases}$$

où $u(x,y)$ est une fonction des 2 variables x et y

et soient, avec $h = \frac{1}{N+1}$,

la grille de points de coordonnées $x_i = ih$, $y_j = jh$ ($i, j = 0, 1, \dots, N+1$)

et le problème (\mathcal{P}') "approchant" (\mathcal{P}) :

$$(\mathcal{P}') \begin{cases} u_{ij} = \frac{u_{i+1,j} + u_{i-1,j} + u_{i,j+1} + u_{i,j-1}}{4 + h^2} \\ u_{ij} = (ye^{-x}) & \text{sur les bords de (C)} \\ x = x_i, y = y_j \end{cases}$$

La résolution de (\mathcal{P}') par la méthode de Gauss-Seidel a, pour $N = 10$,

donné les résultats suivants où u_{DP} désigne la solution de (\mathcal{P}')

calculée en double précision, et u^1, u^2, u^3, u^4 désignent respectivement

les solutions obtenues par convergence pour :

$$\begin{aligned} (*) &= (*)_{\nabla} & \bar{u}^{(0)} < \bar{u}^{(1)} \\ (*) &= (*)_{\nabla} & \bar{u}^{(0)} > \bar{u}^{(1)} \\ (*) &= (*)_{\Delta} & \bar{u}^{(0)} < \bar{u}^{(1)} \\ (*) &= (*)_{\Delta} & \bar{u}^{(0)} > \bar{u}^{(1)} \end{aligned}$$

		u^1	u^2	u_{DP}
(extraits de tables)	i=5 j=1	0.5770592E-01	0.5770653E-01	0.5770702E-01
	=2	0.1154117E 00	0.1154129E 00	0.1154139D 00
	=3	0.1731173E 00	0.1731190E 00	0.1731206D 00
	=4	0.2308224E 00	0.2308246E 00	0.2308269D 00
	=5	0.2885270E 00	0.2885296E 00	0.2885327D 00
	=6	0.3462309E 00	0.3462340E 00	0.3462376D 00
	=7	0.4039344E 00	0.4039375E 00	0.4039415D 00
	=8	0.4616367E 00	0.4616398E 00	0.4616438D 00
	=9	0.5193383E 00	0.5193407E 00	0.5193443D 00
	=10	0.5770385E 00	0.5770399E 00	0.5770420D 00

u^3	u^4	ye^{-x}
0.5770778E-01	0.5770835E-01	0.5770330E-01
0.1154155E 00	0.1154166E 00	0.1154066E 00
0.1731229E 00	0.1731244E 00	0.1731099E 00
0.2308300E 00	0.2308320E 00	0.2308131E 00
0.2885365E 00	0.2885389E 00	0.2885165E 00
0.3462422E 00	0.3462446E 00	0.3462198E 00
0.4039464E 00	0.4039490E 00	0.4039230E 00
0.4616485E 00	0.4616511E 00	0.4616264E 00
0.5193481E 00	0.5193505E 00	0.5193297E 00
0.5770445E 00	0.5770457E 00	0.5770330E 00

On remarquera que dans le cas étudié, le vecteur u^1 , bien qu'entaché d'une plus forte erreur que u^2 par rapport à u_{DP} , est une meilleure approximation de la solution exacte u de (\mathcal{P}) :

L'erreur commise en substituant (\mathcal{P}') à (\mathcal{P}) et l'erreur due aux calculs, se compensent en effet davantage dans le 1er cas, la relation $u \leq u_{DP}$ résultant du signe de $u_{\frac{4}{x}}^{(4)}(x, y)$ dans le carré unité.

Lorsque la matrice B de l'itération linéaire est négative :

Le choix d'un vecteur initial $x^{(0)}$ tel que $x^{(0)} \leq x^{(2)}$
(resp. $x^{(0)} \geq x^{(2)}$) implique :

$$x \geq x^{(0)}$$

$\{ x^{(2k)} \}$ monotone non décroissante

$\{ x^{(2k+1)} \}$ monotone non croissante

Le calcul de la suite $\{ \bar{x}^{(n)} \}$:

par des opérations D-dirigées pour les itérés d'indice pair

par des opérations G-dirigées pour les itérés d'indice impair

et l'arrêt des calculs dès que la relation d'ordre entre les éléments

des 2 suites $\{ \bar{x}^{(2k)} \}$ et $\{ \bar{x}^{(2k+1)} \}$ n'est plus vérifiée, permet :

- l'accélération de la convergence
- l'obtention d'une bonne précision.

III - RESOLUTION ITERATIVE EN ARITHMETIQUE A VIRGULE FLOTTANTE DE L'EQUATION LINEAIRE

$x = ax + c$

Nous supposons a et c éléments de T_b^s dans l'équation :

$$x = ax + c \quad (IV. 13)$$

L'étude de cette équation, inintéressante du point de vue pratique, permet de souligner l'influence des erreurs d'arrondi suivant les valeurs de a et d'établir des comparaisons entre différentes arithmétiques, relativement à la résolution numérique de (IV. 13) par itérations.

Nous énoncerons trois théorèmes :

Proposition IV. 3.

Si les opérations \otimes et \oplus sont correctes, alors :

si $-1 < a \leq 0$ $\text{Card} \{ \bar{x} \mid \bar{x} = a \otimes \bar{x} \oplus c \} \leq 2$

si $0 < a \leq \frac{1}{2}$ $\text{Card} \{ \bar{x} \mid \bar{x} = a \otimes \bar{x} \oplus c \} \leq 4$, sauf peut-être

si la mantisse de c est de valeur absolue supérieure à $1 - 2b^{1-s}$
(condition C_m).

En effet, si les opérations \otimes et \oplus sont correctes, et si \bar{x} est un élément de T vérifiant $\bar{x} = a \otimes \bar{x} \oplus c$, alors :

$$-1 < a \leq \frac{1}{2} \implies |x - \bar{x}| < \frac{1}{1-a} \epsilon(\bar{x}) \quad (IV. 14), \text{ sauf peut-être}$$

si C_m a lieu

$$\frac{1}{2} < a < 1 \implies |x - \bar{x}| < \frac{2}{1-a} \epsilon(\bar{x}) \quad (IV. 15)$$

(IV. 15) vient immédiatement quel que soit a satisfaisant $|a| < 1$ par l'écriture :

$$a \bar{x} = a \otimes \bar{x} + e_1 \quad \wedge \quad |e_1| < \epsilon(a \bar{x}) \leq \epsilon(\bar{x})$$

$$a \otimes \bar{x} + c = \bar{x} + e_2 \quad \wedge \quad |e_2| < \epsilon(\bar{x})$$

entraînant avec $e = e_1 + e_2$:

$$|x - \bar{x}| = \frac{|e|}{1-a} < \frac{2}{1-a} \varepsilon(\bar{x})$$

Ecrivons :

$$|a\bar{x}| \leq \frac{|a|}{1-a} [|c| + |e|]$$

Si $-1 < a \leq 0$, on aura :

$$\frac{|a|}{1-a} \leq \frac{1}{2}$$

$$|x| \leq |c|$$

$$|a\bar{x}| < \frac{1}{2} [|c| + 2b \varepsilon(c)] < |c|$$

Si $0 < a \leq \frac{1}{2}$, on aura :

$$|a\bar{x}| < |c| + (2b - 1) \varepsilon(c)$$

assurant

$\varepsilon(a \otimes \bar{x}) \leq \varepsilon(c)$ si C_m n'a pas lieu.

Les 2 cas ci-dessus permettant d'écrire :

$$\varepsilon(a \otimes \bar{x}) \leq \varepsilon(c)$$

on restreindra la borne de e_2 par :

$$|e_2| \leq \varepsilon(\bar{x}) - \varepsilon(a \otimes \bar{x})$$

et

$$|e| < \varepsilon(\bar{x})$$

implique (IV. 14)

La proposition IV. 3. se déduit de (IV. 14) par :

$$-1 < a \leq 0 \implies \bar{x} = \nabla x \vee \Delta x$$

$$0 < a \leq \frac{1}{2} \implies \bar{x} = (\nabla x)' \vee \nabla x \vee \Delta x \vee (\Delta x)'' \quad (\text{sauf peut-être cas } C_m :$$

ainsi pour $x = 0,50x + 0,92$ et $T = T_{10}^2$, $\{\bar{x}\} = \{1,7; 1,8; 1,9; 2,0; 2,1\}$)

Remarque :

Le résultat $\bar{x} = \nabla x \nabla \Delta x$ lorsque $-1 < a < 0$, n'implique évidemment pas que ∇x (resp. Δx) vérifie $\nabla x = a \otimes \nabla x \oplus c$ (resp. $\Delta x = a \otimes \Delta x \oplus c$) pour l'un au moins des choix "corrects" d'opérations : $(\otimes_{\nabla}, \oplus_{\nabla})$, $(\otimes_{\nabla}, \oplus_{\Delta})$, $(\otimes_{\Delta}, \oplus_{\nabla})$, $(\otimes_{\Delta}, \oplus_{\Delta})$, ni que la suite des itérés $\{\bar{x}^{(n)}\}$ converge -fortement ou faiblement- vers des valeurs satisfaisantes.

Ainsi pour $x = -0,99x + 0,43$ et $T = T_{10}^2$:

$\nabla x = 0,21$ ne vérifie pas $\nabla x = a \otimes \nabla x \oplus c$

La suite $\{\bar{x}^{(n)}\}$, obtenue à partir de $\bar{x}^{(0)} = 0,21$ par $\bar{x}^{(n+1)} = a \otimes_{\Delta} \bar{x}^{(n)} \oplus_{\Delta} c$ si n est pair, $\bar{x}^{(n+1)} = a \otimes_{\nabla} \bar{x}^{(n)} \oplus_{\nabla} c$ si n est impair, converge faiblement en cyclant sur 0 et 0,43.

On a cependant :

Proposition IV. 4.

Si $-1 < a < 0$, alors quel que soit $\bar{x}^{(0)} \neq x$, la suite $\{\bar{x}^{(n)}\}$ des itérés définis par :

$$\bar{x}^{(n+1)} = a \otimes_{\nabla} \bar{x}^{(n)} \oplus_{\nabla} c \quad \text{si } n \text{ pair quand } \bar{x}^{(0)} < x, \\ \text{si } n \text{ impair quand } \bar{x}^{(0)} > x$$

$$\bar{x}^{(n+1)} = a \otimes_{\Delta} \bar{x}^{(n)} \oplus_{\Delta} c \quad \text{si } n \text{ impair quand } \bar{x}^{(0)} < x, \\ \text{si } n \text{ pair quand } \bar{x}^{(0)} > x,$$

détermine avec le test d'arrêt suivant :

$$\bar{x}^{(P+1)} - \bar{x}^{(P)} \leq 0 \quad \text{si } P \text{ pair } \wedge \bar{x}^{(0)} < x \quad \text{ou } P \text{ impair } \wedge \bar{x}^{(0)} > x$$

$$\bar{x}^{(P+1)} - \bar{x}^{(P)} \geq 0 \quad \text{si } P \text{ impair } \wedge \bar{x}^{(0)} < x \quad \text{ou } P \text{ pair } \wedge \bar{x}^{(0)} > x$$

l'un des deux nombres ∇x ou Δx (égal à $\bar{x}^{(P)}$).

Nous supposons, afin de ne pas alourdir la démonstration, que les éléments de T , au voisinage de x , vérifient une condition de régularité relative à leur espacement, c'est-à-dire que nous supposons que la mantisse $m(x)$ de x vérifie :

$$\frac{1}{b} + b^{1-s} \leq |m(x)| < 1 - b^{1-s} \quad (\text{condition } c_x)$$

A chaque itéré $\bar{x}^{(n+1)}$, associons l'itéré $x^{(n+1)}$, calculé sur \mathbb{R} à partir de $\bar{x}^{(n)}$, par :

$$x^{(n+1)} = a \bar{x}^{(n)} + c$$

la disposition de $x^{(n+1)}$ par rapport à x , à partir de la position de $\bar{x}^{(n)}$ par rapport à x , est :

$$\bar{x}^{(n)} < x \implies x^{(n+1)} > x$$

$$\bar{x}^{(n)} > x \implies x^{(n+1)} < x$$

avec de plus par $|a| < 1$:

$$|x^{(n+1)} - x| < |\bar{x}^{(n)} - x|$$

Etudions l'écart entre $x^{(n+1)}$ et $\bar{x}^{(n+1)}$, avec les notations :

$$a \bar{x}^{(n)} = a \otimes \bar{x}^{(n)} + e_1$$

$$a \otimes \bar{x}^{(n)} + c = \bar{x}^{(n+1)} + e_2$$

$$e = e_1 + e_2$$

$$x^{(n+1)} = \bar{x}^{(n+1)} + e$$

Plusieurs cas peuvent se présenter :

i) $a \bar{x}^{(n)}$ et c sont de même signe

lorsque $\varepsilon(a \otimes \bar{x}^{(n)}) \leq \varepsilon(c)$ alors $|e| < \varepsilon(\bar{x}^{(n+1)})$

lorsque $\varepsilon(a \otimes \bar{x}^{(n)}) > \varepsilon(c)$ alors $|e| < 2\varepsilon(\bar{x}^{(n+1)})$, ce cas impliquant

$\varepsilon(\bar{x}^{(n+1)}) > \varepsilon(x)$ par : $\varepsilon(c) \geq \varepsilon(x)$.

ii) a $\bar{x}^{(n)}$ et c sont de signes contraires et une cancellation ne se produit pas lors de la formation de $\bar{x}^{(n+1)}$

$|a \otimes \bar{x}^{(n)}| > |c|$ implique que $x^{(n+1)}$ et $\bar{x}^{(n+1)}$ n'ont pas le signe de x .

$|a \otimes \bar{x}^{(n)}| \leq |c|$ assure $|e| < \varepsilon(\bar{x}^{(n+1)})$

iii) a $\bar{x}^{(n)}$ et c sont de signes contraires et une cancellation se produit lors de la formation de $\bar{x}^{(n+1)}$

$|a \otimes \bar{x}^{(n)}| > |c|$ entraîne la même implication que ci-dessus

$|a \otimes \bar{x}^{(n)}| \leq |c|$ assure $|e| < \varepsilon(c) \leq b \varepsilon(x)$ mais :

lorsque $a < -1 + \frac{1}{b}$, la condition c_x entraîne $|x^{(n+1)}| < x - b \varepsilon(x)$

lorsque $a \geq -1 + \frac{1}{b}$, de $|c| < (2 - \frac{1}{b}) |x|$ et $|a \otimes \bar{x}^{(n)}| > |x|$ on tire

$|e - a \otimes \bar{x}^{(n)}| \leq (1 - \frac{1}{b}) |x|$

En conclusion, l'erreur e est inférieure à $\varepsilon(\bar{x}^{(n+1)})$ dès que $x^{(n+1)}$ est "suffisamment près" de x .

Il résulte de l'étude précédente que :

$\bar{x}^{(n+1)}$ a la même disposition par rapport à x que $x^{(n+1)}$

$|\bar{x}^{(n+1)} - x| < |\bar{x}^{(n)} - x|$

tant que $x^{(n+1)}$ vérifie :

$x^{(n+1)} \leq \forall x \vee x^{(n+1)} \geq \Delta x$

Si $x \in T$, $\bar{x}^{(n)}$ converge fortement vers x .

Si $x \notin T$, avec $\bar{x}^{(0)} < x$ par exemple, il existe un indice $P = 2Q$ (ou $P = 2Q + 1$) satisfaisant : $\forall x < x^{(2Q)} < x$ (ou $x < x^{(2Q+1)} < \Delta x$) .

Considérons le 1^{er} cas : $\bar{x}^{(2Q)} = \Delta x$

le test d'arrêt est alors réalisé :

au rang P si $\bar{x}^{(P-1)} = \Delta x$

au rang P + 1 sinon.

Remarques :

- 1) la suite $\{\bar{x}^{(n)}\}$ converge plus rapidement que la suite $\{x^{(n)}\}$.
- 2) L'étude de l'erreur e établit que l'élément \bar{x} de T , le plus proche de x -ou les deux s'il y a équidistance- vérifie certainement $\bar{x} = a \otimes \bar{x} \oplus c$ pour $\otimes = \otimes_{\Delta}$ si $\bar{x} = \Delta x$, $\otimes = \otimes_{\nabla}$ si $\bar{x} = \nabla x$.

Nous donnerons aussi à titre d'exemple, la :

Proposition IV. 5.

Si $a = 1 - b^{-s}$, alors

- 1) si les opérations \otimes et \oplus sont correctes et I-dirigées, $\forall \bar{x} \in T_b^s \cap]b^{e(x)-2}, b^{e(x)-1}] \times \text{signe}(x)$ si le 1er chiffre de mantisse de x est 1, pour $\bar{x} = b^{e(x)-1} \times \text{signe}(x)$ sinon, on a:

$$\bar{x} = a \otimes \bar{x} \oplus c$$
- 2) si les opérations \otimes et \oplus sont correctes et E-dirigées, $\forall \bar{x}^{(0)} \in [0, b^{e(x)}]$ la suite $\{\bar{x}^{(n)}\}$ converge vers $b^{e(x)} \times \text{signe}(x)$. (voir r) $\times \text{signe}(x)$

$$x = b^s c$$

Nous considérons $c > 0$.

1) s'établit par :

$$a \otimes_{\nabla} \bar{x} = (\bar{x})'$$

$$(\bar{x})' \oplus_{\nabla} c = \bar{x} \quad \text{dans les conditions indiquées.}$$

2) s'établit par :

$$a \otimes_{\Delta} y = y \quad \text{sauf si} \quad a \otimes_{\Delta} y = a y$$

(r): Si la condition (C) du ch.II n'est pas vérifiée par l'opération \oplus si $\text{mantisse}(x) < \frac{1}{b} + \frac{1}{b}$

il convient de modifier légèrement la proposition.

IV - CONCLUSION

Il est certain que la possibilité de disposer, sur un calculateur, des 2 opérations approchées $\textcircled{*}_{\nabla}$ et $\textcircled{*}_{\Delta}$ de chacune des opérations élémentaires $*$ définies sur \mathbb{R} , est un élément important pour le contrôle des erreurs d'arrondi : Ce fait est particulièrement vrai dans le cas de méthodes itératives de résolution, les opérateurs étant fréquemment monotones.

D'un point de vue différent, nous mentionnerons dans ce chapitre, l'étude effectuée par M. TIENARI [37] pour le contrôle de la longueur de mantisse, dans les méthodes itératives en arithmétique à virgule flottante.

CHAPITRE V

STABILITE NUMERIQUE
DE DIFFERENTS SCHEMAS DE SOMMATION.

CALCUL DE SOMMES DE SERIES

INTRODUCTION

Les problèmes étudiés dans ce chapitre, relatifs à l'algorithme de la sommation, sont des problèmes de "stratégie" :

Les opérations \otimes définies sur les sous-ensembles finis de \mathbb{R} ne sont pas généralement associatives. La non-associativité de l'opération \oplus pose, pour la sommation, d'éléments donnés, par l'arithmétique approchée, le problème de la disposition des opérations en vue de l'obtention de la meilleure précision.

Les résultats qui suivent sont des théorèmes de comparaison de précision entre les algorithmes les plus usuels réalisant la sommation de N nombres par $N - 1$ additions (coût minimum).

Les théorèmes sont établis pour des sous-ensembles de \mathbb{R} de type flottant, ou des systèmes de nombres à virgule flottante.

Une application particulière des résultats obtenus, est leur utilisation au calcul de sommes de séries numériques.

I - ETUDE COMPARATIVE DE SCHEMAS DE SOMMATION [34]

I. 1. PREMIER THEOREME

Une conséquence directe du lemme I. 2 du chapitre I , relatif aux sous-ensembles de \mathbb{R} de type flottant, est la proposition suivante :

Proposition V. 1.

Si S est un sous-ensemble de \mathbb{R} de type flottant, si l'addition sur S , notée \oplus , est correcte, si x_1, x_2, \dots, x_N sont N éléments de S^+ vérifiant les inégalités :

$$\varepsilon \left(\oplus \sum_{i=1}^p x_i \right) \leq \varepsilon (x_{p+1}) \quad p = 2, \dots, N-1 \quad (V. 1)$$

alors, si l'on désigne par :

$$\sigma = \sum_{i=1}^N x_i \quad \text{la somme des } x_i \text{ sur } \mathbb{R}$$

$$\bar{\sigma}_S = \oplus \sum_{i=1}^N x_i \quad \text{la somme des } x_i \text{ sur } S \text{ par } \oplus \text{ dans l'ordre séquentiel,}$$

il vient :

$$\bar{\sigma}_S = \nabla \sigma \vee \Delta \sigma$$

Désignant par e_i ($i = 1, \dots, N - 1$) l'erreur élémentaire commise lors de la i ème addition dans le calcul de $\bar{\sigma}_S$, on peut affirmer par le lemme I. 2 (chapitre I) :

$$|e_i| \leq \varepsilon \left(\bigoplus_{j=1}^{i+1} x_j \right) - \varepsilon \left(\bigoplus_{j=1}^i x_j \right) \quad i=1, \dots, N-1$$

Par la majoration :

$$|\sigma - \bar{\sigma}_s| \leq \sum_{i=1}^{N-1} |e_i| < \varepsilon (\bar{\sigma}_s)$$

on déduit alors la proposition.

- Pour conclure à l'égalité $\bar{\sigma}_s = \nabla \sigma \vee \Delta \sigma$ dans le cas particulier où $\bar{\sigma}_s$ vérifie $\varepsilon ((\bar{\sigma}_s)') \neq \varepsilon (\bar{\sigma}_s)$, on remarquera que :

a) si $\left(\bigoplus_{i=1}^{N-1} x_i \right) + x_N < \bar{\sigma}_s$ alors :

$$|\sigma - \bar{\sigma}_s| < \varepsilon \left(\left(\bigoplus_{i=1}^{N-1} x_i \right) + x_N \right) \leq \varepsilon ((\bar{\sigma}_s)')$$

b) si $\left(\bigoplus_{i=1}^{N-1} x_i \right) + x_N \geq \bar{\sigma}_s$ alors $e_{N-1} \geq 0$.

$$\sigma \geq \bar{\sigma}_s \implies \bar{\sigma}_s = \nabla \sigma$$

$$\sigma < \bar{\sigma}_s \implies \bar{\sigma}_s - \sigma < \sum_{i=1}^{N-2} e_i < \varepsilon \left(\bigoplus_{i=1}^{N-1} x_i \right) \leq \varepsilon ((\bar{\sigma}_s)')$$

Remarque

Dans les hypothèses suivantes :

S est un sous-ensemble de \mathbb{R} de type flottant,

l'addition sur S, notée \bigoplus , est induite de l'addition sur \mathbb{R} par ∇ (resp. Δ),

x_1, x_2, \dots, x_N sont N éléments de S^+ vérifiant (V. 1),

si σ et $\bar{\sigma}_S$ sont définies comme précédemment et si l'on note $\bar{\sigma}$ une somme quelconque des x_i sur S par \oplus , $\bar{\sigma}_S$ est une meilleure approximation de σ dans l'ensemble $\Sigma = \{ \bar{\sigma} \}$:

$$| \sigma - \bar{\sigma}_S | = \text{Min}_{\bar{\sigma} \in \Sigma} | \sigma - \bar{\sigma} |$$

L'ordre séquentiel croissant est optimal pour la sommation des suites $\{x_i\}$ à croissance rapide (vérifiant V. 1) .

Exemple : Sommation par \oplus_{∇} des nombres suivants écrits en base 10 avec 2 chiffres de mantisse : 0,030 ; 0,10 ; 0,20 ; 0,30 ; 0,67 ; 1,0 ; 1,7 ; 14.
 $\sigma = 18$

Somme calculée par sommation dans l'ordre séquentiel croissant : 18

Somme calculée par sommation dans l'ordre séquentiel décroissant : 16

Somme calculée par sommation suivant un schéma dichotomique : 16

Le corollaire suivant de la proposition V. 1, relatif à la sommation de trois nombres, sera utile pour la suite :

Corollaire V. 1.

Si S est un sous-ensemble de \mathbb{R} de type flottant, pour tout triplet $(x, y, z) \in (S^+)^3$ et vérifiant l'inégalité :

$$\varepsilon(x + y) \leq \varepsilon(z) \tag{V. 2}$$

on a :

$$(x \oplus_{\nabla} y) \oplus_{\nabla} z \geq \max \{ (x \oplus_{\nabla} z) \oplus_{\nabla} y, (y \oplus_{\nabla} z) \oplus_{\nabla} x \}$$

$$(x \oplus_{\Delta} y) \oplus_{\Delta} z \leq \min \{ (x \oplus_{\Delta} z) \oplus_{\Delta} y, (y \oplus_{\Delta} z) \oplus_{\Delta} x \}$$

I. 2. THEOREMES GENERAUX DE COMPARAISON DES SOMMES SEQUENTIELLE ET DICHOTOMIQUE

I. 2. 1. Sommation d'éléments à valeurs positives et décroissantes

Lorsque les éléments à sommer sont classés dans l'ordre décroissant, on peut énoncer (l'écriture ci-dessous supposant que N est une puissance de 2 : $N = 2^k$, afin de ne pas compliquer inutilement les notations) :

Proposition V. 2.

Si S est un sous-ensemble de \mathbb{R} de type flottant, si l'addition sur S , notée \oplus , est induite de l'addition sur \mathbb{R} par ∇ (resp. Δ),

si x_1, x_2, \dots, x_N sont N éléments de S^+ rangés dans l'ordre décroissant,

$$\text{alors, les sommes } \sigma = \sum_{i=1}^N x_i, \bar{\sigma}_s = \oplus \sum_{i=1}^N x_i \text{ et } \bar{\sigma}_d = \sigma_{1k} \text{ où}$$

σ_{1k} est défini par la récurrence suivante :

$$\sigma_{i0} = x_i \quad i = 1, \dots, N$$

$$\sigma_{ij} = \sigma_{2i-1, j-1} \oplus \sigma_{2i, j-1} \quad i = 1, \dots, N/2^j$$

$$j = 1, \dots, k$$

vérifient :

$$|\sigma - \bar{\sigma}_d| \leq |\sigma - \bar{\sigma}_s|$$

Démonstration :

L'algorithme de sommation définissant $\bar{\sigma}_d$ est le schéma de sommation dichotomique .

Etablissons la proposition V. 2 dans le cas où $\oplus = \oplus_{\nabla}$

Il suffit de montrer $\bar{\sigma}_d \geq \bar{\sigma}_s$.

De l'inégalité :

$$\sigma_{2i-1, 0} + \sigma_{2i, 0} \geq \sigma_{2i+1, 0} + \sigma_{2i+2, 0}$$

on déduit par la monotonie de l'arrondi

$$\sigma_{i, 1} \geq \sigma_{i+1, 1} \quad 1 \leq i < N/2$$

et pour tout j :

$$\sigma_{i, j} \geq \sigma_{i+1, j} \quad 1 \leq i < N/2^j$$

Montrons alors l'inégalité suivante :

$$\oplus_{\nabla} \sum_{i=1}^{N/2^{j+1}} \sigma_{i, j+1} \geq \oplus_{\nabla} \sum_{i=1}^{N/2^j} \sigma_{i, j} \quad \forall j \quad (V. 3)$$

Ecrivons :

$$\begin{aligned} \oplus_{\nabla} \sum_i \sigma_{i, j+1} &= \sigma_{1, j+1} \oplus_{\nabla} (\sigma_{3, j} \oplus_{\nabla} \sigma_{4, j}) \oplus_{\nabla} (\sigma_{5, j} \oplus_{\nabla} \sigma_{6, j}) \oplus_{\nabla} \dots \\ &\quad \oplus_{\nabla} (\sigma_{N/2^{j-1}, j} \oplus_{\nabla} \sigma_{N/2^j, j}) \end{aligned}$$

En remarquant :

$$\varepsilon(\sigma_{3, j} + \sigma_{4, j}) \leq \varepsilon(\sigma_{1, j+1})$$

il vient, par le corollaire V. 1,

$$\sigma_{1, j+1} \oplus_{\nabla} (\sigma_{3, j} \oplus_{\nabla} \sigma_{4, j}) \geq \sigma_{1, j} \oplus_{\nabla} \sigma_{2, j} \oplus_{\nabla} \sigma_{3, j} \oplus_{\nabla} \sigma_{4, j} = A$$

puis

$$A \oplus_{\nabla} (\sigma_{5, j} \oplus_{\nabla} \sigma_{6, j}) \geq A \oplus_{\nabla} \sigma_{5, j} \oplus_{\nabla} \sigma_{6, j}$$

et ainsi de suite, ce qui entraîne, en tenant compte de la propriété de monotonie de l'arrondi, l'inégalité (V. 3).

Par l'application répétée de l'inégalité (V. 3), on déduit :

$$\sigma_{1k} \geq \oplus_{\nabla} \sum_i \sigma_{i,j} \quad \forall j$$

et en particulier :

$$\bar{\sigma}_d \geq \bar{\sigma}_s$$

Remarques :

Lorsque N n'est pas une puissance de 2, pour certains j , le nombre d'éléments de la couche j est du type $2\lambda - 1$, impair, définissant $\sigma_{\lambda, j+1} = \sigma_{2\lambda - 1, j}$.

Les hypothèses de la proposition V. 2 ne sont pas suffisamment fortes pour assurer l'optimalité d'un schéma de sommation particulier. A titre d'exemple, la sommation des nombres décimaux, de 2 chiffres de mantisse :

9,4 8,5 7,5 6,3 5,8 4,7 4,5 4,4

donne les résultats suivants :

$$\sigma = 51,1$$

$$\text{si } \oplus = \oplus_{\nabla} :$$

$$x_1 \oplus (x_2 \oplus x_3) \oplus (x_4 \oplus x_5) \oplus (x_6 \oplus x_7) \oplus x_8 = 50$$

$$\bar{\sigma}_d = 48 \quad \bar{\sigma}_s = 47 \quad \oplus \sum_{i=N}^1 x_i = 48$$

$$\text{si } \oplus = \oplus_{\Delta} :$$

$$\bar{\sigma}_d = 52 \quad \bar{\sigma}_s = 54 \quad \oplus \sum_{i=N}^1 x_i = 54$$

I. 2. 2. Sommation d'éléments à valeurs positives quelconques

Etablissons le lemme préliminaire suivant :

Lemme V. 1.

Si S est un sous-ensemble de \mathbb{R} de type flottant, pour tout triplet $(x, y, z) \in (S^+)^3$, les inégalités suivantes sont vérifiées :

$$(x)' \oplus_{\nabla} (y \oplus_{\nabla} z) \geq ((x \oplus_{\nabla} y) \oplus_{\nabla} z)' \quad (V. 4)$$

$$(x)'' \oplus_{\Delta} (y \oplus_{\Delta} z) \leq ((x \oplus_{\Delta} y) \oplus_{\Delta} z)'' \quad (V. 4')$$

Démonstration :

Démontrons (V. 4) .

Il suffit d'établir $(x)' \oplus_{\nabla} (y \oplus_{\nabla} z) \geq (\nabla(x+y+z))'$.

L'égalité :

$$x + y + z - [(x)' \oplus_{\nabla} (y \oplus_{\nabla} z)] = \varepsilon((x)') + e_1 + e_2$$

où :

$$e_1 = y + z - (y \oplus_{\nabla} z)$$

$$e_2 = (x)' + (y \oplus_{\nabla} z) - [(x)' \oplus_{\nabla} (y \oplus_{\nabla} z)]$$

permet, avec les inégalités déduites du lemme I. 2 :

$$|e_1| < \varepsilon(y+z)$$

$$|e_2| < \varepsilon[(x)' + (y \oplus_{\nabla} z)] - \varepsilon(y \oplus_{\nabla} z) \quad \text{si } \varepsilon(y \oplus_{\nabla} z) \leq \varepsilon((x)')$$

$$|e_2| < \varepsilon[(x)' + (y \oplus_{\nabla} z)] - \varepsilon((x)') \quad \text{si } \varepsilon(y \oplus_{\nabla} z) > \varepsilon((x)')$$

d'affirmer dans les 2 cas, avec la notation $\tilde{\sigma} = (x)' \oplus_{\nabla} (y \oplus_{\nabla} z)$

$$(x + y + z) - \tilde{\sigma} < 2 \varepsilon (\tilde{\sigma})$$

donc :

$$\tilde{\sigma} \geq (\nabla (x + y + z))'$$

Il vient alors :

Proposition V. 3.

Si S est un sous-ensemble de \mathbb{R} de type flottant,
 si l'addition sur S , notée \oplus , est induite de l'addition sur \mathbb{R} par ∇
 (resp. Δ),
 alors, étant donnés N éléments x_i ($1 \leq i \leq N$) de S^+ , les deux sommes
 définies sur S par :

$$\bar{\sigma}_s = \oplus \sum_{i=1}^N x_i$$

$$\bar{\sigma} = \left| \begin{array}{l} \oplus \sum_{j=1}^{N/2} (x_{2j-1} \oplus x_{2j}) \end{array} \right. \quad \text{si } N \text{ pair}$$

$$\left. \begin{array}{l} \oplus \sum_{j=1}^{\lfloor N/2 \rfloor} (x_{2j-1} \oplus x_{2j}) \end{array} \right\} \oplus x_N \quad \text{si } N \text{ impair}$$

vérifient :

$$\bar{\sigma} \geq (\bar{\sigma}_s)' \quad \text{si } \oplus = \oplus_{\nabla}$$

$$\bar{\sigma} \leq (\bar{\sigma}_s)'' \quad \text{si } \oplus = \oplus_{\Delta}$$

Démonstration :

La démonstration est immédiate par récurrence.

Prenons le cas où $\oplus = \oplus_{\nabla}$ et posons :

$$\bar{\sigma}_s^{2i} = \oplus_{\nabla} \sum_{j=1}^{2i} x_j$$

$$\bar{\sigma}^{2i} = \oplus_{\nabla} \sum_{j=1}^i (x_{2j-1} \oplus_{\nabla} x_{2j})$$

- Pour $i = 2$, la propriété $\bar{\sigma}^{2i} \geq (\bar{\sigma}_s^{2i})'$ est bien vraie.

- Pour i fixé ≥ 2 , supposons $\bar{\sigma}^{2i} \geq (\bar{\sigma}_s^{2i})'$,

$$\bar{\sigma}^{2i+2} = \bar{\sigma}^{2i} \oplus_{\nabla} (x_{2i+1} \oplus_{\nabla} x_{2i+2})$$

$$\bar{\sigma}_s^{2i+2} = \bar{\sigma}_s^{2i} \oplus_{\nabla} x_{2i+1} \oplus_{\nabla} x_{2i+2}$$

Si $\bar{\sigma}^{2i} \geq \bar{\sigma}_s^{2i}$, on a évidemment, S étant à intervalles non décroissants et l'arrondi monotone, $\bar{\sigma}^{2i+2} \geq (\bar{\sigma}_s^{2i+2})'$.

Si $\bar{\sigma}^{2i} = (\bar{\sigma}_s^{2i})'$, la propriété résulte du lemme précédent.

I. 3. THEOREMES RELATIFS AUX SYSTEMES DE NOMBRES A VIRGULE FLOTTANTE

Ce paragraphe établit des théorèmes plus particuliers de comparaison d'algorithmes de sommation, pour des nombres x_i , écrits en virgule flottante et vérifiant des hypothèses assez restrictives.

Relativement aux systèmes de nombres à virgule flottante de base 2, il vient :

Proposition V. 4

N éléments x_i ($1 \leq i \leq N$) de $(T_2^s)^+$ étant donnés,
si l'addition sur T_2^s est induite de l'addition sur \mathbb{R} par ∇ (resp. Δ),
si les N éléments x_i ont même exposant

alors, les sommes $\bar{\sigma}_s$ et $\bar{\sigma}_d$ calculées respectivement sur T_2^s
par les schémas de sommation séquentiel et dichotomique, vérifient
l'inégalité :

$$|\sigma - \bar{\sigma}_d| \ll |\sigma - \bar{\sigma}_s|$$

La proposition précédente se démontre de façon analogue à la proposition V. 2 :

Si $\oplus = \oplus_{\nabla}$, on établira l'inégalité $\bar{\sigma}_d \geq \bar{\sigma}_s$, en remarquant que
si p désigne l'exposant commun des x_i ,

$$e(x_i) = p \quad i = 1, \dots, N$$

par les hypothèses, les éléments $\sigma_{i,j}$ ont, pour j fixé, tous le même
exposant :

$$e(\sigma_{i,j}) = p + j \quad j = 0, \dots, k$$

les inégalités du type :

$$\varepsilon(\sigma_{3,j} + \sigma_{4,j}) \leq \varepsilon(\sigma_{1,j} + 1)$$

sont alors assurées.

Pour un système de nombres à virgule flottante de base b
quelconque, la démonstration de la proposition V. 4 reste valable, si l'on
remplace la sommation dichotomique par la sommation b à b d'éléments
de T_b^s ayant même exposant.

(Le schéma de sommation b à b se définit pour $N = b^k$ par la récurrence :

$$\sigma_{i0} = x_i \quad i = 1, \dots, N$$

$$\sigma_{ij} = \oplus \sum_{l=b(i-1)+1}^{bi} \sigma_{l,j-1} \quad \begin{array}{l} i = 1, \dots, b^{k-j} \\ j = 1, \dots, k \end{array}$$

Proposition V. 4'.

Pour un système de nombres flottants de base b , lorsque l'addition flottante est induite de l'addition sur \mathbb{R} par ∇ (resp. Δ), la sommation b à b de flottants de même signe et de même exposant, calcule une meilleure approximation de la somme de ces nombres sur \mathbb{R} que la sommation séquentielle.

En pratique, cette méthode de sommation pour une base $b \neq 2$, n'est pas usitée.

Si, dans le système T_b^s , $b \neq 2$, on compose la sommation par le schéma dichotomique à la sommation b à b , de flottants de même exposant, la première sommation se comporte moins "régulièrement", par rapport à la sommation séquentielle, que la seconde ; toutefois, dans la plupart des cas, le schéma dichotomique est meilleur que la sommation b à b .

Un cas d'application courante (surtout lorsque p , ci-dessous, est nul) où il est possible de classer ces deux procédés, indépendamment des x_i , est donné par la :

Proposition V. 5.

Si l'addition sur $T = T_b^s$ ($b > 3$) est induite de l'addition sur \mathbb{R} par ∇ (resp. Δ), alors, étant donnés les éléments x_i de T vérifiant :

$$\frac{1}{2} b^p \leq x_i < b^p \quad \forall i \quad (p \text{ entier relatif fixé}),$$

la sommation par le schéma dichotomique de ces nombres donne une meilleure approximation de $\sum_i x_i$, que la sommation b à b .

La démonstration de la proposition V. 5, pour $\oplus = \oplus_{\nabla}$, lorsque b est une puissance de 2, $b = 2^\alpha$ ($\alpha > 1$), s'établit en remarquant l'inégalité suivante :

$$\sigma_{i,j}^{(b)} \leq \sigma_{i,\alpha j}^{(2)} \quad \forall i$$

où les notations $\sigma_{ij}^{(b)}$ et $\sigma_{ij}^{(2)}$ désignent les éléments des couches j dans les calculs par \oplus suivant les schémas de sommation respectivement b à b et 2 à 2.

La proposition V. 5 n'est pas valable lorsque $b = 3$. Ainsi :

Exemple 1 : $b = 3$, $s = 2$, $\oplus = \oplus_{\nabla}$, $N = 6$

$$x_1 = x_4 = 0,12 \quad x_2 = x_5 = 0,21 \quad x_3 = x_6 = 0,20 \quad \sigma = 11$$

Par les sommations : b à b $\bar{\sigma}_b = 11$ 2 à 2 : $\bar{\sigma}_d = 10$

Lorsque $b = 10$, la proposition V. 5 se démontre par la succession d'inégalités du type suivant, écrites pour $N = 16$:

$$\sigma_{1,1}^{(10)} \leq \sigma_{1,3}^{(2)} \oplus \sigma_{5,1}^{(2)}$$

$$\oplus \sum_{i=11}^{16} x_i \leq \sigma_{6,1}^{(2)} \oplus \sigma_{7,1}^{(2)} \oplus \sigma_{8,1}^{(2)}$$

Par les hypothèses sur les éléments x_i :

$$\sigma_{6,1}^{(2)} \oplus \sigma_{7,1}^{(2)} \oplus \sigma_{8,1}^{(2)} = \sigma_{6,1}^{(2)} + \sigma_{7,1}^{(2)} + \sigma_{8,1}^{(2)} = \sigma_{6,1}^{(2)} \oplus \sigma_{4,2}^{(2)}$$

et

$$\begin{aligned} \bar{\sigma}_b &\leq (\sigma_{1,3}^{(2)} \oplus \sigma_{5,1}^{(2)}) \oplus (\sigma_{6,1}^{(2)} \oplus \sigma_{4,2}^{(2)}) \leq \sigma_{1,3}^{(2)} \oplus (\sigma_{5,1}^{(2)} \oplus (\sigma_{6,1}^{(2)} \oplus \sigma_{4,2}^{(2)})) \\ &\leq \sigma_{1,3}^{(2)} \oplus ((\sigma_{5,1}^{(2)} \oplus \sigma_{6,1}^{(2)}) \oplus \sigma_{4,2}^{(2)}) = \bar{\sigma}_d \end{aligned}$$

Nous donnons enfin ci-dessous un exemple où la base est une puissance de 2, mais où la condition $1/2 b^p \leq x_i < b^p \quad \forall i$, n'est pas satisfaite :

Exemple 2 :

$b = 4 \quad s = 4 \quad N = 16$ avec :

$x_1 = 0,1103$	$x_2 = 0,1012$	$x_3 = 0,2013$	$x_4 = 0,1330$
$x_5 = 0,1003$	$x_6 = 0,1100$	$x_7 = 0,1131$	$x_8 = 0,2320$
$x_9 = 0,1323$	$x_{10} = 0,1330$	$x_{11} = 0,2131$	$x_{12} = 0,2222$
$x_{13} = 0,1323$	$x_{14} = 0,2002$	$x_{15} = 0,1323$	$x_{16} = 0,2200$

La somme exacte de ces nombres est $\sigma = 0,132122 \times 4^2$

Les sommes calculées en arithmétique flottante par \oplus_{∇} suivant les schémas de sommation, respectivement séquentiel, b à b et dichotomique sont :

$$\bar{\sigma}_s = 0,1312 \times 4^2$$

$$\bar{\sigma}_b = 0,1321 \times 4^2$$

$$\bar{\sigma}_d = 0,1313 \times 4^2$$

I. 4. CONCLUSION

1) On peut remarquer qu'étant donné un algorithme, succession d'opérations élémentaires, l'étude de la minimisation des erreurs d'arrondi pour un coût fixé, c'est-à-dire, suivant la définition de I. BABUSKA [1] pour un nombre d'opérations fixé, la nécessité de mémoires supplémentaires ou une plus grande difficulté de programmation étant négligées, peut être approchée sous plusieurs aspects ; un algorithme étant défini et son coût fixé, on peut, en effet, envisager :

* une étude maximaliste, c'est-à-dire la recherche de la stratégie minimisant la borne supérieure de l'erreur d'arrondi, ce maximum étant calculé lorsque les données initiales parcourent leurs espaces de définition.

* une étude probabiliste, afin de minimiser, dans les mêmes hypothèses que pour l'étude précédente, la moyenne de l'erreur d'arrondi, considérée comme variable aléatoire.

* une étude arithmétique : recherche de la stratégie minimisant l'erreur d'arrondi pour un coût donné et pour des données initiales fixées.

C'est ce dernier point de vue que nous avons adopté, auquel répondent les propositions V. 1 à V. 5.

Nous constatons que les résultats obtenus par l'étude arithmétique rejoignent les conclusions établies par l'étude maximaliste. En effet :

J. H. WILKINSON [42 p.17] a établi que la borne supérieure de l'erreur d'arrondi entachant une somme algébrique, est minimum si les termes sont sommés dans l'ordre des valeurs absolues croissantes.

P. LINZ [22] a démontré que les bornes supérieures de l'erreur d'arrondi entachant une somme de N nombres, sont, pour les schémas séquentiel et dichotomique dans le rapport $N/2 \log_2 N$, lorsque les hypothèses suivantes sont satisfaites : les N nombres sont écrits en virgule flottante normalisée, la base étant 2 et la règle définissant l'addition étant une règle de troncature.

I. BABUSKA [1] établit que le schéma dichotomique est optimal, du point de vue maximaliste, pour un coût fixé à $N - 1$ opérations.

2) Cette première partie du chapitre V comporte peu d'exemples numériques, portant sur des sommes à nombre d'éléments faible : la comparaison expérimentale entre les schémas de sommation séquentiel et dichotomique a été réalisée par plusieurs auteurs ([1] , [22]). Nous signalerons l'étude expérimentale de J. GREGORY [11] , comparant précision, rapidité d'exécution et encombrement relatifs des schémas de sommation, séquentiel, dichotomique, et d'une méthode de "correction". Enfin, pour la programmation du schéma dichotomique, on pourra voir R. J. WALKER [41] .

3) Nous soulignerons l'application des résultats précédents au calcul de sommes de séries positives : la recherche d'une meilleure précision dans ce calcul, implique de ne pas utiliser systématiquement l'algorithme usuel, de la sommation séquentielle de N termes de la série dans l'ordre des indices croissants, particulièrement dans le cas de séries monotones.

II - CALCUL DE LA SOMME D'UNE SERIE NUMERIQUE ALTERNÉE DONT LE TERME GENERAL DECROIT

II . 1 . INTRODUCTION ET NOTATIONS

T désignant le système T_b^S de nombres à virgule flottante, considérons le problème suivant :

Une série numérique alternée de terme général décroissant étant donnée sur T , soit à calculer la somme des N premiers termes de la série :

$$S = x_1 - x_2 + x_3 - \dots + (-1)^{N-1} x_N$$

avec : $0 < x_{i+1} < x_i$ $1 \leq i \leq N-1$

Le calcul sur T de la somme S définie s'effectue généralement par la sommation séquentielle des termes x_i , dans l'ordre des indices croissants.

Le but de l'étude ci-dessous est d'établir que l'algorithme de sommation sur T des termes dans l'ordre des modules croissants, calcule une meilleure approximation de S sur T .

Nous donnerons des conditions sous lesquelles, tous les chiffres de la somme calculée par ce second algorithme, sont significatifs.

Pour la commodité des notations, nous poserons :

$$X_i = (-1)^{N-i} x_{N-i+1}$$

soit

$$S = \sum_{i=1}^N X_i$$

Notons :

$$F_i = \bigoplus_{j=1}^i X_j \quad i = 1, \dots, N \text{ (sommation dans l'ordre séquentiel)}$$

$$E_i = F_i + X_{i+1} - F_{i+1} \quad i = 1, \dots, N-1$$

$$p_i = e(F_i) \quad i = 1, \dots, N \text{ (exposant de } F_i)$$

L'écart entre S et la somme $\bar{S} = F_N$, calculée sur T dans l'ordre des modules croissants, est :

$$E = S - F_N$$

vérifiant :

$$E = \sum_{i=1}^{N-1} E_i$$

II . 2 . ETUDE PRELIMINAIRE

Nous excluons les possibilités d'underflow en admettant qu'est réalisée la condition suffisante $p_1 \geq m_2 + s-1$.

Supposons par exemple X_1 positif.

La suite des sommes partielles $\sum_{j=1}^i X_j$, i impair, est positive et croissante.

La suite des sommes partielles $\sum_{j=1}^i X_j$, i pair, est négative et décroissante.

Ces propriétés sont conservées par les sommes partielles F_i calculées sur T ; en effet :

Si les conditions suivantes sont réalisées :

$$\left[\begin{array}{l} \text{L'opération } \oplus \text{ est correcte} \\ |v_i, |X_{i+1}| > (|X_i|)'' \text{ où } (|X_i|)'' \text{ désigne le successeur de } |X_i| \end{array} \right.$$

dans T , on peut affirmer que, pour tout i :

$$F_i \text{ a le signe de } (-1)^{i-1}$$

$$|F_i| < |X_{i+1}|$$

$$|F_{i+2}| > |F_i|$$

Les deux premiers résultats s'établissent conjointement par récurrence, la troisième inégalité est en conséquence de ceux ci.

Remarque : La nécessité de l'hypothèse: $|X_{i+1}| \geq (|X_i|)^s$ est illustrée par l'exemple suivant :

$b = 10, s = 2, \oplus = \oplus_{\nabla}, S = 0,26 - 1,4 + 1,4 - 11 + 11$
où $F_3 = 0,20 \quad F_5 = 0$.

D'autre part, rappelant :

Si X et Y sont deux éléments de T , d'exposants respectifs $e(X)$ et $e(Y)$ avec $e(X) \geq e(Y)$, l'erreur $\varepsilon = X + Y - (X \oplus Y)$ entachant l'addition de ces deux nombres sur T par une opération \oplus correcte, vérifie :

$$\varepsilon = \lambda b^{e(Y)-s} \quad \text{où } \lambda \text{ est un entier tel que } 0 \leq |\lambda| \leq b^{e(X+Y)-e(Y)-1}$$

λ est nul dans les cas de soustractions arithmétiques où $e(X+Y) < e(Y)$ (cancellation).

nous remarquerons que, dans le problème étudié, le phénomène de cancellation ne peut se produire pour 2 additions successives $F_{i-1} + X_i$ et $F_i + X_{i+1}$, en raison de la non-décroissance des exposants p_i de même parité, et nous énoncerons ici :

Lemme V . 2 .

Si, dans le calcul des termes F_i , une cancellation se produit au rang i , c'est-à-dire si $p_i < p_{i-1}$, l'erreur E_i au rang suivant vérifie :

$$E_i = \lambda b^{p_{i-1}-s} \quad \text{où } \lambda \text{ est un entier tel que } 0 \leq |\lambda| \leq b^{p_{i+1}-p_{i-1}-1}.$$

Soit, en effet :

$$F_i = F_{i-1} \oplus X_i$$

$$F_{i+1} = F_i \oplus X_{i+1}$$

Une cancellation lors du calcul de F_i (la condition $p_{i-1} = e(X_i) \vee e(X_i) - 1$ est nécessaire pour cela) implique, si t est l'entier positif défini par $p_i = p_{i-1} - t$, que les t derniers chiffres de la mantisse de F_i sont nuls.

L'erreur E_i entachant le calcul de F_{i+1} s'écrit alors

$$E_i = \lambda b^{p_i - s + t} \quad \text{avec} \quad 0 \leq |\lambda| \leq b^{p_{i+1} - p_i - t - 1}$$

Extrayons de la suite (p_1, p_2, \dots, p_N) la sous-suite $(p_{i_1}, p_{i_2}, \dots, p_{i_K})$ définie par les conditions suivantes :

$$\begin{cases} p_{i_1} = p_1 \\ p_{i_j} > p_{i_{j-1}} \quad (\text{strict}) \quad j \geq 2 \end{cases}$$

L'étude précédente permet d'écrire les erreurs élémentaires E_{i_j-1} et l'erreur globale E sous la forme :

$$E_{i_j-1} = \lambda_{i_j-1} b^{p_{i_j} - s} \quad \text{où} \quad 0 \leq |\lambda_{i_j-1}| \leq b^{p_{i_j} - p_{i_{j-1}} - 1} \quad 2 \leq j \leq K$$

$$E = \sum_{i=1}^{N-1} E_i = \sum_{j=2}^K E_{i_j-1}$$

II . 3 . THEOREME

On peut énoncer

Proposition V . 6 .

Si X_1, X_2, \dots, X_N sont N éléments de T , de signes alternés et de modules strictement croissants,

Si l'addition sur T , notée \oplus , est correcte,

Alors, si $\bar{S} = \bigoplus_{i=1}^N X_i$ et $S = \sum_{i=1}^N X_i$ sont respectivement,

la somme des X_i sur T et R , l'une au moins des égalités suivantes est vérifiée :

$$\bar{S} = \nabla S \vee \Delta S \quad \text{ou} \quad \bar{S} - X_N = \nabla(S - X_N) \vee \Delta(S - X_N)$$

Si $\text{Max} \{e(S), e(S - X_N)\} \leq e(X_1) + s$ (V. 5)

l'erreur $E = S - \bar{S}$ est élément de T .

Si l'inégalité (V.5) est stricte et si l'une des conditions suivantes est réalisée :

1) $e(\sum_{j=1}^{2i} X_j) \leq e(\sum_{j=1}^{2i-1} X_j) \quad \forall i$, et \bigoplus I-dirigée
(resp. \geq)

2) $k_1 + k_2 \leq s-1$ où $k_1 = \text{max} \{e(S), e(S - X_N)\} - e(X_1)$ et k_2 désigne le nombre maximum de chiffres de cancellation lors des sommes partielles

$$\bigoplus_{j=1}^i X_j, \quad 2 \leq i \leq N$$

on a :

$$E = - (\bar{S} \ominus X_N \ominus X_{N-1} \ominus \dots \ominus X_1) \quad (\text{V. 6})$$

Démonstration :

$$E = \sum_{j=2}^K E_{i_j-1}$$

La majoration :

$$|E| \leq \sum_{j=2}^K b^{p_{i_j}-s} - b^{p_{i_{j-1}}-s} = b^{p_{i_K}-s} - b^{p_1-s} < b^{p_{i_K}-s}$$

entraîne :

l'égalité $\bar{S} = F_N = \nabla S \vee \Delta S$ si $p_{i_K} = p_N$

l'égalité $F_{N-1} = \nabla(S - X_N) \vee \Delta(S - X_N)$ sinon,
d'où la première partie de la proposition.

D'autre part, E s'écrit :

$$E = \left(\sum_{j=2}^K \lambda_{i,j-1} b^{p_{i,j-1} - p_1} \right) b^{p_1 - s} = M b^{p_1 - s}$$

où

$$M = \sum_{j=2}^K \lambda_{i,j-1} b^{p_{i,j-1} - p_1} \text{ est un entier satisfaisant, si (V. 5) est}$$

vérifiée à l'inégalité :

$$|M| < b^s$$

E est donc élément de T .

Pour la troisième partie de la proposition, notons pour $1 \leq i \leq N$:

$$F_i = l_i b^{p_i - s} \quad b^{s-1} \leq |l_i| < b^s$$

Remarquons que E_i , erreur entachant l'addition correcte de deux nombres de signes contraires F_i et X_{i+1} où $e(F_i) \leq e(X_{i+1})$ vérifie :

$$F_i - E_i = \left(\lfloor l_i b^{p_i - p_{i+1}} \rfloor \vee \lceil l_i b^{p_i - p_{i+1}} \rceil \right) b^{p_{i+1} - s} \text{ si } p_{i+1} \geq p_i \quad (V.7)$$

Le calcul de E par la formule (V.6), lorsque l'une des conditions suffisantes 1) ou 2) est réalisée, s'établit en montrant que, pour $i = N - 1, \dots, 0$ les nombres :

$$U_i = F_N - \sum_{j=N}^{j=i+1} X_j$$

sont éléments de T .

Ces nombres s'écrivent :

$$U_i = F_i - \left(\sum_{j=i}^{N-1} E_j \right) \quad (V.8)$$

Nous distinguerons deux cas, suivant les valeurs relatives de p_i et p_{i+1} :

a) $p_{i+1} > p_i$

$p_{i+1} = p_i + u$ où u est un entier positif.

On décomposera :

$$U_i = F_i - E_i - \left(\sum_{j=i+1}^{N-1} E_j \right)$$

Un calcul semblable à celui effectué sur E pour montrer que E est élément de T , permet d'écrire :

$$\sum_{j=i+1}^{N-1} E_k = M b^{p_{i+1}-s} \quad \text{où } M \text{ est entier et vérifie } |M| < b^{s-1}$$

(L'inégalité V.5 est stricte)

D'autre part, par la formule (V.7) :

$$F_i - E_i = L b^{p_{i+1}-s} \quad \text{où } L \text{ est entier et vérifie } b^{s-1-u} \leq |L| \leq b^{s-u}$$

De l'écriture :

$$U_i = (L - M) b^{p_{i+1}-s}$$

où l'entier relatif $L - M$ satisfait à $|L - M| < b^s$

on déduit $U_i \in T$.

b) $p_{i+1} \leq p_i$

Nécessairement $e(X_{i+1}) = p_i \vee p_i + 1$. Distinguons les 2 cas :

$$\underline{(b-1)e(X_{i+1}) = p_i + 1}$$

On a donc :

$$e(X_j) \geq p_i + 1 \quad \text{pour } j \geq i + 1.$$

Soit λ , s'il existe, le plus petit indice supérieur à i vérifiant $p_{\lambda+1} > p_i$

S'il n'existe pas un tel indice dans l'ensemble $\{ i + 1 , \dots , N - 1 \}$ le problème est terminé puisque $U_i = F_i$, sinon décomposons (V.8) sous la forme :

$$U_i = (F_\lambda - \sum_{j=\lambda}^{N-1} E_j) - (\sum_{j=i+1}^{\lambda} X_j)$$

Pour le premier terme :

$p_{\lambda+1} > p_i \geq p_\lambda$ implique par le cas a) :

$$F_\lambda - \sum_{j=\lambda}^{N-1} E_j = m b^{p_{\lambda+1}-s}$$

et en notant $p_{\lambda+1} = p_i + t$ $t > 0$

$$F_\lambda - \sum_{j=\lambda}^{N-1} E_j = m b^{t-1} \cdot b^{p_i+1-s} \quad \text{avec } |m b^{t-1}| < b^s$$

- Cette dernière majoration résulte de l'étude faite en a) , la notation $p_{\lambda+1} = p_\lambda + u$ ($u \geq t$) permettant d'écrire :

$$|m| < b^{s-u} + b^{s-1-t} < b^{s-t} + b^{s-1-t} \quad \text{si } F_\lambda - \sum_{j=\lambda}^{N-1} E_j \text{ a le signe de } F_\lambda$$

$$|m| < b^{s-1-t} \quad \text{si " est du signe contraire à } F_\lambda -$$

Pour le second terme, remarquons :

$$e(X_j) = p_i + 1 \quad i + 1 \leq j \leq \lambda$$

donc :

$$\sum_{j=i+1}^{\lambda} X_j = n b^{p_i+1-s} \quad \text{où } n \text{ est entier}$$

De l'écriture :

$$\sum_{j=i+1}^{\lambda} X_j = (\sum_{j=i+1}^{\lambda-1} X_j) + X_\lambda$$

où le terme entre parenthèse est du signe de $X_{\lambda-1}$ et de module inférieur à X_{λ} ,

on tire :

$$|n| < b^s - b^{s-1}$$

$U_i \in T$

$$\underline{b - 2) e(X_{i+1}) = p_i}$$

U_i étant considéré sous la forme (V.8), on remarquera que l'on peut écrire :

$$\sum_{j=i}^{N-1} E_j = M b^{p_i - s} \quad \text{où } M \text{ est entier et satisfait à } |M| < b^{s-1}$$

Il faut alors vérifier que dans la soustraction :

$$F_i - \sum_{j=i}^{N-1} E_j = (l_i - M) b^{p_i - s} \quad (V.9)$$

l'entier relatif $l_i - M$ satisfait à $|l_i - M| < b^s$.

Distinguons deux cas :

$$b - 2 - 1) \quad p_{i+1} = p_i$$

De :

$$F_{i+1} = F_i + X_{i+1}, \quad p_{i+1} = p_i = e(X_{i+1}) \quad \text{on déduit :}$$

$$|l_i| < (b-1) b^{s-1}$$

d'où le résultat $U_i \in T$.

$$b - 2 - 2) \quad p_{i+1} < p_i$$

Suivant les conditions 1) ou 2) de la proposition V.6 :

* Si la condition 1) est vérifiée avec par exemple :

$$e\left(\sum_{j=1}^{2i} X_j\right) \leq e\left(\sum_{j=1}^{2i-1} X_j\right)$$

alors l'inégalité $p_{i+1} < p_i$ ne peut avoir lieu que pour des valeurs impaires de l'indice i .

Pour tout i ,

$$E_{2i-1} = 0$$

E_{2i} a le signe de F_{2i+1}

Dans la formule (V.9), les éléments F_i et $\sum_{j=i}^{N-1} E_j$ sont de même signe et $U_i \in T$.

* Si la condition 2) est vérifiée, on a :

$$|M| < b^{k_1}$$

D'autre part :

$$|F_i| = |X_{i+1} - F_{i+1}|$$

et

$$|l_i| < b^s - b^{s-k_2-1}$$

La condition $k_1 + k_2 \leq s - 1$ assure alors $|l_i - M| < b^s$ dans (V.9)

Pour tout i , $F_N - \sum_{j=N}^{j=i+1} X_j$ est élément de T , donc, \oplus étant

correcte, on obtient successivement :

$$F_N \ominus X_N = F_N - X_N$$

puis

$$(F_N \ominus \sum_{j=N}^{i+1} X_j) \ominus X_i = F_N - \sum_{j=N}^i X_j \quad i = N - 1, \dots, 1$$

d'où le résultat.

REMARQUES :

1) Les conditions suffisantes 1) et 2) de la proposition V.6 sont assez restrictives. Relativement à ce point, nous donnerons les deux exemples suivants pour lesquels la formule (V.6) n'est pas exacte.

Exemple 1.

$b = 10$, $s = 2$, $S = 0,29 - 0,63 + 0,64 - 0,71 + 1,6 - 2,2 + 12$, \oplus I-dirigée.

La condition suffisante 1) est vérifiée mais l'inégalité V.5 n'est pas stricte.

$$\bar{S} = 10 \quad E = 0,99 \quad \bar{S} \ominus \sum_{i=N}^1 X_i = - 0,96$$

Exemple 2.

$b = 10$, $s = 3$, $\oplus = \oplus_{\Delta}$, $S = 0,984 - 0,995 + 0,998 - 2,23 + 25$

$k_1 = 2$, $k_2 = 1$

$$\bar{S} = 23,8 \quad E = - 0,043 \quad \bar{S} \ominus \sum_{i=N}^1 X_i = 0,046$$

2) Si la formule (V.6) n'est pas applicable dans tous les cas, nous remarquerons que dès que E est élément de T , cette erreur est directement calculable par la formule :

$$E = \sum_{j=2}^K E_{i_j-1}$$

Il sera fait à chaque pas i une comparaison des exposants p_{i+1} et $\max_{j \leq i} p_j$, les erreurs élémentaires, lorsqu'elles seront non nulles,

seront calculées par la formule - chapitre II - :

$$E_i = - (F_{i+1} \ominus X_{i+1} \ominus F_i)$$

II. 4. EXEMPLES NUMERIQUES ET MODE OPERATOIRE

1) Le calcul des 2 sommes $\sigma_3 = \sum_{i=1}^N \frac{(-1)^{i-1}}{i}$ et

$\sigma_4 = \sum_{i=1}^N \frac{(-1)^{i-1}}{i^2}$, sur une I.B.M. 360, a donné les résultats joints

(à comparer à ceux du paragraphe I. 5 du chapitre III, p.62):

Dans les 2 tableaux de nombres, A -nombres écrits en base 16- et B -nombres écrits en base 10- :

$\bar{\sigma}_3$ et $\bar{\sigma}_4$ sont les sommes calculées sur T dans l'ordre des modules croissants

E_{σ_3} et E_{σ_4} sont les erreurs entachant $\bar{\sigma}_3$ et $\bar{\sigma}_4$, calculées par la formule (V. 6)

σ_{3DP} et σ_{4DP} sont les valeurs exactes de σ_3 et σ_4 , calculées en double précision.

(Les tableaux A et B, représentant les mêmes résultats dans 2 bases différentes, ont été tous deux reproduits pour conserver précision et compréhension des résultats. Relativement au tableau B nous soulignerons que, bien que les valeurs $\bar{\sigma}_3$ et $\bar{\sigma}_4$ soient des éléments de T, leur sortie est effectuée en double précision pour éviter une catastrophique erreur de changement de base).

On vérifie bien :

- que tous les chiffres de mantisse de $\bar{\sigma}_3$ et $\bar{\sigma}_4$ sont exacts

- que E_{σ_3} et E_{σ_4} donnent la valeur exacte des erreurs $\sigma_3 - \bar{\sigma}_3$

et $\sigma_4 - \bar{\sigma}_4$: nous avons souligné les chiffres de mantisse d'ordre supérieur à $s = 6$ dans la double précision ainsi que la mantisse du terme d'erreur.

N	σ_3	E_{σ_3}	σ_{3DP}	σ_4	E_{σ_4}	σ_{4DP}
100	40B02C0E	80000000	40B02C0E00000000	40D289F4	3A8BF000	40D289F48BF00000
200	40B0CEAB	BA000000	40B0CEAA40000000	40D28C62	3A59F000	40D28C6259F00000
300	40B1050D	BA1B0000	40B1050CE5000000	40D28CD6	3A1F2000	40D28CD61F200000
400	40B12048	BAFA0000	40B1204706000000	40D28CFE	3A851000	40D28CFE85100000
500	40B130A0	RA1F0000	40B1309FE1000000	40D28C11	3A87C200	40D28D1187C20000
600	40B13B87	B9E00000	40B13B86F2000000	40D28D1B	3AC13800	40D28D1BC1380000
700	40B14352	BADC0000	40B1435124000000	40D28D21	3AEC2900	40D28D21EC290000
800	40B1492A	BAE70000	40B1492919000000	40D28D25	3AED4600	40D28D25ED460000
900	40B14DB5	BAZE0000	40B14DB4D2000000	40D28D28	3AAC4400	40D28D28AC440000
1000	40B15158	BA190000	40B15157E7000000	40D28D2A	3AA33600	40D28D2AA3360000

N	σ_3	E_{σ_3}	σ_{3DP}	σ_4	E_{σ_4}	σ_{4DP}
100	0.68817222118378	0.0	0.68817222118378	0.82241749763489	0.3258174E-07	0.82241753021663
200	0.69065350294113	-0.4470348E-07	0.69065345823765	0.82245457172394	0.2094021E-07	0.82245459266414
300	0.69148331880569	-0.6286427E-08	0.69148331251927	0.82246148586273	0.7246854E-08	0.82246149310959
400	0.69189882278442	-0.5820766E-07	0.69189876457676	0.82246387004852	0.4291451E-07	0.82246391296303
500	0.69214820861816	-0.7217750E-08	0.69214820140041	0.82246500253677	0.3160858E-07	0.82246503414535
600	0.69231456518173	-0.3259629E-08	0.69231456192210	0.82246559858322	0.4498725E-07	0.82246564357C47
700	0.69243347644806	-0.5122274E-07	0.69243342522532	0.82246595621109	0.5498532E-07	0.82246601119641
800	0.69252264499664	-0.5378388E-07	0.69252259121276	0.82246619462967	0.5524453E-07	0.82246624987420
900	0.69259196519852	-0.1071021E-07	0.69259195448831	0.82246637344360	0.4010872E-07	0.82246641355232
1000	0.69264745712280	-0.5820766E-08	0.69264745130204	0.82246649265289	0.3800051E-07	0.8224665306534C

2) Le sous-programme reproduit ci-dessous:

SOMSER (U, N, S, E) réalise, à partir du tableau U des N premiers termes d'une série numérique alternée de terme général décroissant, les

$$\text{calculs de } \bar{S} = \oplus \sum_{i=N}^1 U(i) \text{ et de } E = \left(\sum_{i=1}^N U(i) \right) - \bar{S}.$$

Ce dernier calcul est effectué par la méthode de comparaison à chaque pas des exposants p_{i+1} et $\max_{j \leq i} p_j$. La fonction auxiliaire NEXP (Z)

isole l'exposant du nombre flottant Z.

```
SUBROUTINE SOMSER(U,N,S,E)
DIMENSION U(N)
S=0.0
E=0.0
NO=NEXP(U(N))
DO 1 I=1,N
L=N+1-I
T=S+U(L)
NP=NEXP(T)
IF(NP.EQ.NO)GOTO 2
NO=NP
E=E-(T-U(L))-S
2 CONTINUE
S=T
1 CONTINUE
RETURN
END
```

```
INTEGER FUNCTION NEXP(Z)
EQUIVALENCE(M,U)
L=16777216
U=ABS(Z)
M=M/L
NEXP=M
RETURN
END
```

Mode opératoire

La recherche d'une meilleure précision dans le calcul de la somme d'une série numérique alternée, par la sommation des termes dans l'ordre croissant, entraîne une modification au niveau de la programmation : Il est en effet nécessaire, avant de procéder à la sommation proprement dite, de déterminer l'ordre N du plus petit terme, en valeur absolue, à calculer.

Ceci pourra se faire par un balayage rapide, à l'aide d'un indice k augmentant à partir de 1, par "pas" large de P (par exemple $P = 25$ ou 50), une comparaison de x_k à la précision cherchée étant établie pour chaque pas. Ce balayage et, pour une précision donnée, le nombre d'additions supplémentaires de termes de la série que peut entraîner par rapport à la méthode classique, l'utilisation du pas P pour déterminer l'arrêt des calculs, est largement compensé par la précision obtenue.

Relativement au critère d'arrêt lors du calcul numérique de la limite d'une suite, mention doit être faite du critère de convergence numérique établi par K. NICKEL et K. RITTER [30] : Le critère donné par ces auteurs, pour la détermination, en fonction de la longueur s de la représentation du mot-machine, de l'indice $N(s)$ du dernier terme à calculer, permet d'assurer la convergence numérique sous des conditions suffisantes assez faibles.

CHAPITRE VI

UTILISATION CONJOINTE
DE DEUX ARITHMETIQUES CORRECTES

RODUCTION

Soit S un sous-ensemble fini de \mathbb{R} contenant 0 .

A l'opération $*$ appartenant à l'ensemble $\{+, -, \times, /\}$ des opérations élémentaires définies sur \mathbb{R} , il a été associé les 2 opérations correctes, définies sur S , \otimes_{∇} et \otimes_{Δ} . $\forall (x, y) \in S^2 \wedge x * y \in \mathbb{R}_S$:

$$x \otimes_{\nabla} y = \text{Max} \{ z \mid z \in S \wedge z \leq x * y \}$$

$$x \otimes_{\Delta} y = \text{Min} \{ z \mid z \in S \wedge z \geq x * y \}$$

Une suite $G = \{*_i, i=1, \dots, L\}$ d'opérations sur \mathbb{R} étant donnée, il est possible de faire correspondre à G , 2^L suites d'opérations correctes sur S :

$$O = \{ \otimes_{\square_i} \} \text{ où } \square_i \text{ prend la valeur } \nabla \text{ ou } \Delta.$$

Deux problèmes se posent alors :

- Etude de l'ensemble des résultats obtenus lors de l'application des 2^L suites précédentes à une même suite finie de $L + 1$ opérandes.

- Possibilité de déterminer, parmi les suites d'opérations sur S , 2 suites particulières calculant, à partir de la suite donnée d'opérandes, deux valeurs d'encadrement de la valeur définie sur \mathbb{R} par l'application de G aux données initiales.

I - DEFINITIONS

$N \in \mathbb{N}$, étant donné un vecteur $x = (x_1, x_2, \dots, x_N)$ de S^N et une suite d'opérations sur \mathbb{R} , $\mathcal{G} = \{\ast_1, \ast_2, \dots, \ast_{N-1}\}$ définissant le réel, noté $i_{\mathcal{G}}(x)$, $r = r_N$ par : $r_1 = x_1$

$$r_{i+1} = r_i \ast_i x_{i+1} \quad i=1, \dots, N-1,$$

soit $\mathcal{O} = \{\circledast_1, \circledast_2, \dots, \circledast_{N-1}\}$ une suite d'opérations sur S , non nécessairement correctes correspondant aux opérations de \mathcal{G} .

On appellera "image de x par \mathcal{O} " et on notera $i_{\mathcal{O}}(x)$ le nombre $\bar{r} = \bar{r}_N$ défini par : $\bar{r}_1 = x_1$

$$\bar{r}_{i+1} = \bar{r}_i \circledast_i x_{i+1} \quad i=1, \dots, N-1.$$

Nous poserons :

Définition VI. 1.

La suite d'opérations sur S , $\mathcal{O} = \{\circledast_1, \dots, \circledast_{N-1}\}$ sera dite une G - suite (resp. D - suite) de \mathcal{G} sur la partie \mathcal{P} de S^N si :

$$i_{\mathcal{O}}(x) \leq i_{\mathcal{G}}(x) \quad , \forall x \in \mathcal{P}$$

(resp. $i_{\mathcal{O}}(x) \geq i_{\mathcal{G}}(x) \quad , \forall x \in \mathcal{P}$)

Définition VI. 2.

Une G - suite $\tilde{\mathcal{O}} = \{\tilde{\circledast}_1, \dots, \tilde{\circledast}_{N-1}\}$ de \mathcal{G} sur \mathcal{P} , sera dite minimisante (resp. maximisante) sur \mathcal{P} si :

- les opérations $\tilde{\circledast}_i$ sont correctes
- pour toute autre G - suite $\mathcal{O} = \{\circledast_1, \dots, \circledast_{N-1}\}$ de \mathcal{G} sur \mathcal{P} d'opérations correctes, on a :

$$i_{\mathcal{O}}(x) \leq i_{\tilde{\mathcal{O}}}(x) \quad , \forall x \in \mathcal{P}$$

(resp. $i_{\mathcal{O}}(x) \geq i_{\tilde{\mathcal{O}}}(x) \quad , \forall x \in \mathcal{P}$)

De façon analogue, on définit une D - suite minimisante (resp. maximisante) de \mathcal{G} sur \mathcal{P} .

Définition VI. 3.

$O = \{ \circledast_i \}$ étant une suite d'opérations correctes sur S , correspondant à une suite $\mathcal{G} = \{ *_{i} \}$ d'opérations définies sur \mathbb{R} , on appelle suite conjuguée de O , et on note $\bar{O} = \{ \bar{\circledast}_i \}$ la suite définie

$$\text{par : } \bar{\circledast}_i = \begin{cases} \circledast_i_{\Delta} & \text{si } \circledast_i = \circledast_i_{\nabla} \\ \circledast_i_{\nabla} & \text{si } \circledast_i = \circledast_i_{\Delta} \end{cases}$$

Les notations et définitions précédentes s'étendent lorsque la relation d'ordre définie dans l'espace image n'est pas une relation d'ordre totale -succession d'opérations élémentaires faisant correspondre, à un vecteur de S^N , un élément d'une puissance de S , vecteur ou matrice-.

Etant donnée une suite \mathcal{G} d'opérations sur \mathbb{R} , l'existence, dans l'ensemble des suites d'opérations correctes sur S correspondant à \mathcal{G} d'une G - suite (resp. D - suite) de \mathcal{G} n'est pas assurée.

II - SUITES MAXIMISANTES D'OPERATIONS SUR S RELATIVES AUX ALGORITHMES USUELS [35]

II. 1. SOMME ALGEBRIQUE ET PRODUIT SCALAIRE

2 N éléments x_i et y_i de S étant donnés ($i=1, \dots, N$),

soient :

$$r = \sum_{i=1}^N x_i \quad , \quad r' = \sum_{i=1}^N x_i y_i$$

Désignant par \mathcal{G}_1 et \mathcal{G}_2 les suites d'opérations définissant r et r' , on a immédiatement le résultat suivant :

Proposition VI. 1.

Les deux suites $\tilde{\mathcal{O}}_G = \{ \textcircled{*}_i \nabla \}$ et $\tilde{\mathcal{O}}_D = \{ \textcircled{*}_i \Delta \}$ sont, respectivement, une G - suite et une D - suite maximisantes de \mathcal{G}_1 sur S^N , de \mathcal{G}_2 sur S^{2N} .

Conséquences pratiques :

Un cas où la mise en oeuvre des suites maximisantes d'opérations précédentes, sur un système $T_b^S = T$ de nombres à virgule flottante, est très simplement résolue, est celui de la somme algébrique lorsqu'on dispose de l'opération

$$\textcircled{+} \nabla \text{ (resp. } \textcircled{+} \Delta \text{)} :$$

Considérons le cas où l'opération $\textcircled{+}$, d'addition sur T , est définie par $\textcircled{+} \nabla$.

La symétrie ($z \rightarrow -z$) n'étant pas un automorphisme de T , on formera :

$$r^1 = \textcircled{+} \sum_{i=1}^N x_i$$

$$r^2 = - \{ \textcircled{+} \sum_{i=1}^N (-x_i) \}$$

Alors :

$$r^1 = i \tilde{\mathcal{O}}_G(x), \quad r^2 = i \tilde{\mathcal{O}}_D(x) \quad x = (x_1, \dots, x_N).$$

Dans les autres cas, il sera préférable de générer directement les opérations $\textcircled{*} \nabla$ et $\textcircled{*} \Delta$. Nous notons ci-dessous, à titre d'exemple,

les encadrements des sommes $\sigma_1 = \sum_{i=1}^N \frac{1}{i}$ et $\sigma_3 = \sum_{i=1}^N \frac{(-1)^{i-1}}{i}$, calculés

sur une IBM 360 ($\sigma_{1 D P}$ et $\sigma_{3 D P}$ désignent les valeurs obtenues en double précision).

	N	100	200	300	400	500
σ_1	$\sigma_{1 G}$	5,187340	5,877946	6,282538	6,569756	6,792601
	$\sigma_{1 D}$	428	8129	816	70128	3069
	$\sigma_{1 D P}$	377	8030	663	69929	2822
σ_3	$\sigma_{3 G}$	0,6881702	0,6906486	0,6914756	0,6918882	0,6921345
	$\sigma_{3 D}$	42	579	904	9088	610
	$\sigma_{3 D P}$	22	534	833	8987	482

II. 2. EVALUATION NUMERIQUE D'UN POLYNOME

Un polynôme P étant défini sur S par ses coefficients a_i ($i = 0, \dots, N$), soit r la valeur de ce polynôme en un point $x \in S$:

$$r = P(x) = a_0 x^N + \dots + a_N$$

Désignant par \mathcal{G} la suite d'opérations

$\mathcal{G} = \{ *_{2i-1} = x, *_{2i} = +, i = 1, \dots, N \}$ qui permet le calcul de r par le schéma de Horner, on peut énoncer :

Proposition VI. 2

Pour le calcul de la valeur d'un polynôme en un point par le schéma de Horner, une G - suite et une D - suite maximisantes de \mathcal{G} sont respectivement :

$$\begin{aligned} & \text{sur } S^{N+1} \times S^+ \\ \tilde{\mathcal{O}}_G &= \{ \textcircled{*}_i = \textcircled{*}_i \nabla \mid i = 1, \dots, 2N \} & \tilde{\mathcal{O}}_D &= \overline{\tilde{\mathcal{O}}_G} \\ & \text{sur } S^{N+1} \times S^- \\ \tilde{\mathcal{O}}_G &= \mathcal{O}_1 \text{ si } N \text{ impair, } \mathcal{O}_2 \text{ si } N \text{ pair} & \tilde{\mathcal{O}}_D &= \overline{\tilde{\mathcal{O}}_G} \end{aligned}$$

avec :

$$\begin{aligned} \mathcal{O}_1 &= \{ \textcircled{*}_i = \textcircled{*}_i \nabla \mid \text{si } i \equiv 1 \text{ ou } 2 \pmod{4} \\ & \quad = \textcircled{*}_i \Delta \mid \text{si } i \equiv 3 \text{ ou } 0 \pmod{4} \} \\ \mathcal{O}_2 &= \{ \textcircled{*}_i = \textcircled{*}_i \Delta \mid \text{si } i \equiv 1 \text{ ou } 2 \pmod{4} \\ & \quad = \textcircled{*}_i \nabla \mid \text{si } i \equiv 3 \text{ ou } 0 \pmod{4} \} \end{aligned}$$

Exemple : valeurs d'encadrement du polynôme $P_N(x) = \prod_{j=1}^N (x-j)$ au point $x = N - 0,01$ (r et \bar{r} désignent respectivement, la valeur exacte et la valeur obtenue par l'arithmétique d'une IBM 360 - 30).

N	r_G	r_D	r	\bar{r}
5	- 0,2359009	- 0,2335205	- 0,2350349	- 0,2339172
6	- 1,233643	- 1,117676	- 1,172824	- 1,171875
7	- 9,113281	- 4,878906	- 7,025215	- 6,410156
8	- 108,4023	- 8,023438	- 49,10625	- 88,42188
9	- 2594	2417,438	- 392,3589	- 88

La mise en évidence du mauvais conditionnement de la plus grande racine de $P_N(x)$, est ainsi particulièrement nette.

II. 3. CALCUL DE LA RACINE CARREE D'UN NOMBRE PAR LA METHODE DE NEWTON

Etant donné un élément c de S^+ , la détermination de la racine carrée de c par la méthode de Newton, est réalisée, à partir d'une valeur initiale x_0 , à l'aide de la formule de récurrence :

$$x_{n+1} = \frac{1}{2} \left(x_n + \frac{c}{x_n} \right) \quad (\text{VI. 1})$$

Désignons par G la suite d'opérations qui permet, à partir des données c et x_0 , le calcul de x_{n+1} , et supposons, sans restreindre le problème, $x_0 \geq \sqrt{c}$.

On peut énoncer :

Proposition VI. 3.

Si S est un système de nombres à virgule flottante de base 2, pour la détermination du $(n+1)$ ième itéré de (VI.1) à partir d'une valeur initiale $x_0 \geq \sqrt{c}$, une G -suite et une D -suite d'opérations sur S , maximisantes de la suite des opérations sur \mathbb{R} , sont respectivement :

$$\tilde{O}_G = \{ \textcircled{*}_i = \textcircled{x}_i \nabla \} \quad , \quad \tilde{O}_D = \{ \textcircled{*}_i = \textcircled{x}_i \Delta \}$$

La proposition VI. 3. s'établit par récurrence sur l'ordre n .
Nous ferons au préalable la remarque suivante :

Si les opérations $\textcircled{+}$ et $\textcircled{/}$ définies sur $S = T_2^S$ sont correctes, les éléments v et \bar{v} définis respectivement sur \mathbb{R} et T_2^S , à partir d'un élément $u \geq \sqrt{c}$ de T_2^S par :

$$v = \frac{1}{2} \left(u + \frac{c}{u} \right) \quad \bar{v} = \frac{1}{2} \left(u \textcircled{+} c \textcircled{/} u \right)$$

vérifient :

$$\bar{v} = \nabla v \vee \Delta v \quad (\text{VI. 2})$$

(VI. 2) est obtenue par :

a) $|c/u - c \oslash u| < \varepsilon(c/u)$

b) $u \geq \sqrt{c} \Rightarrow c \oslash u \leq u$ donc $|(u + c \oslash u) - (u \oplus c \oslash u)| \leq \varepsilon(u + c \oslash u) - \varepsilon(c \oslash u)$

assurant ainsi :

$$|v - \bar{v}| < \varepsilon(\bar{v})$$

- Pour conclure à l'égalité $\bar{v} = \nabla v \vee \Delta v$ dans le cas particulier où \bar{v} vérifie $\varepsilon(\bar{v}) \neq \varepsilon(\bar{v})$, la démonstration est analogue à la démonstration établie lors du cas particulier similaire de la proposition V. 1 du chapitre V -

Il vient alors, pour la suite $\tilde{\mathcal{O}}_G$ par exemple :

Le couple $\tilde{\mathcal{O}}_G = \{\oslash_{\nabla}, \oplus_{\nabla}\}$ définit, à l'ordre 1, une G - suite maximisante d'opérations de \mathcal{G} , et l'image \tilde{x}_1 de x_0 par $\tilde{\mathcal{O}}_G$ vérifie $\tilde{x}_1 \geq \nabla(\sqrt{c})$.

Supposons qu'à l'ordre n, la suite $\tilde{\mathcal{O}}_G = \{\otimes_i\}_{\nabla}$ soit une G - suite maximisante de \mathcal{G} et que l'inégalité $\tilde{x}_n \geq \nabla(\sqrt{c})$ soit satisfaite : les nièmes itérés \tilde{x}_n et \bar{x}_n de (VI. 1), calculés sur T_2^s par respectivement, $\tilde{\mathcal{O}}_G$ et une suite \mathcal{O} d'opérations correctes sur T_2^s correspondant à \mathcal{G} , vérifient $\bar{x}_n \geq \tilde{x}_n$

1) Si $\tilde{x}_n > \sqrt{c}$, les valeurs $\bar{x}_{n+1} = \frac{1}{2}(\bar{x}_n + \frac{c}{\bar{x}_n})$ et $\tilde{x}_{n+1} = \frac{1}{2}(\tilde{x}_n + \frac{c}{\tilde{x}_n})$

satisfont :

$$\bar{x}_{n+1} \geq \tilde{x}_{n+1} > \sqrt{c},$$

en raison du sens de variation de la fonction $y = \frac{1}{2}(x + \frac{c}{x})$ qui de plus admet, au point $x = \sqrt{c}$, un minimum égal à \sqrt{c} .

(VI. 2) assure alors :

$$\bar{x}_{n+1} \geq \tilde{x}_{n+1} \geq \nabla(\sqrt{c})$$

2) Si $\tilde{x}_n = \nabla(\sqrt{c}) \leq \sqrt{c}$, le processus itératif de Newton est terminé. En effet, montrons que les conditions considérées entraînent l'une des deux possibilités :

$\tilde{x}_{n+1} = \nabla(\sqrt{c})$ donc convergence forte de $\{\tilde{x}_n\}$ vers $\nabla(\sqrt{c})$.

$\tilde{x}_{n+1} = \Delta(\sqrt{c}) > \sqrt{c}$ et convergence faible de $\{\tilde{x}_n\}$, les itérés cyclant sur les 2 valeurs $\nabla(\sqrt{c})$ et $\Delta(\sqrt{c})$ - la propriété de G - suite de \tilde{O}_G n'est plus vérifiée dans ce cas limite -

De :

$$\frac{1}{2} (\nabla \sqrt{c} + c/\nabla \sqrt{c}) - \Delta \sqrt{c} = [c - (\Delta \sqrt{c} + \varepsilon(\sqrt{c})) \nabla \sqrt{c}] / 2\nabla \sqrt{c} \leq 0$$

on conclut :

$$\tilde{x}_{n+1} = \nabla \sqrt{c} \vee \Delta \sqrt{c},$$

la condition $c < (\Delta \sqrt{c} + \varepsilon(\sqrt{c})) \nabla \sqrt{c}$ étant une condition nécessaire et suffisante de convergence forte.

Une conséquence de la démonstration précédente est alors pour la suite $\{\tilde{x}_n\}$, des itérés calculés par \tilde{O}_G , décroissante tant que $\tilde{x}_n > \nabla(\sqrt{c})$:

Proposition VI. 4.

Un élément c étant donné sur T_2^S , la suite des itérés du calcul de la racine carrée de c sur T_2^S par la méthode de Newton lorsque les opérations d'addition et de division sur T_2^S sont respectivement \oplus_{∇} et \oslash_{∇} , converge, quelle que soit la valeur initiale $x_0 \in T_2^S$:
 soit fortement vers $\nabla(\sqrt{c})$
 soit faiblement en cyclant sur les valeurs $\nabla(\sqrt{c})$ et $\Delta(\sqrt{c})$.

J. M. YOHE [45] a établi la convergence forte vers $(\nabla \sqrt{c})''$ de la suite des itérés de Newton, calculés à partir d'une valeur initiale $x_0 \neq \sqrt{c}$, lorsque les opérations d'addition et de division sur T_2^S sont respectivement \oplus_{Δ} et \oslash_{Δ} .

Exemples :

1) $b = 2 \quad s = 3 \quad c = 0,100 \quad \sqrt{c} = 0,10110 \dots 0 = \overset{\sim}{0}_G$
 Convergence forte de $\{ \tilde{x}_n \}$ vers $0,101$.

2) $b = 2 \quad s = 3 \quad c = 0,110 \quad \sqrt{c} = 0,11011 \dots 0 = \overset{\sim}{0}_G$
 Convergence faible de $\{ \tilde{x}_n \}$ qui cycle sur $0,110$ et $0,111$.

Remarque

Si le calcul de la racine carrée d'un nombre par la méthode de Newton, est effectué en arithmétique flottante de base b différente de 2 :

La proposition VI. 3 reste valable, si le nième itéré de (VI. 1) calculé par $\overset{\sim}{0}_G$ satisfait à $\tilde{x}_n \geq \sqrt{c}$.

* Il suffit en effet, de reprendre le raisonnement par récurrence de la proposition, en remarquant que l'hypothèse :

$$\sqrt{c} \leq \tilde{x}_n \leq \bar{x}_n$$

qui implique $\tilde{x}_n + c/\tilde{x}_n \leq \bar{x}_n + c/\bar{x}_n$,

$$\text{assure } \tilde{x}_n \oplus_{\nabla} c \oslash_{\nabla} \tilde{x}_n \leq \bar{x}_n \oplus c \oslash \bar{x}_n$$

par l'égalité $u \oplus c \oslash u = \nabla(u + c/u) \oslash_{\Delta} \Delta(u + c/u) \quad \forall u \in T_b^S \wedge u \geq \sqrt{c}$
 la monotonie de l'arrondi permet alors de conclure $\tilde{x}_{n+1} \leq \bar{x}_{n+1}$ *

La convergence des itérés de Newton vers les 2 nombres flottants entourant \sqrt{c} n'est par contre plus assurée, la formule (VI. 2) étant à remplacer, avec les mêmes notations, pour une base b multiple de 2, par l'inégalité plus large : $|v - \bar{v}| < \frac{b}{2} \varepsilon(\bar{v})$.

La propriété de convergence forte de la suite des itérés \tilde{x}_n calculés par $\overset{\sim}{0}_D$ est toujours satisfaite lorsque b est multiple de 2.

*En effet, $\tilde{x}_n > \sqrt{c}$ assurant $\bar{x}_n + c / \bar{x}_n < 2 \bar{x}_n$, 2 cas peuvent se présenter:

$$\Delta (\bar{x}_n + c/\bar{x}_n) < 2\bar{x}_n \text{ impliquant } \tilde{x}_{n+1} \leq \bar{x}_n$$

$$\Delta (\bar{x}_n + c/\bar{x}_n) = \Delta(2\bar{x}_n) \text{ impliquant } \tilde{x}_{n+2} = \tilde{x}_{n+1} = \frac{1}{2} \Delta (2\bar{x}_n)$$

$$(\frac{1}{2} \Delta (2u) \in T, \forall u \in T, \tilde{x}_{n+1} \geq \bar{x}_n \Rightarrow \tilde{x}_n + c/\tilde{x}_n \leq \tilde{x}_{n+1} + c/\tilde{x}_{n+1} < 2 \tilde{x}_{n+1}) *$$

II.4. DECOMPOSITION LR ET INVERSE D'UNE M - MATRICE

II.4.1. Décomposition LR d'une M - matrice

Soit $A = (a_{ij})$ une M-matrice donnée sur S Rappelons :

Définition. $A \in \mathcal{X}_{n,n}(\mathbb{R})$ est une M-matrice si les deux axiomes suivants sont vérifiés :

- les éléments hors diagonaux de A sont négatifs ou nuls.
- tous les mineurs principaux de A sont positifs.

La décomposition d'une M-matrice A en un produit LR où $L = (\ell_{ij})$ est triangulaire inférieure ne comportant que des "1" sur la diagonale, et $R = (r_{ij})$ triangulaire supérieure, est toujours possible par les formules classiques :

$$r_{p+1,i} = a_{p+1,i} - \sum_{k=1}^p \ell_{p+1,k} r_{k,i} \quad , p+1 \leq i \leq n \quad (VI. 3)$$

$$\ell_{i,p+1} = (a_{i,p+1} - \sum_{k=1}^p \ell_{i,k} r_{k,p+1}) / r_{p+1,p+1} \quad , p+2 \leq i \leq n$$

$p = 0, 1, \dots, n-1.$

On a le résultat suivant [7]

Si A est une M - matrice, alors dans la décomposition LR de A :

- les éléments diagonaux de R sont positifs.
- les éléments hors-diagonaux de L et R sont négatifs ou nuls.

Désignant par \mathcal{G} la suite d'opérations (VI.3) définissant L et R , considérons les deux suites particulières suivantes de calcul sur S de L et R , de résultats notés \tilde{L} et \tilde{R} :

$$\tilde{r}_{p+1,i} = a_{p+1,i} \ominus_{\Delta} \left\{ \oplus_{\nabla} \sum_{k=1}^p \tilde{\ell}_{p+1,k} \otimes_{\nabla} \tilde{r}_{k,i} \right\} \quad (O_D)$$

$$\tilde{\ell}_{i,p+1} = \left[a_{i,p+1} \ominus_{\Delta} \left\{ \oplus_{\nabla} \sum_{k=1}^p \tilde{\ell}_{i,k} \otimes_{\nabla} \tilde{r}_{k,p+1} \right\} \right] \oslash_{\Delta} \tilde{r}_{p+1,p+1}$$

et

$$\alpha_G = \alpha_D$$

On peut énoncer :

Proposition VI. 5.

La suite \tilde{O}_D est une D - suite maximisante de \mathcal{G} sur son domaine de définition.

La suite \tilde{O}_G est une G - suite maximisante de \mathcal{G} sur son domaine de définition si la condition suivante est réalisée :

$$\tilde{r}_{pp} > 0 \quad \forall p$$

Soit par exemple, à établir que \tilde{O}_D est une D - suite maximisante de \mathcal{G} : Désignant par $O = \{\otimes\}$ une suite d'opérations correctes sur S correspondant à \mathcal{G} et donnant pour image de A, \bar{L} et \bar{R} , il suffit de raisonner par récurrence :

Supposant que les p premières lignes de R, \tilde{R} et \bar{R} et les p premières colonnes de L, \tilde{L} et \bar{L} vérifient les inégalités suivantes :

$$\begin{aligned} r_{ij} &\leq \tilde{r}_{ij} \leq 0 & , & & \bar{r}_{ij} &\leq \tilde{r}_{ij} & , & & 1 &\leq i &\leq p & , & i+1 &\leq j &\leq n \\ 0 &< r_{ii} &\leq \tilde{r}_{ii} & , & \bar{r}_{ii} &\leq \tilde{r}_{ii} & , & & 1 &\leq i &\leq p & & & & & \\ \bar{r}_{ij} &\leq \tilde{r}_{ij} \leq 0 & , & & \bar{r}_{ij} &\leq \tilde{r}_{ij} & , & & 1 &\leq j &\leq p & , & j+1 &\leq i &\leq n , \end{aligned}$$

on établit que ces inégalités sont encore vérifiées pour les $(p+1)^e$ lignes de R, \tilde{R} et \bar{R} et les $(p+1)^e$ colonnes de L, \tilde{L} et \bar{L} : Ainsi, de :

$$\bar{r}_{p+1,k} \bar{r}_{k,i} \geq \tilde{r}_{p+1,k} \tilde{r}_{k,i} > 0 \quad , \quad 1 \leq k \leq p$$

on déduit par la monotonie de l'arrondi :

$$\oplus \sum_{k=1}^p \bar{r}_{p+1,k} \otimes \bar{r}_{k,i} \geq \oplus \sum_{k=1}^p \tilde{r}_{p+1,k} \otimes \tilde{r}_{k,i}$$

ce qui entraîne, $\{\otimes_j = \otimes_j\}_\nabla$ étant une G - suite maximisante pour le produit scalaire :

$$\oplus \sum_{k=1}^p \bar{r}_{p+1,k} \otimes \bar{r}_{k,i} \geq \oplus \nabla \sum_{k=1}^p \tilde{r}_{p+1,k} \otimes \nabla \tilde{r}_{k,i}$$

Il vient alors :

$$\begin{aligned} \bar{r}_{p+1,i} &\leq \tilde{r}_{p+1,i} \leq 0 & , & & p+2 &\leq i &\leq n \\ \bar{r}_{p+1,p+1} &\leq \tilde{r}_{p+1,p+1} \end{aligned}$$

et, sous la condition $\bar{r}_{p+1,p+1} > 0$ (nous ne considérons que les suites O pour lesquelles cette condition $\bar{r}_{pp} > 0 \quad \forall p$ est satisfaite) :

$$\bar{r}_{i,p+1} \leq \tilde{r}_{i,p+1} \leq 0 \quad , \quad p+2 \leq i \leq n .$$

Ces inégalités sont vérifiées pour la première ligne de R , \tilde{R} et \bar{R} et la première colonne de L , \tilde{L} et \bar{L} .

II. 4. 2. Résolution du système linéaire $AX = B$, A étant une M -matrice, B un second membre positif (resp. négatif)

Soient respectivement, $B = (b_i)$ et $X = (x_i)$, le second membre, donné sur S , et la solution du système linéaire $AX = B$.

La décomposition de A en un produit LR étant effectuée, le calcul de X est ramené à la résolution des deux systèmes triangulaires :

$$LY = B$$

$$RX = Y$$

définissant la suite d'opérations θ .

Soient, sur S :

la suite (θ_G) suivante, associant à L, B et R les images \tilde{Y} et \tilde{X} :

$$\tilde{y}_i = b_i \ominus_{\nabla} \left\{ \oplus_{\Delta} \sum_{j=1}^{i-1} l_{ij} \otimes_{\Delta} \tilde{y}_j \right\} \quad i = 1, \dots, n$$

$$\tilde{x}_i = \left[\tilde{y}_i \ominus_{\nabla} \left\{ \oplus_{\Delta} \sum_{j=i+1}^n r_{ij} \otimes_{\Delta} \tilde{x}_j \right\} \right] \oslash_{\nabla} r_{ii} \quad i = n, \dots, 1 \quad (\theta_G)$$

et $\theta_D = \bar{\theta}_G$.

Il vient :

Proposition VI. 6.

A étant une M - matrice,

- Si $B \geq 0$,

$(\tilde{\theta}_D, \tilde{\theta}_G)$ est une G - suite de (G, θ)

$(\tilde{\theta}_G, \tilde{\theta}_D)$ est une D - suite de (G, θ) si la condition $\tilde{r}_{pp} > 0 \forall p$, est réalisée.

- Si $B \leq 0$,

$(\tilde{\theta}_G, \tilde{\theta}_G)$ est une G - suite de (G, θ) si la condition $\tilde{r}_{pp} > 0 \forall p$, est réalisée.

$(\tilde{\theta}_D, \tilde{\theta}_D)$ est une D - suite de (G, θ) .

L'application des résultats précédents à l'algorithme de détermination de l'inverse A^{-1} de A par résolution de n systèmes linéaires, permet d'énoncer en particulier lorsque $S = T_b^S$ -les résultats précédents ne sont pas modifiés par le fait que \oplus ci-dessous est non correcte- :

Proposition VI. 6'.

Si, sur $T_b^S = T$,

- l'opération \oplus est définie par une troncature sans chiffre de garde, les nombres négatifs étant représentés avec signe et valeur absolue,

- les opérations \otimes et \oslash sont correctes et I - dirigées $\forall (x,y) \in T^2$,

alors les inverses C et \tilde{C} d'une M - matrice A, calculées par décomposition LR de A et résolution de n systèmes linéaires, sur, respectivement \mathbb{R} et T, vérifient :

$$\tilde{C} \leq C$$

Exemple de résolution de système linéaire :

Nous avons écrit ci-dessous, les solutions \tilde{X} , \bar{X} et X obtenues respectivement, par l'algorithme $(\mathcal{O}_D, \mathcal{O}_G)$, par l'arithmétique IBM 360 - 30, et par la résolution exacte sur \mathbb{R} de $AX = B$, pour la matrice A et le second membre B de l'exemple p.94 du chapitre IV :

0,9999908	0,9999927	1,000000
" 17	" 28	"
" 18	" 31	"
" 20	" 34	"
" 18	" 33	"
" 14	" 30	"
" 26	" 39	"
" 18	" 33	"
\tilde{X}	\bar{X}	X

II. 4. 3. Remarque : Algorithme classique de Gauss appliqué à la résolution de $AX = B$, A M - matrice, $B \geq 0$.

La résolution, sur \mathbb{R} , du système linéaire $AX = B$
 - par la méthode précédente,
 - ou par l'algorithme classique de Gauss, consistant à former 2 suites $\{A^k\}$ et $\{B^k\}$ ($1 \leq k \leq n$) de matrices et de vecteurs pour remplacer le système donné par le système triangulaire $RX=Y$ détermine dans les 2 cas l'unique solution du système donné.

Il n'en est pas de même pour la résolution sur S , suivant une même arithmétique, de $AX = B$ par les 2 méthodes considérées, lorsque, comme il en est l'usage, le calcul des éléments $r_{p+1,i}$ et $l_{i,p+1}$ inclut

le calcul des produits scalaires $\oplus \sum_k l_{p+1,k} \otimes r_{k,i}$ et

$$\oplus \sum_k l_{i,k} \otimes r_{k,p+1} .$$

Nous avons étudié la résolution de $AX = B$, A M - matrice, $B \geq 0$, relativement à l'algorithme de décomposition LR de A .
 Pour l'algorithme classique de Gauss, on aura immédiatement :

Désignant par \mathcal{G}' la suite d'opérations :

$$a_{ij}^{k+1} = a_{ij}^k - \frac{a_{ik}^k}{a_{kk}^k} a_{kj}^k \quad i, j = k + 1, \dots, n$$

$$b_i^{k+1} = b_i^k - \frac{a_{ik}^k}{a_{kk}^k} b_k^k \quad i = k + 1, \dots, n$$

$$k = 1, \dots, n - 1$$

et définissant la suite (\tilde{O}'_1) d'opérations sur S , suivante :

$$\tilde{a}_{ij}^{k+1} = \tilde{a}_{ij}^k \ominus_{\Delta} \left(\tilde{a}_{ik}^k \oslash_{\Delta} \tilde{a}_{kk}^k \right) \otimes_{\nabla} \tilde{a}_{kj}^k \quad (\tilde{O}'_1)$$

$$\tilde{b}_i^{k+1} = \tilde{b}_i^k \ominus_{\nabla} \left(\tilde{a}_{ik}^k \oslash_{\Delta} \tilde{a}_{kk}^k \right) \otimes_{\Delta} \tilde{b}_k^k$$

ainsi que $\tilde{O}'_2 = \tilde{O}'_1$, il vient de même que précédemment :

\tilde{O}'_1 est une suite maximisante de \mathcal{G}' sur son domaine de définition

\tilde{O}'_2 est une suite maximisante de \mathcal{G}' sur son domaine de définition, si la condition de stabilité $\tilde{a}_{kk}^k > 0$ est réalisée,

et les images (R, Y) et (\tilde{R}, \tilde{Y}) de (A, B) par \mathcal{G}' et $O' = \tilde{O}'_1 \vee \tilde{O}'_2$

vérifient :

$$\begin{aligned} R &\leq \tilde{R} \wedge \tilde{Y} \leq Y && \text{si } O' = \tilde{O}'_1 \\ R &\geq \tilde{R} \wedge \tilde{Y} \geq Y && \text{si } O' = \tilde{O}'_2 \end{aligned}$$

et aussi :

Dans les hypothèses relatives à l'arithmétique flottante de la proposition VI.6'. l'inégalité :

$$\tilde{C} \leq c$$

est vérifiée lorsque \tilde{C} est calculée sur T par l'algorithme de Gauss classique .

où l'opération \otimes_I a la définition : $x \otimes_I y = \begin{cases} x \otimes_{\Delta} y & \text{si } x * y \geq 0 \\ x \otimes_{\Lambda} y & \text{si } x * y < 0 \end{cases}$

On peut énoncer, relativement à $\tilde{0}_I$, et plus généralement en excluant l'hypothèse de correction de la soustraction arithmétique.

Proposition VI.7.

Si, sur S :

- l'opération \oplus est E - dirigée $\forall (x,y) \in S^2_{\wedge} \ x y < 0$, et correcte et I - dirigée $\forall (x,y) \in S^2_{\wedge} \ x y \geq 0$,
- les opérations \otimes et \oslash sont correctes et I - dirigées $\forall (x,y) \in S^2$, alors les inverses C et \tilde{C} d'une matrice tridiagonale, symétrique définie positive, obtenues par l'algorithme de Gauss, sur, respectivement R et S , vérifient :

$$|\tilde{c}_{ij}| \leq |c_{ij}| \quad i, j = 1, \dots, n$$

Démonstration

On vérifie successivement les inégalités suivantes :

$$\tilde{\alpha}_i \geq \alpha_i \quad i = 2, \dots, n$$

$$|\tilde{\epsilon}_i| \leq |\epsilon_i| \quad i = j + 1, \dots, n$$

L'inégalité :

$$|\tilde{x}_i| \leq |x_i| \quad i = n, \dots, 1$$

s'établit alors par récurrence sur i , en remarquant que pour $i = n - 1, \dots, j$, les nombres ϵ_i et $b_i x_{i+1}$ sont de signes contraires (la récurrence permet d'établir que \tilde{x}_i a le signe $\tilde{\epsilon}_i$ qui reste toujours celui de ϵ_i , les opérations I - dirigées étant supposées correctes, et on tient compte de l'égalité $\epsilon_{i+1} = -(b_i / \alpha_i) \epsilon_i$ où $\alpha_i > 0$).

Remarques

1) Sous la condition de stabilité $\tilde{\alpha}_i > 0 \ \forall i$, $\tilde{0}_E = \tilde{0}_I$ donne pour image de A , \tilde{C} vérifiant , $\forall i, j : |\tilde{c}_{ij}| \geq |c_{ij}|$.

2) La proposition VI.7 est encore valable avec les hypothèses suivantes sur A :

A est tridiagonale, ses mineurs principaux sont positifs et tout couple d'éléments de A, disposés symétriquement par rapport à la diagonale principale, est formé d'éléments de même signe.

Exemple : Les inverses \tilde{C} et C de la matrice A définie par :
 $a_i = 3$ ($1 \leq i \leq 6$), $b_i = 1$ ($1 \leq i \leq 5$), sont respectivement :

$$C = \frac{1}{377} \begin{vmatrix} 144 & -55 & 21 & -8 & 3 & -1 \\ -55 & 165 & -63 & 24 & -9 & 3 \\ 21 & -63 & 168 & -64 & 24 & -8 \\ -8 & 24 & -64 & 168 & -63 & 21 \\ 3 & -9 & 24 & -63 & 165 & -55 \\ -1 & 3 & -8 & 21 & -55 & 144 \end{vmatrix} \quad [6 \text{ p.152}]$$

$$\tilde{C} = \frac{1}{377} \begin{vmatrix} 144 & -54,99995 & 20,99997 & -7,999987 & 2,999993 & -0,9999985 \\ -54,99995 & 164,9999 & -62,99991 & 23,99995 & -8,999985 & 2,999995 \\ 20,99997 & -62,99994 & 167,9998 & -63,99994 & 23,99995 & -7,999990 \\ -7,999989 & 23,99995 & -63,99994 & 167,9999 & -62,99994 & 20,99997 \\ 2,999993 & -8,999984 & 23,99995 & -62,99994 & 164,9999 & -54,99995 \\ -0,9999979 & 2,999993 & -7,999987 & 20,99997 & -54,99995 & 143,9999 \end{vmatrix}$$

II . 6 . INVERSION PAR LA METHODE DES SOUS-MATRICES, D'UNE MATRICE A BLOCS DIAGONAUX M-MATRICES, BLOCS EXTRA-DIAGONAUX POSITIFS.

Soient A une matrice d'ordre n, d'inverse Z, et un partitionnement des deux matrices effectué selon le schéma ci-dessous :

$$A = \begin{vmatrix} B & S \\ R & C \end{vmatrix} \begin{matrix} \left. \vphantom{\begin{matrix} B \\ R \end{matrix}} \right\} p \\ \left. \vphantom{\begin{matrix} S \\ C \end{matrix}} \right\} m \end{matrix} \qquad Z = \begin{vmatrix} D & U \\ T & E \end{vmatrix} \begin{matrix} \left. \vphantom{\begin{matrix} D \\ T \end{matrix}} \right\} p \\ \left. \vphantom{\begin{matrix} U \\ E \end{matrix}} \right\} m \end{matrix}$$

Rappelons [9] , qu'avec les notations :

$$Q = C - R B^{-1} S \quad F = B^{-1} S \quad G = R B^{-1}$$

l'algorithme d'obtention de Z , sous la condition d'inversibilité de B et de Q , est :

$$\begin{aligned} E &= Q^{-1} \\ T &= - E G \\ U &= - F E \\ D &= B^{-1} - F T \end{aligned} \tag{VI.4}$$

Définissant, si $A = (a_{ij})$ et $B = (b_{ij})$ sont 2 matrices éléments de $\mathcal{M}_{n,n}(S)$, les matrices $A \oplus_{\nabla} B$, $A \oplus_{\Delta} B$, $A \otimes_{\nabla} B$, $A \otimes_{\Delta} B$ par les égalités :

$$\begin{aligned} A \oplus_{\nabla} B &= (a_{ij} \oplus_{\nabla} b_{ij}) & A \oplus_{\Delta} B &= (a_{ij} \oplus_{\Delta} b_{ij}) \\ A \otimes_{\nabla} B &= \left(\oplus_{\nabla} \sum_{k=1}^n a_{ik} \otimes_{\nabla} b_{kj} \right) & A \otimes_{\Delta} B &= \left(\oplus_{\Delta} \sum_{k=1}^n a_{ik} \otimes_{\Delta} b_{kj} \right) \end{aligned}$$

on peut énoncer par exemple , $\tilde{\mathcal{O}}_D$ et $\tilde{\mathcal{O}}_G$ ayant les définitions du § II.4 :

Proposition VI.8.

Si dans le partitionnement de A les conditions suivantes sont réalisées :

R et S matrices positives, B et Q M-matrices, alors, une G - suite de la suite d'opérations définissant (D , E , -T , -U) sur R par (VI.4) est :

$$\begin{aligned} \tilde{B}^{-1} &= \tilde{\mathcal{J}}_{\tilde{\mathcal{O}}_D, \tilde{\mathcal{O}}_G} (B) \\ \tilde{Q} &= C \oplus_{\Delta} R \otimes_{\nabla} B^{-1} \otimes_{\nabla} S \quad \tilde{F} = \tilde{B}^{-1} \otimes_{\nabla} S \quad \tilde{G} = R \otimes_{\nabla} B^{-1} \\ \tilde{E} &= \tilde{\mathcal{J}}_{\tilde{\mathcal{O}}_D, \tilde{\mathcal{O}}_G} (\tilde{Q}) \\ \tilde{T} &= - \tilde{E} \otimes_{\nabla} \tilde{G} \\ \tilde{U} &= - \tilde{F} \otimes_{\nabla} \tilde{E} \\ \tilde{D} &= \tilde{B}^{-1} \oplus_{\nabla} \tilde{F} \otimes_{\nabla} (-\tilde{T}) \end{aligned}$$

Un cas particulier où Q a la propriété de M -matrice est donné par le :

Lemme VI . 1

Si les conditions suivantes sont réalisées :

- la sous-matrice B de A est diagonale
- la matrice A est diagonalement dominante et vérifie :

$$a_{ij} \leq 0 \quad i \neq j \quad \wedge \quad p+1 \leq i, j \leq n$$

$$a_{ij} \geq 0 \quad \text{sinon ,}$$

alors Q définie en (VI.4) est une M -matrice.

Le lemme VI.1 s'établit en montrant que Q vérifie les deux propriétés suivantes :

- 1) $q_{ij} \leq 0 \quad i \neq j$
- 2) Q est à diagonale positive dominante.

Ces 2 propriétés impliquent alors pour Q la propriété de M -matrice [7].

L'inégalité $q_{ij} \leq 0$ pour $i \neq j$ vient immédiatement.

D'autre part, de $q_{ij} = c_{ij} - \sum_{k=1}^p \frac{r_{ik} s_{kj}}{b_k}$ on déduit :

$$q_{ii} - \sum_{j \neq i} |q_{ij}| = c_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^m |c_{ij}| - \sum_{j=1}^m \sum_{k=1}^p \frac{r_{ik} s_{kj}}{b_k} =$$

$$c_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^m |c_{ij}| - \sum_{k=1}^p \frac{|r_{ik}|}{b_k} \sum_{j=1}^m |s_{kj}|$$

La matrice A étant à diagonale positive dominante :

$$b_k > \sum_{j=1}^m |s_{kj}| \quad k = 1, \dots, p$$

donc :

$$q_{ii} - \sum_{j \neq i} |q_{ij}| > c_{ii} - \left(\sum_{\substack{j=1 \\ j \neq i}}^m |c_{ij}| + \sum_{k=1}^p |r_{ik}| \right) > 0 \quad \text{par hypothèse, d'où le}$$

résultat.

Relativement aux systèmes de nombres à virgule flottante, nous énoncerons :

Proposition VI . 9 .

Si

- 1) les conditions du lemme VI.1 sont réalisées
- 2) $S = T_b^S$, avec l'arithmétique suivante :
 - l'opération \oplus est définie par une troncature sans chiffre de garde, les nombres négatifs étant représentés avec signe et valeur absolue,
 - les opérations \otimes et \oslash sont correctes et I - dirigées $\forall (X,Y) \in T^2$,
 alors, les inverses Z et \tilde{Z} de A , calculées par la méthode des sous-matrices, sur, respectivement \mathbb{R} et T_b^S , vérifient :

$$|\tilde{z}_{ij}| \leq |z_{ij}| \quad i,j = 1, \dots, n$$

Cette proposition vient immédiatement : on remarquera en particulier que pour l'arithmétique considérée, l'hypothèse $c_{ij} \leq 0 \quad i \neq j$ est nécessaire pour conclure à $\tilde{q}_{ij} \geq q_{ij} \quad \forall i,j$.

II . 7 . REMARQUES

1) Les suites maximisantes d'opérations sur S que nous avons mises en évidence aux paragraphes II.1 à II.6 précédents, exclusion faite du paragraphe II.3 , définissent 2 schémas de calcul identiques aux schémas générés, pour les algorithmes considérés, par l'arithmétique d'intervalle définie par R.E. MOORE [29] :

a, b, c, d étant éléments de S :

$$[a,b] \oplus [c,d] = [\vee(a+c) , \Delta(b+d)]$$

$$[a,b] \ominus [c,d] = [\vee(a-d) , \Delta(b-c)]$$

$$[a,b] \otimes [c,d] = [\vee(\min\{ac,ad,bc,bd\}) , \Delta(\max\{ac,ad,bc,bd\})]$$

et si $0 \notin [c,d]$:

$$[a,b] \oslash [c,d] = [\vee(\min\{a/c,a/d,b/c,b/d\}) , \Delta(\max\{a/c,a/d,b/c,b/d\})]$$

Pour les algorithmes cités ci-dessus, les bornes, de la solution, calculées par l'arithmétique d'intervalle peuvent être effectivement atteintes par le calcul sur S .

Il n'en est pas ainsi, sauf cas particulier de la valeur x_0 initiale, pour l'algorithme d'obtention du nième itéré du calcul de la racine carrée d'un nombre par la méthode de Newton.

2) Nous donnons ci-dessous un exemple très simple de non-existence de G - suite d'opérations correctes sur S correspondant à une suite donnée d'opérations sur \mathbb{R} :

Pour $S = T_{10}^2$, $r = (x y - z u) (x y - v w)$, $x = 0,57$, $y = 0,63$,
 $z = 0,58$, $u = 0,62$, $v = 0,65$, $w = 0,55$

les valeurs \bar{r} calculées par les différentes suites correctes d'opérations sur T_{10}^2 ont été :

0 et 0,0001

alors que $r = - 0,0000008$.

III - UN EXEMPLE DE SUITES NON MAXIMISANTES D'OPERATIONS SUR S GENERANT DES IMAGES ENCADRANT LA SOLUTION REELLE.

$N, M \in \mathbb{N}$, étant donnés un vecteur x de S^N et une suite \mathcal{O} d'opérations sur \mathbb{R} , faisant correspondre à x un vecteur $r = (r_i)$ de \mathbb{R}^M , introduisons la définition suivante :

Définition VI.4.

Deux suites \mathcal{O}^1 et \mathcal{O}^2 d'opérations sur S , correspondant à \mathcal{O} , sont dites suites d'encadrement de \mathcal{O} sur la partie \mathcal{P} de S^N , si $\forall x \in \mathcal{P}$, les images r, r^1 et r^2 de x par $\mathcal{O}, \mathcal{O}^1$ et \mathcal{O}^2 satisfont :

$$r_i^1 \leq r_i \leq r_i^2 \vee r_i^1 \geq r_i \geq r_i^2 \quad i = 1, \dots, M$$

La définition précédente s'applique lorsque \mathcal{O} fait correspondre à x de S^N , une matrice de M lignes et P colonnes, par l'isomorphisme de $\mathcal{M}_{M,P}(\mathbb{R})$ sur $\mathbb{R}^{M \times P}$.

III.1. RESOLUTION D'UN SYSTEME TRIANGULAIRE - NOTATIONS

Soient $A = (a_{ij})$ une matrice triangulaire d'ordre n , non singulière, donnée sur S , et $b = (b_i)$ un vecteur donné de S^n . Supposons par exemple, A triangulaire inférieure, et considérons la résolution du système linéaire :

$$A x = b \quad (\text{VI.5})$$

et la suite \mathcal{O} d'opérations sur \mathbb{R} définissant x .

Désignons respectivement par $x = (x_i)$ la solution exacte de (VI.5), $\bar{x} = (\bar{x}_i)$ la solution calculée sur S , et $\epsilon = (\epsilon_i)$ le vecteur erreur $x - \bar{x}$.

Définissant pour $i = 1, \dots, n$, les nombres e_i et $\bar{\bar{x}}_i$ suivants :

$$\bar{\bar{x}}_i = (b_i - \sum_{j=1}^{i-1} a_{ij} \bar{x}_j) / a_{ii} \quad 1 \leq i \leq n$$

$$e_i = \bar{\bar{x}}_i - \bar{x}_i \quad 1 \leq i \leq n$$

(e_i est l'erreur "générée au rang i " et soit $e = (e_i)$),

rappelons l'égalité suivante vérifiée par \bar{x} ([8], [42]) :

$$A \bar{x} = b - (\text{Diag } A) e$$

soit encore avec la notation $A^{-1} = C = (c_{ij})$:

$$\epsilon_i = \sum_{j=1}^i c_{ij} a_{jj} e_j \quad (\text{VI.6})$$

III.2. DETERMINATION DE SUITES D'ENCADREMENT DE \mathcal{O}

III . 2 . 1 L'erreur e_i

L'expression de e_i en fonction des erreurs élémentaires de multiplication, d'addition et de division, entachant le calcul de \bar{x}_i sur S , a été rappelée au § II.2.2 du chapitre III.

\mathcal{O}_i étant la suite d'opérations sur \mathbb{R} définissant \bar{x}_i par :

$$\bar{x}_i = \left[b_i - \left(\sum_{j=1}^{i-1} a_{ij} \bar{x}_j \right) \right] / a_{ii}$$

deux suites maximisantes d'encadrement de \mathcal{O}_i sont respectivement générées par les choix suivants d'opérations sur S :

$$C_E = \left\{ \otimes_{\Delta}, \oplus_{\Delta}, \ominus_{\nabla}, \emptyset = \begin{cases} \emptyset_{\nabla} & \text{si } a_{ii} > 0 \\ \emptyset_{\Delta} & \text{si } a_{ii} < 0 \end{cases} \right\}, C_{E'} = \bar{C}_E$$

$$(C_E \rightarrow a_{ii} e_i > 0 ; C_{E'} \rightarrow a_{ii} e_i \leq 0)$$

III . 2 . 2 . Théorème

Proposition VI . 10

Si dans la résolution d'un système triangulaire $A x = b$ donné sur S , la matrice inverse de A , $C = A^{-1}$ est de la forme :

$$C = D P \Delta$$

où :

P est une matrice non négative

$D = (d_i)$ et $\Delta = (\delta_i)$ sont 2 matrices diagonales d'éléments égaux à 1 ou -1

alors, deux suites \mathcal{O}^1 et \mathcal{O}^2 d'encadrement de \mathcal{O} sont définies par les choix C_i^1 et C_i^2 suivants, d'opérations sur S pour le calcul de \bar{x}_i , $i = 1, \dots, n$:

$$C_i^1 = \begin{cases} C_E & \text{si } \delta_i > 0 \\ C_{E'} & \text{si } \delta_i < 0 \end{cases}, C_i^2 = \bar{C}_i^1$$

En effet, avec la notation $P = (p_{ij})$, (VI.6) s'écrit :

$$\varepsilon_i = \sum_{j=1}^i d_i p_{ij} \delta_j a_{jj} e_j = d_i \sum_{j=1}^i p_{ij} \delta_j a_{jj} e_j$$

ce qui avec :

$$\delta_j a_{jj} e_j \geq 0 \quad \forall j \quad \text{pour } 0^1$$

$$\delta_j a_{jj} e_j \leq 0 \quad \forall j \quad \text{pour } 0^2$$

implique la proposition.

III . 2 . 3 Exemple et contre-exemple

1) Considérons la résolution sur T_{10}^2 du système triangulaire :

$$\begin{vmatrix} 3 & & & & \\ -4 & 5 & & & \\ 1 & -2 & 1 & & \\ -2 & 1 & -1 & 1 & \end{vmatrix} \begin{vmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{vmatrix} = \begin{vmatrix} 1 \\ 0 \\ 0 \\ 0 \end{vmatrix}$$

La condition d'application de la proposition VI.10 est vérifiée par la matrice A d'inverse A^{-1} non négative.

Les 2 suites 0^1 et 0^2 d'encadrement de 0 , définies par cette proposition, calculent respectivement les solutions approchées :

$$\bar{x}^1 = (0,33 ; 0,26 ; 0,19 ; 0,59) \quad \text{et} \quad \bar{x}^2 = (0,34 ; 0,28 ; 0,22 ; 0,62)$$

de la solution exacte du système $x = (0,3333.. ; 0,2666.. ; 0,2 ; 0,6)$

On remarquera que 0^1 et 0^2 ne sont pas maximisantes relativement au calcul de la composante x_3 :

La résolution sur T_{10}^2 , du système donné, par l'ensemble des suites correctes d'opérations, s'explicite en effet suivant :

$\bar{x}_1 = 1 \textcircled{1} 3$	0,33	0,34
$\bar{x}_2 = (4 \otimes \bar{x}_1) \textcircled{1} 5$	0,26 0,28	0,26 0,28
$\bar{x}_3 = -\bar{x}_1 \oplus 2 \otimes \bar{x}_2$	0,19 0,23	0,18 0,22
$\bar{x}_4 = 2 \otimes \bar{x}_1 \ominus \bar{x}_2 \oplus \bar{x}_3$	0,59 0,61	0,60 0,62

définissant 4 solutions approchées différentes \bar{x} de x .

On constate que \bar{x}^1 satisfait à :

$$\|x - \bar{x}^1\| = \underset{\bar{x}}{\text{Min}} \|x - \bar{x}\|$$

pour les 3 normes de Holder usuelles.

2) \mathcal{O} étant la suite d'opérations sur \mathbb{R} définie par :

$$\begin{vmatrix} 7 \\ 9 & -2 \\ -2 & 1 & -1 \\ 1 & 0 & 0 & -1 \\ 0 & 2 & -3 & 5 & -1 \end{vmatrix} \begin{vmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{vmatrix} = \begin{vmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{vmatrix}$$

il n'existe pas, dans l'ensemble des suites correctes d'opérations sur T_{10}^2 correspondant à \mathcal{O} , de suites d'encadrement de \mathcal{O} .

$$\begin{array}{l} \bar{x}_1 = 1 \textcircled{1} 7 \\ \bar{x}_2 = (9 \otimes \bar{x}_1) \textcircled{1} 2 \\ \bar{x}_3 = -2 \otimes \bar{x}_1 \oplus \bar{x}_2 \\ \bar{x}_4 = \bar{x}_1 \\ \bar{x}_5 = 2 \otimes \bar{x}_2 \ominus 3 \times \bar{x}_3 \oplus 5 \otimes \bar{x}_4 \end{array} \quad \begin{array}{c} 0,14 \\ / \quad \backslash \\ 0,60 \quad 0,65 \\ | \quad | \\ 0,32 \quad 0,37 \\ | \quad | \\ 0,14 \quad 0,14 \\ | \quad / \quad \backslash \\ 0,94 \quad 0,90 \quad 0,80 \end{array} \quad \begin{array}{c} 0,15 \\ / \quad \backslash \\ 0,65 \quad 0,70 \\ | \quad | \\ 0,35 \quad 0,40 \\ | \quad | \\ 0,15 \quad 0,15 \\ / \quad | \quad \backslash \\ 0,95 \quad 1 \quad 1,1 \quad 0,95 \end{array}$$

la solution exacte du système étant :

$$x = (0,1428.. ; 0,6428.. ; 0,3571.. ; 0,1428.. ; 0,9285..)$$

III . 2. 4. Une autre formulation des hypothèses

Une forme différente de (VI.6) peut être donnée en exprimant les éléments de C en fonction des mineurs de A .

Notant Δ_{ij} le mineur correspondant à a_{ij} , $D(A)$ le déterminant de A ,

et $D(A_{\{j+1, \dots, i\}, \{j, \dots, i-1\}})$ le déterminant de la sous-matrice de A

formée par les lignes $j+1, \dots, i$ et les colonnes $j, \dots, i-1$

pour $j \leq i - 1$:

Si dans la résolution d'un système triangulaire $A x = b$ donné sur S , la matrice quasi-triangulaire d'ordre $n-1$ définie par :

$$M_{n-1} = \begin{pmatrix} a_{ij} \\ a_{ii} \end{pmatrix} ; i = 2, \dots, n ; j = 1, \dots, n-1$$

est de la forme :

$$M_{n-1} = D' P' \Delta'$$

où

$P' = (p'_{ij})$ est une matrice vérifiant $D(P'_{\{p, \dots, q\}, \{p, \dots, q\}}) > 0 \quad 1 \leq p, q \leq n-1$

$D' = (d'_i)$ et $\Delta' = (\delta'_i)$ sont 2 matrices diagonales d'éléments égaux à 1 ou -1

alors, deux suites O^1 et O^2 d'encadrement de \mathcal{O} sont définies par les choix C_i^1 et C_i^2 suivants, d'opérations sur S pour le calcul de \bar{x}_i , $i = 1, \dots, n$:

$$C_1^1 = C_G$$

$$C_i^1 = \begin{cases} C_G & \text{si } (-1)^{i-1} \prod_{j=1}^{i-1} d'_j \delta'_j > 0 \\ C_D & \text{" " " " } < 0 \end{cases}, \quad C_i^2 = \bar{C}_i^1$$

III . 2 . 5 . Cas où les suites d'encadrement de \mathcal{O} sont maximisantes.

Un cas particulier où les 2 suites d'encadrement de \mathcal{O} définies précédemment sont maximisantes, est donné par la :

Proposition VI . 11

Si tous les termes du développement de $D(P'_{\{p, \dots, q\}, \{p, \dots, q\}})$ ($1 \leq p, q \leq n-1$) sont de même signe, alors O^1 et O^2 sont 2 suites maximisantes de \mathcal{O} .

Il suffit de raisonner pour $M_{n-1} = P'$.

$\bar{x} = (\bar{x}_i)$ désignant une solution du système calculée sur S par une suite O d'opérations correctes correspondant à \mathcal{O} , la double inégalité :

$$\begin{aligned} \bar{x}_i^1 < \bar{x}_i < \bar{x}_i^2 & \quad i \text{ impair} \\ \bar{x}_i^1 \geq \bar{x}_i \geq \bar{x}_i^2 & \quad i \text{ pair} \end{aligned} \quad (\text{VI.8})$$

s'établit par récurrence sur l'indice i .

Pour $i = 1$, (VI.8) est satisfaite.

Pour $i = 2$, les hypothèses impliquent $\frac{a_{21}}{a_{22}} > 0$ ce qui, avec les définitions de 0^1 et 0^2 , assure (VI.8).

Supposant alors :

$$\left[\begin{array}{l} \frac{a_{ij}}{a_{ii}} \text{ du signe de } (-1)^{i+j-1} \quad (2 \leq i \leq l \wedge 1 \leq j \leq i-1) \\ \text{(VI.8) vérifiée } \forall i \leq l \end{array} \right.$$

on vérifie que ces 2 propriétés sont encore satisfaites pour $i = l + 1$.

Le développement, par rapport à sa dernière ligne de $D(M_l)$:

$$D(M_l) = (-1)^{l+1} \frac{a_{l+1,1}}{a_{l+1,l+1}} + \sum_{j=2}^l (-1)^{l+j} \frac{a_{l+1,j}}{a_{l+1,l+1}} D(M_{j-1})$$

assure la 1ère propriété.

Ce résultat, les définitions de 0^1 et 0^2 et la 2ème hypothèse de récurrence, impliquant alors (VI.8) pour $i = l + 1$.

Remarque : La matrice triangulaire A satisfaisant $M_{n-1} = \begin{pmatrix} a_{ij} \\ a_{ii} \end{pmatrix} = P'$ avec tous les termes du développement de $D(P'_{\{p, \dots, q\}, \{p, \dots, q\}})$ ($1 \leq p, q \leq n-1$) de même signe, vérifie :

$$Z = D A \Delta$$

où :

Z est une M - matrice triangulaire

$D = (d_i)$ et $\Delta = (\delta_i)$ sont 2 matrices diagonales d'élément égaux à :

$$d_i = (-1)^i \times \text{signe}(a_{ii}) \quad \delta_i = (-1)^i$$

Les 2 suites O^1 et O^2 sont maximisantes dans le cas où A est une M -matrice, résultat que nous rapprocherons des résultats énoncés au § II.4 (Résolution d'un seul système triangulaire donné sur S à second membre quelconque dans ce paragraphe, et résolution successive, après décomposition LR d'une matrice donnée sur S , de deux systèmes triangulaires à second membre de signe constant dans le § II.4.2).

III . 2 . 6 . Exemple numérique

Nous avons considéré la résolution du système $A x = b$ pour :

$$A = \begin{array}{c|cccccccc} 3 & & & & & & & & \\ 7 & 5 & & & & & & & \\ 1 & 5 & 3 & & & & & & \\ 1 & 1 & 12 & 7 & & & & & \\ 3 & 2 & 2 & 16 & 7 & & & & \\ 3 & 5 & 5 & 1 & 20 & 5 & & & \\ 1 & 2 & 1 & 3 & 2 & 15 & 3 & & \\ 2 & 0 & 1 & 3 & 0 & 1 & 10 & 1 & \end{array} \quad b = \begin{array}{c|c} 1 \\ 4 \\ 3 \\ 7 \\ 10 \\ 13 \\ 9 \\ 6 \end{array}$$

Les éléments de A vérifient

$$a_{i,i-1} > \sum_{\substack{j=1 \\ j \neq i-1}}^i |a_{ij}| \quad i = 2, \dots, n$$

Par le théorème de Gerschgorin-Hadamard, les valeurs propres réelles de toute sous-matrice $M_{\{p,\dots,q\},\{p,\dots,q\}}$ ($1 \leq p, q \leq n-1$) de M_{n-1} , sont positives.

M_{n-1} est une matrice de type P' .

La solution exacte X du système, et les solutions X^1 et X^2 calculées sur T_{16}^6 par les 2 suites d'opérations O^1 et O^2 sont respectivement (\bar{X} est la solution calculée par l'arithmétique IBM 360) :

$X =$	$\begin{array}{ l} 1/3 \\ 1/3 \\ 1/3 \\ 1/3 \\ 1/3 \\ 1/3 \\ 1/3 \\ 1/3 \end{array}$	$X^1 =$	$\begin{array}{ l} 0,3333333 \\ " 34 \\ " 27 \\ " 47 \\ " 03 \\ " 464 \\ " 2682 \\ " 9691 \end{array}$	$X^2 =$	$\begin{array}{ l} 0,3333334 \\ " 2 \\ " 40 \\ " 22 \\ " 59 \\ " 225 \\ " 871 \\ 28085 \end{array}$	$\bar{X} =$	$\begin{array}{ l} 0,3333333 \\ " 4 \\ " 6 \\ " 0 \\ " 41 \\ " 05 \\ " 476 \\ " 1947 \end{array}$
-------	--	---------	--	---------	---	-------------	---

IV - THEOREMES GENERAUX

Les théorèmes qui suivent sont relatifs à l'ensemble des 2^{N-1} images que les suites d'opérations correctes sur S , correspondant à la suite $\mathcal{O} = \{ \star_i, i=1, \dots, N-1 \}$ d'opérations sur \mathbb{R} , associent à une même suite finie de N opérandes donnés sur S .

Dans le cas de la sommation, c'est-à-dire $\mathcal{O} = \{ +, i=1, \dots, N-1 \}$, on peut énoncer :

Proposition VI . 12 .

Etant donné le vecteur $X = (x_1, \dots, x_N)$ de $(S^+)^N$, si S est un sous-ensemble de \mathbb{R} à intervalles non décroissants, alors :

- L'ensemble E , des images de X par les 2^{N-1} suites \mathcal{O} d'additions correctes sur S , $\mathcal{O} = \{ \bigoplus_i, i = 1, \dots, N-1 ; \square_i = \nabla, \Delta \}$,

est formé d'au plus N éléments consécutifs de S .

- Lorsque $\text{Card } E = N$, le $j^{\text{ième}}$ élément s_j de E vérifie :

$$n(s_j) = \text{Card} \{ \mathcal{O} \mid i_0(X) = s_j \} = C_{N-1}^{j-1} \quad j = 1, \dots, N$$

Démonstration :

Il suffit de raisonner par récurrence.

Pour la 1^{ère} partie de la proposition VI. 12. :

. La propriété est vraie au rang 2

. Supposons la vérifiée au rang k , c'est-à-dire supposons que l'ensemble E_k des images de $X_k = (x_1, \dots, x_k)$ par les 2^{k-1} suites d'additions correctes sur S , soit formé de λ éléments consécutifs de S : $s_1^k, \dots, s_\lambda^k$ avec $\lambda \leq k$.

Lorsqu'on forme les nombres $s_i^k + x_{k+1}$, il vient :

$$\mu = \text{Card} \{ z \mid z \in [s_1^k + x_{k+1}, s_\lambda^k + x_{k+1}] \wedge z \in S \} \leq \lambda,$$

l'égalité $\mu = \lambda$ impliquant $s_1^k + x_{k+1} \in S \wedge s_\lambda^k + x_{k+1} \in S$.

Dans ce dernier cas, l'égalité $\text{Card } E_{k+1} = \lambda$, et dans le cas général, l'inégalité $\text{Card } E_{k+1} \leq \mu + 2 \leq \lambda + 1$ assurent bien la propriété au rang $k + 1$, les éléments de E_{k+1} étant, par ailleurs, consécutifs sur S .

La 2^{ème} partie de la proposition s'établit alors :

. Pour $N = 2$ et $\text{Card } E_2 = 2$, on a bien $n(s_1^2) = n(s_2^2) = 1 = C_1^0 = C_1^1$

. Supposons alors $\text{Card } E_k = k$ et, $n(s_j^k) = C_{k-1}^{j-1}$ ($j = 1, \dots, k$).

(L'égalité $\text{Card } E = N$ entraîne en effet, $\text{Card } E_k = k, \forall k$, d'après l'inégalité $\text{Card } E_{k+1} \leq \text{Card } E_k + 1$ obtenue précédemment).

Raisonnant par l'absurde, et utilisant à nouveau la propriété de non décroissance des longueurs $z_{i+1} - z_i$ pour l'ensemble des éléments z_i de S^+ classés dans l'ordre naturel, on établit : $\text{Card } E_{k+1} = k + 1$ implique qu'il existe un et un seul élément de S , s_j^{k+1} dans tout intervalle $[s_{j-1}^k + x_{k+1}, s_j^k + x_{k+1}[$, $j = 2, \dots, k$ et qu'aucun des nombres $s_j^k + x_{k+1}$ n'est élément de S .

On a alors immédiatement :

$$n(s_j^{k+1}) = n(s_{j-1}^k) + n(s_j^k) = C_{k-1}^{j-2} + C_{k-1}^{j-1} = C_k^{j-1} \quad j = 2, \dots, k$$

$$n(s_1^{k+1}) = n(s_1^k) = 1$$

$$n(s_{k+1}^{k+1}) = n(s_k^k) = 1$$

On peut, pour une somme algébrique, énoncer :

Proposition VI . 13 .

Etant donné le vecteur $X = (x_1, \dots, x_N)$ de S^N et l'ensemble E des images de X par les suites O d'additions correctes sur S , $O = \{ \bigoplus_{\Delta}^i, i = 1, \dots, N-1 ; \Delta_i = \nabla, \Delta \}$, alors

si la condition suivante est réalisée :

$$\varepsilon\left(\bigoplus_{\nabla}^i \sum_{j=1}^i x_j\right) = \varepsilon\left(\bigoplus_{\Delta}^i \sum_{j=1}^i x_j\right) = \text{constante indépendante de } i \quad \forall i$$

E est formé d'éléments consécutifs de S et, avec la notation $P = \text{Card } E$, le $j^{\text{ième}}$ élément s_j de E vérifie :

$$n(s_j) = \text{Card} \{ O \mid i_0(X) = s_j \} = 2^{N-P} C_{P-1}^{j-1} \quad j = 1, \dots, P$$

Démonstration :

Raisonnons encore par récurrence :

. La propriété est vraie pour $N = 2$

. Supposons la vérifiée au rang k , c'est-à-dire, en reprenant les notations de la démonstration précédente, supposons que E_k soit formé d'éléments consécutifs de S , et que, si $\lambda = \text{Card } E_k$, on ait :

$$n(s_j^k) = 2^{k-\lambda} C_{\lambda-1}^{j-1} \quad j = 1, \dots, \lambda$$

Les nombres $s_j^k + x_{k+1}$ ayant le même intervalle de précision que les nombres s_j^k , deux possibilités seules peuvent se présenter :

1) $s_j^k + x_{k+1} \in S \quad \forall j = 1, \dots, \lambda$

alors, $\text{Card } E_{k+1} = \lambda$ et $n(s_j^{k+1}) = 2 n(s_j^k)$

2) il existe un et un seul élément de S intérieur à tout intervalle $[s_{j-1}^k + x_{k+1}, s_j^k + x_{k+1}]$

alors, $\text{Card } E_{k+1} = \lambda + 1$ et la propriété vient immédiatement par :

$$n(s_j^{k+1}) = n(s_{j-1}^k) + n(s_j^k) \quad j = 2, \dots, \lambda$$

Il vient aussi :

Proposition VI . 14 .

Etant donné le vecteur $X = (x_1, \dots, x_N)$ de $(T_b^S)^N$, une suite O de $N-1$ additions correctes sur T_b^S et la suite \bar{O} conjuguée de O , alors si la condition suivante est réalisée :

$$\varepsilon(\bigoplus_{\nabla} \sum_{j=1}^i x_j) = \varepsilon(\bigoplus_{\Delta} \sum_{j=1}^i x_j) \leq \varepsilon(\sum_{j=1}^{i-1} x_j) \quad \forall i \geq 2$$

on peut affirmer :

$$i_O(X) + i_{\bar{O}}(X) = \text{constante indépendante de } O .$$

Avec la notation $X_k = (x_1, \dots, x_k)$, supposons qu'au rang $k < n$, toute suite O_k de $k-1$ additions correctes vérifie :

$$i_{O_k}(X_k) + i_{\bar{O}_k}(X_k) = C \text{ indépendante de } O_k$$

Les conditions énoncées dans la proposition :

$$\varepsilon(i_{O_k}(X_k)) = \alpha \text{ indépendante de } O_k$$

$$\varepsilon(i_{O_k}(X_k) + X_{k+1}) = \beta \leq \alpha \quad \forall O_k$$

impliquent que tous les nombres $i_{O_k}(X_k) + X_{k+1}$ sont les images par une translation d'amplitude e avec $0^k \leq e < \beta$ d'éléments de T_b^S .

Si $e \neq 0$, on a :

$$i_{0_k}(X_k) \oplus x_{k+1} = i_{0_k}(X_k) + x_{k+1} + \begin{cases} e & \text{si } \oplus = \oplus_{\Delta} \\ -(\beta-e) & \text{si } \oplus = \oplus_{\nabla} \end{cases}$$

et par suite :

$$i_{0_{k+1}}(X_{k+1}) + i_{0_{k+1}}(X_{k+1}) = C + 2x_{k+1} + 2e - \beta$$

Pour le produit scalaire, on démontre des théorèmes analogues aux théorèmes précédents, ainsi :

Si S est un sous-ensemble de \mathbb{R} à intervalles non décroissants, l'ensemble des images des vecteurs $X = (x_1, \dots, x_N)$ et $Y = (y_1, \dots, y_N)$ de $(S^+)^N$ par les 2^{2N-1} produits scalaires "corrects" différents, est formé d'au plus $2N$ éléments consécutifs de S .

Enfin, nous signalerons que si S est un sous-ensemble de $[-1, 1]$ à intervalles constants, la proposition VI. 12 est vérifiée pour tout vecteur de S^N lorsqu'on remplace l'opération $+$ par l'opération \times . On a aussi :

Si S est un sous-ensemble de $[-1, 1]$ à intervalles constants, si \mathcal{K} est un algorithme ne comportant que des opérations élémentaires d'addition, de soustraction ou de multiplication, l'ensemble des solutions de \mathcal{K} calculées sur S par des opérations correctes, est formée d'éléments consécutifs de S .

(Exemple : Résolution d'un système triangulaire à éléments diagonaux unité).

V - CONCLUSION

La détermination sur T_b^S de suites d'encadrement, de la suite \mathcal{G} d'opérations sur \mathbb{R} associée à un algorithme, n'est pas aisée en général. De plus, l'existence de telles suites n'est pas assurée.

Les résultats énoncés dans ce chapitre auront essentiellement leur application :

- dans la connaissance de la stabilité numérique du problème étudié, lorsque les suites d'encadrement de \mathcal{G} sont maximisantes.
- dans la connaissance, pour certaines arithmétiques flottantes courantes, des positions respectives des solutions calculées sur \mathbb{R} et T .

CHAPITRE VII

NE ETUDE SUR LA REPARTITION DE L'ERREUR D'ARRONDI

INTRODUCTION

Nous considérons dans ce chapitre des systèmes de nombres à virgule flottante.

Dans la première partie, est étudiée l'hypothèse de l'uniforme répartition de l'erreur d'arrondi élémentaire sur son domaine de définition, lorsqu'on considère cette erreur comme variable aléatoire :

En arithmétique à virgule flottante de base b et de nombre de chiffres de mantisse s , cette hypothèse généralement admise s'exprime, avec les notations suivantes :

$$x \in T_b^s, y \in T_b^s, p = e(x \star y), E = x \star y - (x \otimes y)$$

en écrivant que E , considérée comme variable aléatoire, est uniformément répartie :

- dans l'intervalle $] - b^{p-s}, b^{p-s}[$ si \otimes est correcte
- dans l'intervalle $[0, b^{p-s}[$ si $\otimes = \otimes_{\nabla}$

Nous proposons quelques remarques sur cette hypothèse, tenant compte de la répartition discrète - et non continue - des nombres.

Dans la seconde partie de ce chapitre, est mentionné le problème du meilleur choix, de la base d'une représentation des nombres en virgule flottante, relativement au critère de la précision.

I - FONCTION DE REPARTITION DE L'ERREUR D'ARRONDI CONSIDEREE COMME VARIABLE ALEATOIRE

I . 1 . FORMULATION DU PROBLEME

Nous étudierons l'hypothèse d'équirépartition sur son domaine de définition, de l'erreur d'arrondi entachant une opération élémentaire en arithmétique à virgule flottante, en considérant le problème suivant :

Etant donné un système $T_b^s = T$ de nombres à virgule flottante, une opération \oplus appartenant à l'ensemble $\{ \oplus, \ominus, \otimes, \oslash \}$, et un exposant p , on désigne par Σ l'ensemble :

$$\Sigma = \{ (x,y) \mid (x,y) \in T^2 \wedge e(x \oplus y) = p \}$$

Soit à étudier l'image \mathcal{E} de Σ par l'application de T^2 dans \mathbb{R} :

$$(x,y) \rightarrow E = x \star y - (x \oplus y) .$$

Nous traiterons ce problème pour l'addition et la soustraction arithmétiques, dans le cas de la TRONCATURE avec ou sans chiffre de garde. Cette règle, souvent utilisée, est en effet mal connue du point de vue que nous adoptons (influence des cas de cancellation, dissymétrie apportée à l'erreur par la présence du chiffre de garde dans l'addition de 2 nombres de signes contraires).

Le cas de la multiplication est abordé.

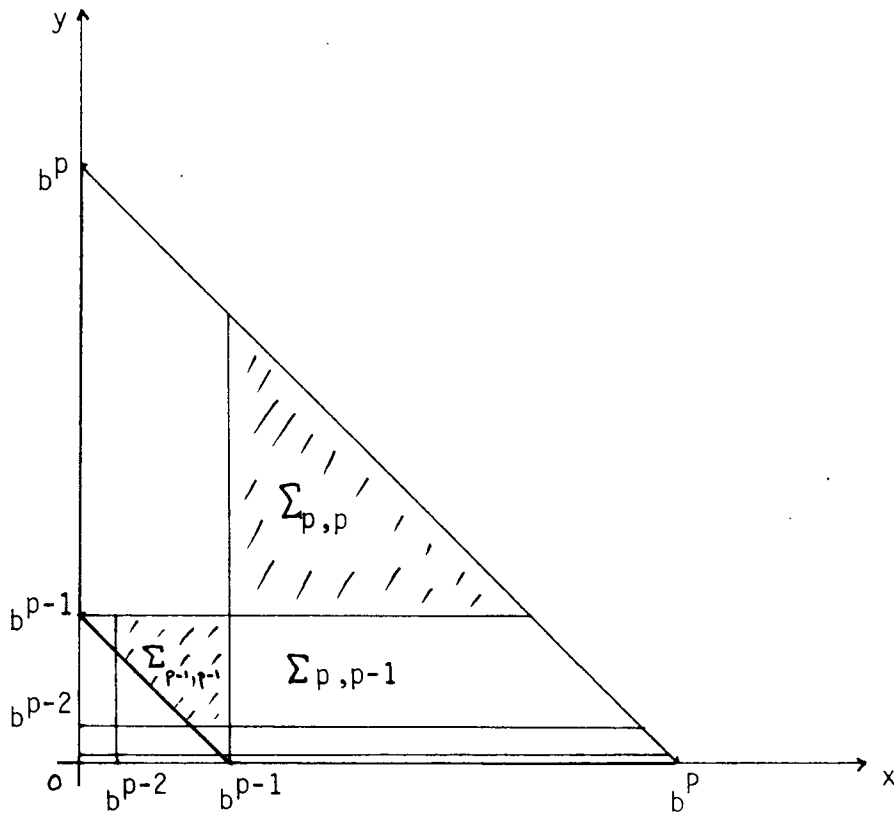
En raison des symétries, pour une représentation des nombres par signe et valeur absolue ou en notation de complément à 2, nous travaillerons avec la représentation des nombres par signe et valeur absolue, et nous considérerons uniquement les couples (x,y) tels que $x \star y$ est positif. Nous noterons respectivement Σ et \mathcal{E} les sous-ensembles de Σ et \mathcal{E} considérés.

I . 2 . L'ADDITION ARITHMETIQUE

Dans le cas d'une règle de troncature, les chiffres de garde, même s'ils existent, ne participent pas à l'opération d'addition de deux nombres de même signe. Un seul cas est donc à considérer dans ce paragraphe.

CARACTERISATION DES ENSEMBLES Σ ET \mathcal{E} .

Σ est l'ensemble des points de T^2 intérieurs au trapèze défini par les inégalités $b^{p-1} \leq x + y < b^p$ où seule la grande base est exclue.
 $x \geq 0, y \geq 0$



Pour la caractérisation de \mathcal{E} :

Posons $e(x) \geq e(y)$.

En remarquant qu'une condition nécessaire, pour que l'exposant de $x \oplus y$ prenne la valeur fixée p , est $e(x) = p \vee p - 1$, étudions ces deux possibilités :

$$1) \quad \underline{e(x) = p}$$

Représentons x et y sous la forme :

$$x = l b^{p-s}$$

$$b^{s-1} \leq l, m < b^s \quad (\text{ou } m = 0)$$

$$y = m b^{q-s}$$

E se déduit de (x, y) par :

$$E = r b^{q-s}$$

où r est défini par la division arithmétique $m = b^{p-q} m' + r \wedge r < b^{p-q}$.

La partie \sum_p de Σ à considérer ici, est l'union des $s + 1$ domaines $\sum_{p,q}$ définis,

pour $q = p, p - 1, \dots, p - s + 1$ par :

$$\sum_{p,q} = \{ (x, y) \mid (x, y) \in T^2 \wedge e(x) = p \wedge e(y) = q \wedge x + y < b^p \}$$

pour $q = p - s$ par :

$$\sum_{p,p-s} = \{ (x, y) \mid (x, y) \in T^2 \wedge e(x) = p \wedge e(y) \leq p - s \}$$

Parcourons \sum_{pq} ($p - s < q \leq p$) sur les droites $x = \text{constante}$.

Les entiers m caractérisant les ordonnées des éléments de \sum_{pq} situés sur la droite $x = \text{constante} = l b^{p-s}$ (avec $(b^s - 1) b^{p-q} > b^{s-1}$) sont les entiers successifs vérifiant l'inégalité :

$$b^{s-1} \leq m < \min(b^s, (b^s - 1) b^{p-q})$$

c'est-à-dire une suite d'entiers successifs compris entre deux multiples de b^{p-q} , le premier multiple seul étant inclus.

Le parcours de \sum_{pq} à l fixé implique donc l'obtention avec la même fréquence de chacune des valeurs de r .

L'image $\xi_{pq} = \{ r b^{q-s} \wedge r \in \mathbb{N} \wedge 0 \leq r < b^{p-q} \}$ de \sum_{pq} admet donc des éléments de valeurs équiprobables si l'on affecte une probabilité uniforme à l'ensemble des points de T^2 .

L'image ξ_p de \sum_p est l'union :

$$\left(\bigcup_{q=p}^{p-s+1} \xi_{pq} \right) \cup \{ y \mid y \in T \wedge y < b^{p-s} \}$$

$$2) \underline{e(x) = p - 1}$$

$$x = l b^{p-1-s}$$

$$y = m b^{q-s}$$

Par la division arithmétique $b^{p-q-1} l + m = b^{p-q} m' + r \wedge r < b^{p-q}$, on a la représentation de E sous la forme $E = r b^{q-s}$.

Ayant défini pour $q = p - 1, \dots, p - s$ les domaines :

$$\sum_{p-1,q} = \{ (x,y) \mid (x,y) \in T^2 \wedge e(x) = p - 1 \wedge e(y) = q \wedge x + y \geq b^{p-1} \},$$

parcourons ces domaines sur les droites $x = \text{constante}$.

Le parcours de $\sum_{p-1,q}$ à l fixé implique

* si $b^s(1 - b^{q-p}) < l < b^s$ l'équifréquence des restes

* si $b^s(1 - b^{q-p+1}) < l \leq b^s(1 - b^{q-p})$ l'obtention, dès que ρ défini par :

$$l = b l' + \rho \wedge p < b$$

est non nul, d'une majoration de 1 pour la fréquence des restes

$0, 1, 2, \dots, \rho b^{p-1-q} - 1$ sur la fréquence uniforme.

$$\text{L'image de } \sum_{p-1} = \bigcup_{q=p-1}^{p-s} \sum_{p-1,q} \text{ est } \mathcal{E}_{p-1} = \bigcup_{q=p-1}^{p-s} \mathcal{E}_{p-1,q}$$

VALEUR MOYENNE DE E

Les hypothèses : Considérer E comme variable aléatoire implique une hypothèse sur la probabilité d'obtention des différents points de T^2 .

Nous choisirons l'hypothèse d'équiprobabilité des points de T^2 .

Cette hypothèse de première approximation serait, pour une étude plus fine, à remplacer [12] par l'hypothèse, réalisant une meilleure approche du schéma réel, d'une distribution logarithmique des mantisses des points de T d'exposant fixé (Voir aussi [38]).

Soulignons toutefois que dans ce dernier cas, l'obtention des restes sur une droite $x = \text{constante}$ s'effectue avec des différences négligeables de probabilité de ces différents restes pour les ensembles $\sum_{p,q}, \sum_{p-1,q}$ lorsque q est proche de p ; par contre, pour l'ensemble $\{y \mid y \in T, y < b^{p-s}\}$ il existe une différence significative des valeurs moyennes calculées suivant les deux hypothèses.

Si le calcul qui suit - valeur moyenne de E quand (x,y) , de composantes positives, parcourt \sum - peut paraître d'un intérêt limité, dans le cas où les nombres sont représentés avec signe et valeur absolue, il n'en est pas de même dans le cas où les nombres négatifs sont représentés en notation de complément à 2:

Pour l'ensemble des couples (x,y) satisfaisant à $e(x \oplus y) = p$ avec des composantes de signes quelconques, il vient en effet, pour l'erreur entachant l'addition définie par une règle de troncature, une moyenne nulle dans le 1er cas, la symétrie étant un automorphisme de T.

Le calcul ci-dessous s'applique donc uniquement aux sommes arithmétiques dans la première hypothèse ; joint aux calculs du paragraphe suivant, il s'applique aux sommes algébriques dans la seconde hypothèse.

Le dénombrement des éléments des différents sous-ensembles entraîne :

$$\text{Card} \left(\sum_{p,p} \right) = \frac{1}{2}(b^s - 2 b^{s-1}) (b^s - 2 b^{s-1} + 1) \sim \frac{1}{2} b^{2s-2} (b - 2)^2$$

et pour $t = 1, \dots, s - 1$:

$$\begin{aligned} \text{Card} \left(\sum_{p,p-t} \right) &= (b^s - b^{s-1}) (b^s - b^{s-t} - b^{s-1}) + \frac{1}{2} (b^s - b^{s-1}) \left(\frac{b^s - b^{s-1}}{b^t} + 1 \right) \\ &\sim b^{2s-2} (b - 1) (b - 1 - b^{1-t}) + \frac{1}{2} b^{2s-2-t} (b - 1)^2 \end{aligned}$$

où

le 1er terme représente le nombre d'éléments de la partie rectangulaire de $\sum_{p,p-t}$ et le second terme le nombre d'éléments de la partie triangulaire de cet ensemble.

$$\begin{aligned} \text{Card} \left(\sum_{p-1,p-1} \right) &= (b^s - b^{s-1})^2 - \frac{1}{2} (b^s - 2 b^{s-1}) (b^s - 2 b^{s-1} + 1) \\ &\sim b^{2s-2} \left(\frac{b^2}{2} - 1 \right) \end{aligned}$$

et pour $t = 2, \dots, s$

$$\begin{aligned} \text{Card} \left(\sum_{p-1,p-t} \right) &= (b^s - b^{s-1}) (b^{s-t} - 1) + \frac{1}{2} (b^s - b^{s-1}) \left(\frac{b^s - b^{s-1}}{b^{t-1}} + 1 \right) \\ &\sim \frac{1}{2} b^{2s-1-t} (b^2 - 1) \end{aligned}$$

Nous pouvons remarquer que l'existence des restes en surnombre associés aux domaines $\sum_{p-1,q}$ est négligeable.

Ecrivons la valeur moyenne de E sous la forme :

$$\bar{E} = \frac{1}{2} b^{p-s} \left(1 - \sum_{t=0}^s \alpha_t \right)$$

avec :

$$\alpha_0 = \frac{\text{Card } \sum_{pp}}{\text{Card } \sum}, \quad \alpha_1 = b^{-1} \frac{\text{Card } \sum_{p-1,p-1} + 2 \text{Card } \sum_{p,p-1}}{\text{Card } \sum}$$

$$\alpha_t = 2 b^{-t} \frac{\text{Card } \sum_{p,p-t} + \text{Card } \sum_{p-1,p-t}}{\text{Card } \sum} \quad 1 < t < s$$

Le biais apporté par l'ensemble $\mathcal{E}_{p,p-s}$ est fonction de la distance de l'exposant $p - s$ à m_2 . Les calculs ci-dessous ne considèrent pas cet ensemble.

Il vient :

$$\text{Card} \left(\sum - \sum_{p,p-s} \right) \sim (2s - 1) (b - 1)^2 b^{2s-2}$$

et :

$$\bar{E}_{\mathcal{E}_{p,p-s}} = \frac{1}{2} b^{p-s} (1-\omega) \quad \text{ou} \quad \omega = \frac{b^2 + b - 2}{2(b-1)^2(2s-1)} \quad (\text{VII.1})$$

On constate que la différence, entre les moyennes théorique et pratique, est faible.

Ainsi :

pour $b = 2$, $s = 22$ $\omega = 0,046$

et

pour $b = 16$, $s = 6$ $\omega = 0,054$

Lemme VII.1.

Pour l'ensemble des couples (x,y) de T^2 satisfaisant aux conditions :

$$e(x \oplus y) = p, \quad x > 0, y > 0, \quad e(x) - e(y) < s$$

la moyenne de l'erreur entachant l'addition, définie par une règle de troncature, des deux composantes du couple, est égale à $\frac{1}{2} b^{p-s} (1-\omega)$ où ω a la valeur définie par (VII.1)

I . 3 . L'ADDITION DE DEUX NOMBRES DE SIGNES CONTRAIRES

Nous supposons $x > 0$, $y \leq 0$.

I . 3 . 1 . Le cas de la troncature sans chiffre de garde

Dans ce cas, l'opération \oplus est E-dirigée et non correcte. La valeur absolue de l'erreur d'arrondi élémentaire n'est évidemment plus contenue dans l'intervalle $[0, b^{p-s}]$. Il est cependant intéressant d'évaluer l'importance relative des cas de "cancellation".

CARACTERISATION DE Σ ET \mathcal{E}

Σ est l'ensemble des points $(x = 1 b^{p+t-s} \quad t \geq 0 \quad , \quad y = -m b^{q-s})$ de T^2 dont les coordonnées vérifient l'inégalité :

$$b^{s-1} \leq (1 - \lfloor m b^{q-p-t} \rfloor) b^t < b^s \quad (\text{VII . 2})$$

Définissant les domaines $\sum_{p+t,q}$ suivants :

$$\sum_{p+t,q} = \{(x,y) \mid x \in T^+_{\wedge} \quad y \in T^-_{\wedge} \quad e(x) = p+t \quad e(y) = q \quad b^{p-1} \leq x \oplus y < b^p\} \\ (q > p-s) ,$$

avec dans le cas $t = 0 \quad q = p-s$

$$\sum_{p,p-s} = \{(x,y) \mid x \in T^+_{\wedge} \quad y \in T^-_{\wedge} \quad e(x) = p \quad e(y) \leq p-s \}$$

on a :

$$\Sigma = \bigcup_{t=0}^{s-1} \left(\bigcup_{q \in I_t} \sum_{p+t,q} \right)$$

les ensembles I_t d'indices étant respectivement :

$$I_0 = \{p, p-1, \dots, p-s\}$$

$$I_1 = \{p+1, \dots, p+2-s\}$$

$$I_t = \{p+t, p+t-1\} \quad \text{pour } t = 2, \dots, s-1$$

L'ensemble $\mathcal{E}_{p+t,q}$ se déduit de $\sum_{p+t,q}$ par les égalités :

$$E = -r b^{q-s}$$

où r est défini par la division $m = b^{p+t-q} m' + r$ \wedge $r < b^{p+t-q}$

De (VII.2) on déduit immédiatement que $\sum_{p+t,q}$ est l'ensemble des points $(x = l b^{p+t-s}, y = -m b^{q-s})$ de T^2 dont les coordonnées vérifient :

$$\max (b^{s-1}, b^{p+t-q}(1 - b^{s-t} + 1)) \leq m \leq \min (b^{s-1}, b^{p+t-q}(1 - b^{s-1-t}) + b^{p+t-q-1})$$

Le parcours de $\sum_{p+t,q}$ sur les droites $x = \text{constante}$, implique alors l'équifréquence des éléments de $\mathcal{E}_{p+t,q}$.

VALEUR MOYENNE DE E .

Il vient :

$$\text{Card} \left(\sum_{p,p} \right) = \frac{1}{2} b^{2s-2} (b-2)^2$$

et pour $u = 1, \dots, s-1$:

$$\text{Card} \left(\sum_{p,p-u} \right) = b^{s-1} (b-1) \left(b^s - b^{s-1} - \frac{1}{2} b^{s-u} - \frac{1}{2} b^{s-1-u} + \frac{1}{2} \right)$$

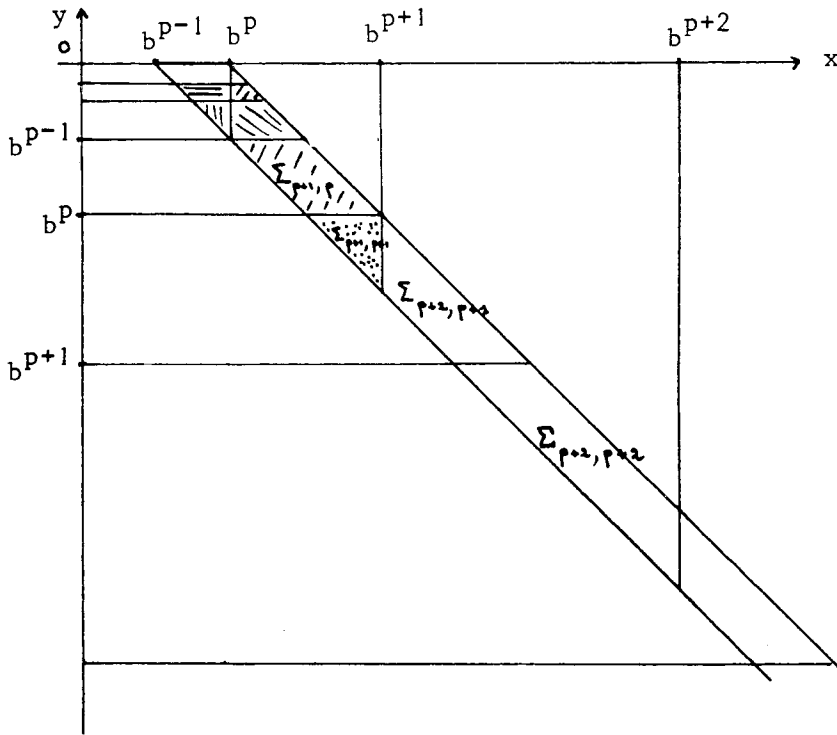
$$\text{Card} \left(\sum_{p+1,p} \right) = \frac{1}{2} b^{2s-3} (b^2 - 2) + \frac{1}{2} b^{s-1} (2 - b)$$

$$\text{Card} \left(\sum_{p+1,p+1-u} \right) = \frac{1}{2} (b^2 - 1) b^{2s-2-u} - \frac{1}{2} (b-1) b^{s-1} \quad 2 \leq u \leq s-1$$

$$\text{Card} \left(\sum_{p+t,p+t} \right) = b^{s-t-1} (b-1) \left(b^s - b^{s-1} - \frac{1}{2} b^{s-t} - \frac{1}{2} b^{s-t-1} + \frac{1}{2} \right) \quad t \geq 1$$

$$\text{Card} (\Sigma_{p+t, p+t-1}) = \frac{1}{2} (b-1) b^{s-t} (b^{s-t} + b^{s-t-1} - 1)$$

$t \geq 2$



Notons $\bar{E}_{p+t, q}$ la valeur moyenne de E sur le sous-ensemble $\mathcal{E}_{p+t, q}$.

La valeur moyenne \bar{E} de E sur l'ensemble \mathcal{E} s'exprime sous la forme :

$$\bar{E} = \left(\sum_{t=0}^{s-1} \sum_{q \in I_t} \bar{E}_{p+t, q} \times \text{Card} \Sigma_{p+t, q} \right) / \text{Card} \Sigma$$

$$\mathcal{E}_{p+t, q} = \{E = -r b^{q-s} \mid r = 0, \dots, b^{p+t-q} - 1 \text{ \u00e9quiprobables}\}$$

$$\bar{E}_{p+t, q} = -\frac{1}{2} b^{p+t-s} (1 - b^{q-p-t})$$

et :

$$\bar{E} = -\frac{1}{2} b^{p-s} \left[\left(\sum_{t=0}^{s-1} \sum_q (1 - b^{q-p-t}) b^t \text{Card} \Sigma_{p+t, q} \right) / \text{Card} \Sigma \right].$$

On obtient successivement :

$$\sum_t \sum_q (1 - b^{q-p-t}) b^t \text{Card} \sum_{p+t,q} \sim (s - \frac{1}{2}) (b - 1)^2 b^{2s-2}$$

et

$$\text{Card} (\sum - \sum_{p,p-s}) \sim (s - \frac{1}{2}) (b - 1)^2 b^{2s-2}$$

soit :

$\bar{E}_{\xi-\xi_{p,p-s}} = -\frac{1}{2} b^{p-s}$	(VII . 3)
--	-----------

Lemme VII . 2 .

Pour l'ensemble des couples (x,y) de T^2 satisfaisant aux conditions :

$$e(x \oplus y) = p \wedge x > 0 \wedge y \leq 0 \wedge |x| > |y| \wedge e(x) - e(y) < s$$

la moyenne de l'erreur entachant l'addition, définie par une règle de troncature sans chiffre de garde, des deux composantes du couple, est, en valeur absolue, exactement égale à $\frac{1}{2} b^{p-s}$.

I . 3 . 2 Le cas de la troncature avec un chiffre de garde

L'addition de deux nombres de signes contraires est correcte mais n'a pas de direction définie.

CARACTERISATION DE Σ ET \mathcal{E}

Avec des définitions analogues à celles du paragraphe précédent :

$$\Sigma = \bigcup_{t=0}^s (\bigcup_{q \in I_t} \sum_{p+t,q})$$

où :

$$I_0 = \{ p, p-1, \dots, p-s-1 \}$$

$$I_1 = \{ p+1, \dots, p+1-s \}$$

$$I_t = \{ p+t, p+t-1 \} \text{ pour } t = 2, \dots, s-1$$

$$I_s = \{ p+s-1 \}$$

L'ensemble \mathcal{E} :

* pour $2 \leq t \leq s$

$$\mathcal{E}_{p+t,q} = \{0\} \quad \forall q$$

* pour $t = 1$

$$\mathcal{E}_{p+1,p+1} = \mathcal{E}_{p+1,p} = \{0\}$$

L'addition pour $q < p$ de $x = 1 b^{p+1-s}$ et $y = -m b^{q-s}$ est

$$E \text{ - dirigée : } x \oplus y = (b \lfloor 1 - \lfloor m b^{q-p} \rfloor \rfloor) b^{p-s}$$

$$\sum_{p+1,q}^{(q < p)} \text{ est l'ensemble des points } (x = 1 b^{p+1-s}, y = -m b^{q-s})$$

de T^2 dont les coordonnées vérifient :

$$\max (b^{s-1}, b^{p-q} (b \lfloor 1 - b^s + 1 \rfloor)) \leq m \leq \min (b^{s-1}, b^{p-q} (b \lfloor 1 - b^{s-1} + 1 \rfloor - 1)) \\ = b^s - 1$$

et dans ce cas :

$$\mathcal{E}_{p+1,q} = \{ E = -r b^{q-s} \mid r = 0, 1, \dots, b^{p-q} - 1 \text{ équiprobables} \}$$

* pour $t = 0$

$$\mathcal{E}_{p,p} = \{0\}, \quad \mathcal{E}_{p,p-s-1} = \{ y \mid y \in T^- \wedge e(y) \leq p-s-1 \}$$

Pour $p-s \leq q \leq p-1$, $\sum_{p,q}$ est l'ensemble des points

($x = 1 b^{p-s}$, $y = -m b^{q-s}$) de T^2 dont les coordonnées vérifient :

$$\left\lfloor \frac{b l - \lfloor m b^{q-p+1} \rfloor}{b} \right\rfloor \geq b^{s-1}$$

ou

$$b^{s-1} \leq m \leq \min (b^s - 1, b^{p-q-1} (b l - b^s + 1) - 1)$$

En rappelant :

$$E = (r_2 b^{p-q-1} - r_1) b^{q-s} \quad \text{avec} \quad m = b^{p-q-1} m_1 + r_1 \wedge r_1 < b^{p-q-1}$$

$$b l - m_1 = b q_1 + r_2 \wedge r_2 < b,$$

le parcours de $\sum_{p,q}$ à l fixé, implique pour les éléments de \mathcal{E}_{pq} :

- si $b^s \leq b^{p-q-1}(bl - b^s + 1)$ une fréquence uniforme
- si $b^s > b^{p-q-1}(bl - b^s + 1)$ l'obtention pour l fixé, des erreurs $0, -b^{q-s}, \dots, -(b^{p-q-1} - 1)b^{q-s}$ en surnombre par rapport à la fréquence uniforme.

$$\mathcal{E}_{p,p-1} = \{ E = r b^{p-1-s} \mid r = 0, \dots, b-1 \}$$

$$\mathcal{E}_{p,q} = \{ E = r b^{q-s} \mid r = -(b^{p-q-1}-1), \dots, 0, 1, \dots, b^{p-q} - b^{p-q-1} \} \quad 1 < p-q < s$$

$$\mathcal{E}_{p,p-s} = \{ E = r b^{p-2s} \mid r = 1, \dots, b^s - b^{s-1} \}$$

VALEURS MOYENNES DE E

Le nombre d'éléments d'un sous-ensemble $\sum_{p+t,q}$ est égal,

au terme $O(b^s)$ près, au nombre d'éléments du sous-ensemble correspondant du § I.3.1 précédent.

L'unique formule de dénombrement à modifier est relative au nombre d'éléments de l'ensemble \sum , qui inclut $\text{Card} \sum_{p,p-s}$:

$$\text{Card} \left(\sum - \sum_{p,p-s-1} \right) \sim \left(s + \frac{1}{2} \right) (b-1)^2 b^{2s-2}$$

Dans les sous-ensembles $\mathcal{E}_{p,q}$, l'écart entre les fréquences des erreurs positives, et négatives ou nulles, est négligeable.

Nous avons étudié l'ensemble \mathcal{E} image de l'ensemble Σ :

$$\Sigma = \{ (x,y) \mid (x,y) \in T^2 \wedge e(x \oplus y) = p \wedge x > 0 \wedge y \leq 0 \wedge |x| > |y| \}$$

avec l'hypothèse de la représentation des nombres négatifs par signe et valeur absolue. Le calcul de la valeur moyenne de E quand (x,y) parcourt Σ sera effectué dans les deux hypothèses de représentation des nombres négatifs.

1) Les nombres négatifs sont représentés par leur signe et leur valeur absolue.

La valeur moyenne de E quand (x,y) parcourt $\Sigma - \sum_{p,p-s-1}$ s'écrit :

$$\bar{E} - \mathcal{E}_{p,p-s-1} = \frac{1}{2} b^{p-s} (1 - \omega) \quad \text{où} \quad \omega = \frac{4s(b-1) + (b-2)^2}{(2s+1)(b-1)b} \quad (\text{VII.4})$$

soit, pour :

b	2	4	8	16
ω	$\frac{2s}{2s+1}$	$\frac{3s+1}{6s+3}$	$\frac{7s+9}{28s+14}$	$\frac{15s+49}{120s+60}$

2) Les nombres négatifs sont représentés en notation de complément à 2

L'opération \oplus est G - dirigée.

Il vient pour les sous-ensembles de \mathcal{E} à éléments non tous nuls :

$$\mathcal{E}_{p+1,q} = \mathcal{E}_{p,q} = \{ E = r b^{q-s} \mid r = 0, 1, \dots, b^{p-q} - 1 \} \quad p-1 \leq q \leq p+1-s$$

$$\mathcal{E}_{p,p-s} = \{ E = r b^{p-2s} \mid r = 1, \dots, b^s - b^{s-1} \}$$

$$\bar{E}_{\mathcal{E}-\mathcal{E}_{p,p-s-1}} = \frac{1}{2} b^{p-s} (1 - \omega) \quad \text{où} \quad \omega = \frac{b^2 + 4b - 2}{(2s+1)(b-1)b} \quad (\text{VII.5})$$

pour $b = 2$, $s = 22$ $\omega = 0,11$

pour $b = 16$, $s = 6$ $\omega = 0,10$

Lemme VII . 3 .

Pour l'ensemble des couples (x,y) de T^2 satisfaisant aux conditions

$$e(x \oplus y) = p \wedge x > 0 \wedge y < 0 \wedge |x| > |y| \wedge e(x) - e(y) \leq s$$

la moyenne de l'erreur entachant l'addition, définie par une règle de troncature avec chiffre de garde, des deux composantes du couple, est égale à $\frac{1}{2} b^{p-s} (1 - \omega)$ où , suivant la représentation des nombres négatifs, par signe et valeur absolue ou en notation de complément à 2 , ω prend les valeurs définies respectivement par (VII.4) et (VII.5).

I . 4 CONCLUSION RELATIVE A L'ADDITION

Les calculs de moyennes effectués aux paragraphes précédents pour respectivement, l'addition et la soustraction arithmétiques, permettent d'évaluer la valeur moyenne de l'erreur entachant l'addition de deux nombres de signes quelconques.

(Nous rappellerons que les éléments des sous-ensembles $\mathcal{E}_{p,p-s}$ - resp. $\mathcal{E}_{p,p-s-1}$ - sont égaux en valeur absolue et de signes opposés pour

l'addition et la soustraction arithmétiques dans le cas d'une troncature sans chiffre de garde - resp. avec chiffre de garde, avec de plus dans ce cas une moyenne de $\frac{1}{2} b^{p-s}$ pour les couples (x,y) tels que $e(x) = p \wedge e(y) = p - s$.

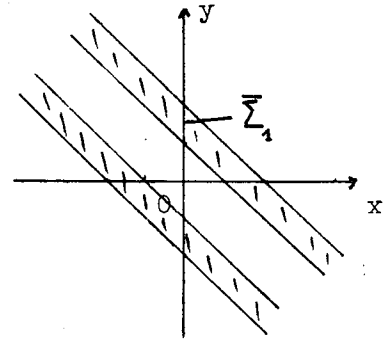
Etant donné le domaine $\bar{\Sigma}$:

$$\bar{\Sigma} = \{ (x,y) \mid (x,y) \in T^2 \wedge e(x \oplus y) = p \}$$

et l'image $\bar{\mathcal{E}}$ de $\bar{\Sigma}$ par l'application :

$$(x,y) \rightarrow E = x + y - (x \oplus y)$$

nous ferons les remarques suivantes :



1) L'écart entre les valeurs moyennes, théorique et calculée exactement, de E sur $\bar{\mathcal{E}}$, est assez faible.

Toutefois, l'existence de ce biais dans le cas d'une représentation des nombres négatifs en complément à 2, se révélera lorsque les sommations s'effectueront sur un grand nombre de termes.

2) Pour des calculs effectués sur des composantes de couples (x,y) appartenant à des sous-ensembles particuliers de $\bar{\Sigma}$, par exemple

$$\bar{\Sigma}_1 = \{ (x,y) \mid (x,y) \in \bar{\Sigma} \wedge x \oplus y > 0 \},$$

l'hypothèse théorique de répartition de l'erreur d'arrondi ne donnera pas des résultats satisfaisants.

Il est à noter que, dans certains cas, le raffinement apporté par le chiffre de garde, s'il diminue la moyenne de l'erreur élémentaire, peut paradoxalement nuire à la précision d'une suite de calculs : Ainsi, les études statistiques de H. KUKI et W.J. CODY [17] établissent, relativement à la somme algébrique, une précision meilleure pour une troncature sans chiffre de garde que pour une troncature avec un chiffre de garde, dans l'hypothèse de la représentation des nombres négatifs avec signe et valeur absolue ; le phénomène de compensation entre les différentes erreurs élémentaires, a diminué plus fortement l'erreur totale dans le 1er cas.

Un cas similaire de "rupture de symétrie" est présenté par la règle d'arrondi - définissant $x \star y$ par l'égalité :

$$|x * y - x \otimes y| = \min_{z \in T} |x * y - z| \quad - \text{ quand l'arrondi se fait}$$

systématiquement loin du zéro (resp. vers le zéro) lorsque l'égalité

$x * y = \frac{\Delta x + \nabla x}{2}$ a lieu. La perte de précision entraînée par la dissymétrie de la définition, a été soulignée par M. URABE [39] et H. KUKI et W.J. CODY [17], qui ont proposé des modifications de la règle simple précédente.

Enfin, dans le cadre de cette étude, nous mentionnerons l'obtention d'une valeur moyenne strictement nulle de l'erreur d'entrée, dans l'arithmétique paritaire définie par S. LINNAINMAA [21].

I. 5. LE CAS DE LA MULTIPLICATION

L'étude de la répartition de l'erreur d'arrondi élémentaire est, dans le cas de la multiplication en arithmétique à virgule flottante, plus complexe que dans le cas de l'addition, et nous poserons ici seulement le problème.

Considérons par exemple, le cas où l'opération \otimes de multiplication définie sur T , est correcte et I - dirigée.

Au couple (x,y) de $(T^+)^2$, de représentation :

$$x = lb^{e(x)-s}$$

$$b^{s-1} \leq 1, m < b^s$$

$$y = mb^{e(y)-s}$$

est associée l'erreur élémentaire E de multiplication :

$$E = rb^{e(x) + e(y) - 2s}$$

où r est défini par la division arithmétique suivante :

$$\left| \begin{array}{l} lm = b^{s-\theta} q + r \wedge r < b^{s-\theta} \\ \text{avec } \theta = 1 \text{ si } lm < b^{2s-1} \\ \theta = 0 \text{ si } lm \geq b^{2s-1} \end{array} \right. \quad (\text{VII. 6})$$

On est ainsi amené à considérer les 2 ensembles Σ^1 et Σ^2 de couples (l,m) d'entiers, d'images intérieures au carré $\left| \begin{array}{l} b^{s-1} \leq l < b^s, \\ b^{s-1} \leq m < b^s \end{array} \right.$,

et vérifiant respectivement :

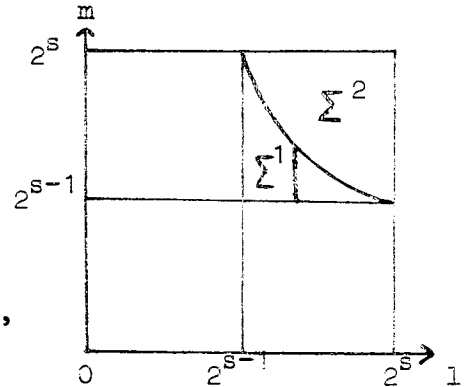
$$lm < b^{2s-1} \text{ pour } \Sigma^1, \quad lm \geq b^{2s-1} \text{ pour } \Sigma^2.$$

Pour chaque ensemble Σ^1 et Σ^2 , il faut étudier les valeurs prises par le reste r de (VII. 6) lorsque le couple (l,m) parcourt l'ensemble.

Dans l'hypothèse par exemple d'une base égale à 2, le parcours de Σ^1 (de manière analogue Σ^2) sur les droites $l = \text{constante}$, entraîne les résultats suivants [13] :

- Pour l fixé impair, m variant de 2^{s-1} à $\lfloor 2^{2s-1}/l \rfloor$, les restes des divisions de lm par 2^{s-1} , prennent leurs valeurs parmi $0, 1, 2, \dots, 2^{s-1}-1$ et forment un système de restes non congrus (mais incomplet).

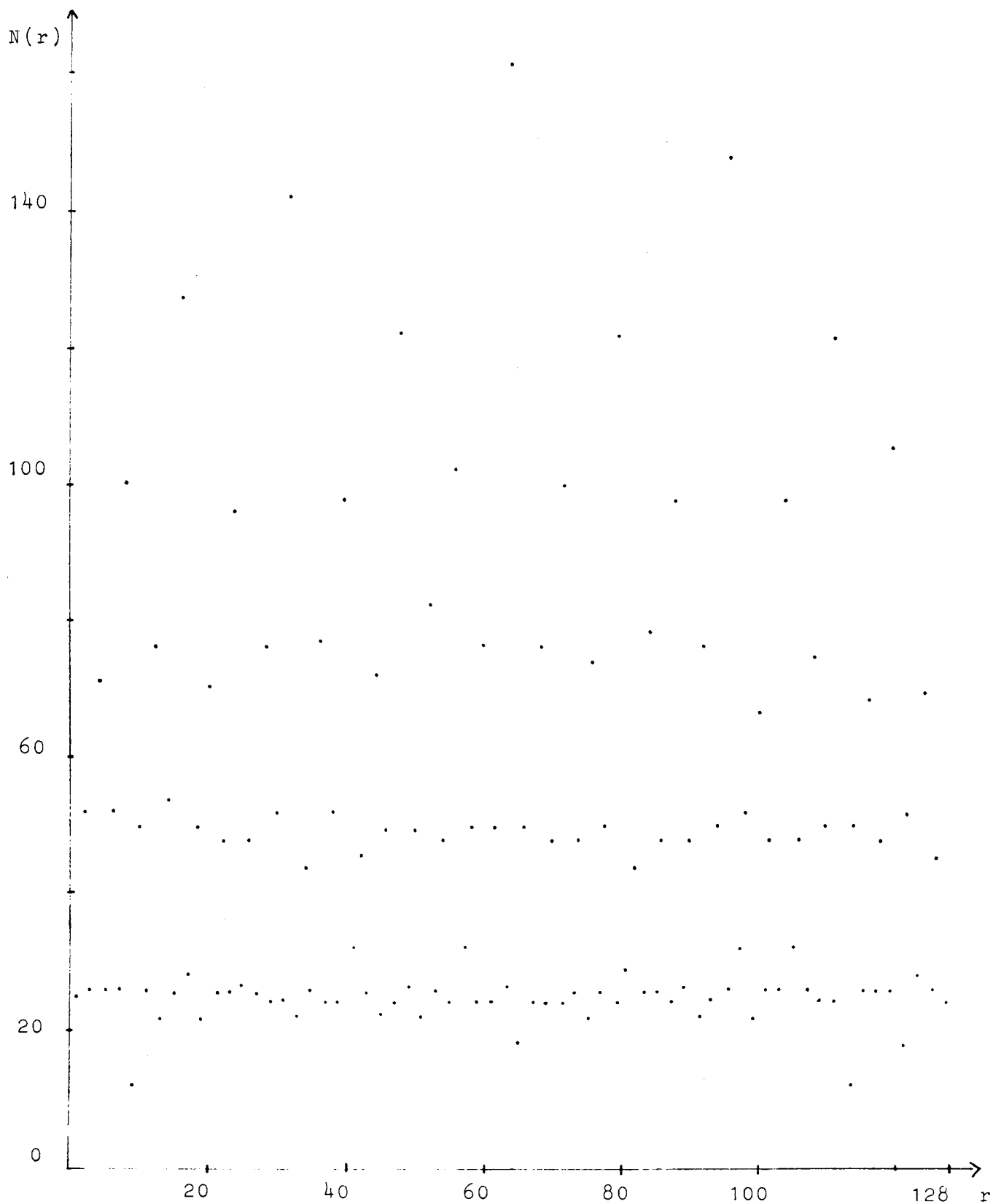
- Pour l fixé pair, $l = 2^k l'$, l' impair, m variant de 2^{s-1} à $\lfloor 2^{2s-1}/l \rfloor$, les restes des divisions prennent leurs valeurs parmi $0, 2^k, 2 \cdot 2^k, \dots, (2^{s-k-1}-1) 2^k$.



L'image \mathcal{E}_k^1 de Σ^1 par l'application : $(l,m) \rightarrow r$ est l'union $\bigcup_{k=0}^{s-1} \mathcal{E}_k^1$, où \mathcal{E}_k^1 désigne l'ensemble image de la partie de Σ^1 obtenue pour les valeurs de l multiples de 2^k .

Malheureusement, montrer que la répartition des éléments de \mathcal{E}_k^1 , $\forall k$, est analogue à une répartition uniforme, se révèle difficile en raison de la limitation $lm < 2^{2s} - 1$ (En arithmétique fixe par contre, on a exactement l'équirépartition des éléments des sous-ensembles de type \mathcal{E}_k^2).

Nous donnons ci-dessous à titre d'exemple , pour $b=2$ et $s=8$,
les fréquences des différents restes obtenus lorsque $(1,m)$ parcourt \sum^4 :



II - LE PROBLEME D'UN MEILLEUR CHOIX DE LA BASE

Il est nécessaire de mentionner dans ce travail, le problème du choix optimal de la base de la représentation des nombres à virgule flottante.

Nous en rappellerons ci-dessous les points essentiels.

Ce problème peut être posé sous plusieurs formes.

Si le nombre de bits affecté à la représentation des nombres à virgule flottante est fixé, on comparera les différentes représentations obtenues lorsque varient, d'une part le choix de la base, d'autre part le nombre de bits affectés respectivement à la mantisse et à l'exposant.

La comparaison entre les représentations doit s'effectuer en tenant compte de :

- la largeur du domaine de R représenté
- la précision de la représentation.

Ce dernier point de vue est assez délicat et différentes mesures de la précision ont été proposées par plusieurs auteurs :

Une mesure de la précision est donnée par l'erreur relative de représentation d'un réel par un nombre à virgule flottante.

Cette erreur est, suivant la terminologie de W.J. CODY, une caractéristique statique de la représentation.

L'étude de W.M. Mc KEEMAN ([23] , 1967) sur la valeur moyenne de cette erreur relative, pour une distribution logarithmique des mantisses, établit qu'il n'existe pas, relativement à ce critère, de différence significative entre les performances des bases puissances de 2, dans un calculateur binaire.

L'étude de W.S. BROWN et P.L. RICHMAN ([3], 1969) sur la borne supérieure de l'erreur relative définie ci-dessus, établit qu'un meilleur choix de la base est donné par $b = r$, r désignant le nombre d'états possibles du dispositif élémentaire composant le mot-mémoire.

[23] et [3] établissent la supériorité de $b = 2$ pour un calculateur binaire, lorsqu'est utilisée, pour une représentation normalisée des nombres, l'écriture implicite du bit le plus significatif de la mantisse -bit toujours égal à 1 sauf pour le nombre nul-.

D.W. MATULA ([27], 1972) a mis en évidence une perte de précision liée au choix d'une large base et entraînée par les 2 points suivants :

- la distribution du chiffre de tête de mantisse étant non uniforme mais logarithmique, accentue le défaut de précision provoqué par les mantisses à bits de tête nuls.
- il existe des intervalles réels, non petits, sur lesquels pour $b > d$, l'espacement entre des éléments adjacents de T_b^s est plus grand que l'espacement entre des éléments adjacents de T_d^s .

Cependant, comme le souligne W.J. CODY, il est utile d'élargir la discussion et de considérer les caractéristiques dynamiques de tout le système arithmétique à virgule flottante, défini par la représentation des nombres et par l'arithmétique utilisée.

R.P. BRENT ([2], 1973) étudie :

- l'écart-type de l'erreur relative d'arrondi d'une part
- les résultats, pour différentes bases et différentes arithmétiques, de la simulation de sommes, résolutions de systèmes linéaires, calculs de valeurs propres, d'autre part,

établit, pour ces critères, que le choix $b = 2$ avec le bit le plus significatif implicite, est le meilleur, suivi par le choix $b = 4$, les opérations flottantes étant optimales.

Une conclusion semblable est donnée par W.J. CODY ([4], 1973), à savoir que compte-tenu des contraintes de l'écriture implicite d'un bit en base 2, une représentation des nombres en base 4 et une arithmétique d'arrondi sur deux chiffres de garde, assurera un meilleur système arithmétique.

Mentionnons aussi, en faveur de la base 4, le travail de T. KREIFELTS ([16], 1973).

Enfin, à coté de l'étude statistique, relative à la précision des systèmes de nombres à virgule flottante, de H. KUKI et W.J. CODY ([17], 1973), nous mentionnerons le travail de J.D. MARASA et D.W. MATULA ([26], 1973) : Ces auteurs, dans une étude simulative de la propagation de l'erreur d'arrondi au cours de longs calculs entre variables corrélées, ont établi que la croissance de l'erreur dépend essentiellement du type des opérations impliquées dans le calcul : pourcentages respectifs d'additions, soustractions, multiplications et divisions ; la soustraction en particulier, conduit en raison des cancellations, aux plus fortes croissances de l'erreur ; l'arrondi, préférable à la troncature en général, n'introduit qu'un facteur multiplicatif, variant peu avec le nombre d'opérations réalisées.

Il apparait enfin que le choix de la base est moins décisif que la possibilité de disposer sur un calculateur de 2 arithmétiques flottantes, l'utilisation parallèle de ces arithmétiques étant particulièrement significative lorsque les opérateurs sont monotones.

(Voir [46] pour la réalisation d'un tel calculateur).

L'importance de la considération attentive de la "stratégie" des calculs sera à nouveau soulignée, en rappelant, par exemple les résultats relatifs à la résolution itérative des systèmes linéaires ou le fait qu'une sommation particulière de termes de signes opposés au chapitre V, a permis d'éviter l'effet traditionnellement néfaste de ces opérations.

ANNEXE

CARACTERISTIQUES DE QUELQUES ARITHMETIQUES
A VIRGULE FLOTTANTE

A titre d'exemple, nous avons noté ci-dessous les caractéristiques de quelques arithmétiques à virgule flottante fréquemment utilisées, et considérées pour des représentations normalisées.

En relation avec le chapitre II, nous avons rappelé, pour chaque mode d'opération, les théorèmes permettant le calcul des erreurs d'arrondi élémentaires.

Nous soulignerons que sur certains calculateurs (CDC 6000-7000, UNIVAC 1108), l'accès direct à l'erreur d'arrondi entachant une opération élémentaire, est prévue par instruction de l'arithmétique flottante : le résultat d'une opération élémentaire est dans ce cas, écrit sur deux nombres à virgule flottante.

Notation : Nous désignerons par \mathcal{C} , le mode de troncature réalisant pour $X \in T_b^s$, $Y \in T_b^s$ l'opération d'addition-soustraction ou de multiplication selon la succession d'étapes suivantes :

1) Pour l'addition :

On décale vers la droite la mantisse du nombre de plus petit exposant, sans perte de précision

On ajoute les 2 mantisses alors cadrées.

Pour la multiplication :

On multiplie les mantisses de X et Y .

Soit $m(R)$ le produit de composition de mantisses obtenu.

2) On normalise $m(R)$ uniquement dans le cas où s'est produit un dépassement de capacité de mantisse, ce qui ne peut avoir lieu qu'en cas d'addition de 2 nombres de même signe.

3) $m(R)$ est tronquée à s chiffres.

le calcul est achevé de manière habituelle.

type de calculateur	base	nombre de chiffres de mantisse	représentation de la mantisse	Mode Addition - Soustraction
PDP 11/45	2	23 bits + 1 bit muet 55 bits + 1 bit muet	Signe et valeur absolue	Deux modes sont possibles: Opération correcte Opération correcte
IBM 360 370	16	6 14 28 poss.	Signe et valeur absolue	Troncature avec un chiffre de garde
cii 10 070	16	6 14	Nombres négatifs en notation de complément à 2	Troncature avec un chiffre de garde Troncature sans chiffre de garde
IRIS 80	16	6 14	Nombres négatifs en notation de complément à 2	Deux modes 1) Troncature avec un chiffre de garde 2) Arrondi sur chiffre de garde Deux modes 1) Troncature sans chiffre de garde 2) Arrondi sur chiffre de garde

<p>opérations</p> <p>multiplication - Division</p>	<p>Calcul des erreurs élémentaires suivant :</p>
<p>ans chacun des formats</p> <p>t I - dirigée</p> <p>t E - dirigée</p>	<p>Propositions II.3 II.9, II.9', II.10 II.13</p>
<p>Opération correcte et I - dirigée</p>	<p>Propositions II.2 II.9, II.11 II. 13</p>
<p>Op. correcte et I - dirigée</p> <p>Op. correcte et I - dirigée</p>	<p>Propositions II.6 II.9, II.9', II.13</p>
<p>sont possibles :</p> <p>Op. correcte et I - dirigée</p> <p>Arrondi</p> <p>sont possibles :</p> <p>Op. correcte et I - dirigée</p> <p>Arrondi</p>	<p>Propositions II.6 II.9, II.9', II.13</p> <p>Voir [5]</p> <p>Voir [5]</p>

type de calculateur	base	nombre de chiffres de mantisse	représentation de la mantisse	Mode Addition - Soustractif
CDC 6600 7600	2	48	Nombres négatifs en notation de complément à 1	Deux modes sont possible: 1) Ecriture en DOUBLE sur 2 nombres flottants Détermination du poids Troncature $\bar{6}$ 2) Ecriture en SIMPLE selon des processus
UNIVAC 1108	2	27 60	Nombres négatifs en notation de complément à 1	Ecriture en DOUBLE sur 2 nombres flottants Détermination du poids Troncature sans chiffre de garde Troncature sans chiffre de garde

Note : Lorsque l'erreur d'arrondi élémentaire est déterminée à l'aide d'une instruction

- Pour une opération d'addition - soustraction, cette erreur est mise à 0 en cas
- Pour une opération de division, la quantité déterminée est non pas l'erreur E

<p>opérations</p> <p>Multiplication - Division</p>	<p>Calcul des erreurs élémentaires suivant :</p>
<p>suivant la précision désirée :</p> <p>Précision du résultat</p> <p>exposants appropriés (sauf division)</p> <p>Part du résultat :</p> <p>correcte et I - dirigée</p> <p>Précision du résultat</p> <p>particuliers d'arrondi</p>	<p>Directement par instructions de l'arithmétique flottante pour l'addition, la soustraction, la multiplication (mode 1 d'opérations)</p>
<p>Précision du résultat</p> <p>exposants appropriés.</p> <p>Part du résultat :</p> <p>opération correcte</p> <p>et I - dirigée</p> <p>) : Troncature ζ</p> <p>) : correcte</p> <p>et I - dirigée</p>	<p>Directement par instructions de l'arithmétique flottante</p> <p>Voir paragraphe I. 3. 1. b</p>

de l'arithmétique flottante :

l'écart trop grand entre les exposants des 2 opérandes.

elle-même, mais la quantité E Y .

B I B L I O G R A P H I E

- [1] BABUSKA, I. : Numerical Stability in Mathematical Analysis.
IFIP Congr. 68, Invited papers, 11-23 (1968).

- [2] BRENT, R.P : On the Precision Attainable with Various Floating-Point
Number Systems.
IEEE Trans. Computers, C-22, 6 (June 1973), 601-607.

- [3] BROWN, W.S and RICHMAN, P.L : The Choice of Base.
Comm. ACM 12,10 (Oct.1969), 560-561.

- [4] CODY, W.J, Jr : Static and Dynamic Numerical Characteristics of
Floating-Point Arithmetic.
IEEE Trans. Computers, C-22, 6 (June 1973), 598-601.

- [5] DEKKER , T.J : A Floating-Point Technique for Extending the
Available Precision.
Numer. Math. 18, 224-242 (1971).

- [6] DURAND, E : Solutions numériques des équations algébriques.
Tome II. Masson (1972).

- [7] FIEDLER, M and PTÁK, V : On matrices with non-positive off-diagonal
elements and positive principal minors.
Чехословацкий математический журнал T.12 (87)
1962, Prague pp. 382-400.

- [8] GASTINEL, N : Matrices du second degré et normes générales en
analyse numérique linéaire.
Publications scientifiques et techniques du
Ministère de l'Air, Paris, N.T.110-S.D.I.T (1960).

- [9] GASTINEL, N : Analyse Numérique linéaire.
Hermann, Paris 1966.

- [10] GENTLEMAN, W. MORVEN and MAROVICH, S.B : More on Algorithms that
Reveal Properties of Floating Point Arithmetic Units.
Short Comm. ACM 17, 5 (May 1974), 276-277.

- [11] GREGORY, J. : A Comparison of Floating-Point Summation Methods.
Short Comm. ACM 15, 9 (Sept. 1972), 838.

- [12] HAMMING, R.W : On the Distribution of Numbers.
Bell Syst. Tech. Journal, Vol 49 (October 1970)
pp. 1609-1625.

- [13] HARDY, G.H and WRIGHT, E.M : An Introduction to the theory of numbers.
4th edition. OXFORD University Press 1962, theorem 56 p 51.
- [14] KAHAN, W : Further remarks on reducing truncation errors.
Comm. ACM 8,1 (January 1965), 40.
- [15] KNUTH, D.E : The Art of Computer Programming. Vol 2 .
Addison Wesley (1969).
- [16] KREIFELTS, T : Optimal choice of basis for a floating-point arithmetic
Computing 11,4 (1973), 353-363 (German) .
- [17] KUKI, H and CODY, W.J : A Statistical Study of the Accuracy of
Floating Point Number Systems.
Comm. ACM 16, 4 (April 1973), 223-230.
- [18] KULISCH, U : An Axiomatic Approach to Rounded Computations.
Numer. Math. 18, 1-17 (1971).
- [19] LA PORTE, M et VIGNES, J : Etude statistique des erreurs dans
l'arithmétique des ordinateurs; application au contrôle
des résultats d'algorithmes numériques.
Numer. Math. 23, 63-72 (1974).
- [20] LA PORTE, M et VIGNES, J : Algorithmes numériques, analyse et mise
en oeuvre.
Technip Paris, 1974.
- [21] LINNAINMAA, S : Analysis of some known methods of improving the
accuracy of floating-point sums.
BIT 14 (1974), 167-202.
- [22] LINZ, P : Accurate Floating-Point Summation.
Comm. ACM 13, 6 (June 1970), 361-362.
- [23] Mc KEEMAN, W.M. : Representation Error for Real Numbers in Binary
Computer Arithmetic.
IEEE Trans. on Electronic Computers (October 1967), 682-683.
- [24] MALCOLM, M.A : On Accurate Floating-Point Summation.
Comm. ACM 14, 11 (November 1971), 731-736.

- [25] MALCOLM, M.A : Algorithms To Reveal Properties of Floating-Point Arithmetic.
Comm. ACM 15, 11 (november 1972), 949-951.
- [26] MARASA, J.D and MATULA, D.W : A Simulative Study of Correlated Error Propagation in Various Finite-Precision Arithmetics.
IEEE Trans. Computers, C-22, 6 (june 1973), 587-597.
- [27] MATULA, D.W : Significant digits : Numerical analysis or numerology.
IFIP Congr. 71, vol 2, pp. 1278-1283 (1972).
- [28] MOLLER, O : Quasi double-precision in floating point addition.
BIT 5 (1965), 37-50.
- [29] MOORE, R.E : Interval analysis.
Prentice Hall 1966.
- [30] NICKEL, K and RITTER, K : Termination Criterion and Numerical Convergence.
Siam J. Numer. Anal. Vol. 9 , No 2, June 1972, pp 277-283.
- [31] PICHAT, M : Correction d'une somme en arithmétique à virgule flottante.
Numer. Math. 19, 400-406 (1972).
- [32] PICHAT, M : Un algorithme, en arithmétique à virgule flottante, pour augmenter la précision d'une somme.
C.R.A.S. t. 274 (1972) série A, pp. 1639-1642.
- [33] PICHAT, M : Correction d'une somme calculée en arithmétique à virgule flottante, à l'aide de l'arithmétique utilisée.
Colloque d'Analyse Numérique Epinal Juin 1972.
- [34] PICHAT, M : Stabilité numérique de différents schémas de sommation dans un sous-ensemble de \mathbb{R} de type flottant.
C.R.A.S t. 277 (1973) série A, pp. 315-317.
- [35] PICHAT, M : Utilisation conjointe de deux arithmétiques flottantes
Séminaire A.N. Grenoble I (Janv. 1975) n° 216.
- [36] ROSS, D.R : Reducing Truncation Errors Using Cascading Accumulators.
Comm. ACM 8 , 1 (January 1965), 32-33.

- [37] TIENARI, M : On the control of floating-point mantissa length in iterative computations.
International Computing Symposium 1973. A. Günther et al.(eds.) 1974, 315 - 322.
- [38] TSAO, N.K : On the Distributions of Significant Digits and Roundoff Errors.
Comm. ACM 17, 5 (May 1974), 269-271.
- [39] URABE, M : Roundoff Error Distribution in Fixed-Point Multiplication and a Remark about the Rounding Rule.
SIAM J. Numer. Anal., Vol. 5, n° 2 (june 1968) pp 202-210.
- [40] VELTKAMP, G.W : Private communications.
Technological University, Eindhoven (1968).
- [41] WALKER, R.J. : Binary Summation.
Short Comm. ACM 14, 6 (june 1971), 417.
- [42] WILKINSON, J.H : Rounding errors in algebraic processes.
Her Majesty's Stationary Office (1963).
- [43] WILKINSON, J.H : The algebraic eigenvalue problem.
Clarendon Press, Oxford, 1965.
- [44] WOLFE, J.M : Reducing truncation errors by programming.
Comm. ACM 7, 6 (June 1964), 355-356.
- [45] YOHE, J.M : Interval Bounds for Square Roots and Cube roots.
Computing, Vol. 11, Fasc. 1, 1973.
- [46] YOHE, J.M : Roundings in Floating-Point Arithmetic.
IEEE Trans. Computers, C-22, 6 (June 1973), 577-586.

INDEX DES TERMES ET NOTATIONS

<u>Termes</u>	<u>Notations</u>	<u>Pages</u>
Algorithme de correction d'une somme	α	52
Arrondi de \mathbb{R} dans un sous-ensemble fini de \mathbb{R}	\square	1
Arrondi G - dirigé	∇	1
Arrondi D - dirigé	Δ	1
Cancellation		9
Convergence numérique faible		87
Convergence numérique forte		86
D - suite d'opérations approchées correspondant à une suite d'opérations sur \mathbb{R}		144
D - suite maximisante d'opérations " " "		144
Exposant d'un nombre flottant	$e(x)$	10
G - suite d'opérations approchées correspondant à une suite d'opérations sur \mathbb{R}		144
G - suite maximisante d'opérations " " "		144
Image de $x \in S^N$ par une suite O d'opérations	$i_0(x)$	144
Intervalle de précision d'un élément x de S	$\epsilon(x)$	5
Nombre flottant de plus petit module	ξ	10
Opération correcte		3
Opération D - dirigée		3
Opération E - dirigée		3
Opération G - dirigée		3
Opération I - dirigée		3
Opération induite de \star sur S par \square	$\star \square$	3
Opération optimale		3
Opération proprement tronquante		19
Prédécesseur de x dans S	$(x)'$	4
Règle de troncature sans chiffre de garde		23
Règle de troncature avec chiffre de garde		28
Règle d'arrondi sur chiffre de garde		31
Schéma de sommation dichotomique		115

Sous-ensemble de $[-1, 1]$ à intervalles constants		5
Sous-ensemble de \mathbb{R} à intervalles non décroissants		7
Sous-ensemble de \mathbb{R} de type flottant		7
Successeur de x dans S	$(x)''$	4
Suite conjuguée d'une suite O d'opérations correctes	\bar{O}	145
Systeme de nombres à virgule flottante	T	9
(de base b , de nombre de chiffres de mantisse s)	T_b^s	

T A B L E D E S M A T I E R E S

	<u>Page</u>
INTRODUCTION	
<u>CHAPITRE I</u> ARRONDIS SUR \mathbb{R}	
I - Arrondis définis de \mathbb{R} dans un sous-ensemble fini de \mathbb{R}	1
II - L'arithmétique "approchée"	2
III - Quelques types de sous-ensembles de \mathbb{R}	4
III.1. Sous-ensembles de $[-1,1]$ à intervalles constants	
III.2. Sous-ensembles de \mathbb{R} à intervalles non décroissants	
III.3. Sous-ensembles de \mathbb{R} de type flottant	
III.4. Systèmes de nombres à virgule flottante	
<u>CHAPITRE II</u> CALCUL AUTOMATIQUE DE L'ERREUR D'ARRONDI ELEMENTAIRE EN ARITHMETIQUE A VIRGULE FLOTTANTE	
I - L'erreur d'arrondi dans l'addition	14
I.1. Théorème préliminaire	
I.2. Expression, par un produit de composition sur T , de l'erreur d'arrondi.	
I.2.1. Première formule	
I.2.2. Généralisation d'un théorème de D.E.KNUTH	
I.3. Quelques exemples de règles d'addition en arithmétique à virgule flottante	
I.3.1. La troncature sans chiffre de garde	
I.3.2. La troncature avec chiffre de garde	
I.3.3. L'arrondi sur chiffre de garde	
I.3.4. Opérations \oplus optimales	
II - L'erreur d'arrondi dans la multiplication	35
II.1. Expression, par un produit de composition sur T , de l'erreur d'arrondi	
II.2. Formules pratiques	
II.3. Etude expérimentale	
III - L'erreur d'arrondi dans la division	43
III.1. Théorème préliminaire	
III.2. Calcul de l'erreur	
IV - CONCLUSION	47
V - APPENDICE . SOUS - PROGRAMMES FORTRAN	48

CHAPITRE III APPLICATIONS

I - Un algorithme de correction d'une somme algébrique	52
I.1. Le problème étudié	
I.2. Convergence de l'algorithme	
I.2.1. Deux théorèmes relatifs à une opération d'addition correcte	
I.2.2. Un théorème de convergence pour une addition définie par une troncature sans chiffre de garde.	
I.3. Remarques	
I.4. Test d'arrêt de l'algorithme	
I.5. Etude numérique	
I.6. Conclusion	
II - Obtention du problème perturbé dans l'analyse à postériori des erreurs d'arrondi	64
II.1. Introduction	
II.2. Un exemple : l'algorithme de Gauss	
II.2.1. Notations	
II.2.2. Rappel des résultats de l'analyse à postériori des erreurs d'arrondi	
II.2.3. Etude expérimentale	
III - Utilisation de méthodes de raffinement itératif	70
III.1. Introduction	
III.2. Un exemple : correction de l'inverse d'une matrice	
III.2.1. Formule théorique	
III.2.2. Calculs en arithmétique à virgule flottante	
III.2.3. Exemple numérique	
IV - Résolution numérique de problèmes mal conditionnés	77
IV.1. Introduction	
IV.2. Evaluation numérique d'un polynôme	
IV.3. L'élimination de Gauss.	
V - Critère d'application pour les méthodes des § III et IV	84

CHAPITRE IV L'ERREUR D'ARRONDI DANS LES METHODES ITERATIVES LINEAIRES

I - Raffinement de la solution numérique d'un système linéaire, calculée par une méthode itérative	86
I.1. Notations	
I.2. Méthode	
I.3. Bornes de l'erreur	
I.4. Evaluation, avant correction, de l'incertitude sur la solution numérique à raffiner.	
I.5. Exemple numérique	

II - Itérations linéaires définies par une matrice positive (resp. négative)	96
III - Résolution itérative en arithmétique à virgule flottante de l'équation linéaire $x=ax+c$	104
IV- Conclusion	110

CHAPITRE V STABILITE NUMERIQUE DE DIFFERENTS SCHEMAS DE SOMMATION.
CALCUL DE SOMMES DE SERIES

I - Etude comparative de schémas de sommation	112
I.1. Premier théorème	
I.2. Théorèmes généraux de comparaison des sommes séquentielle et dichotomique	
I.2.1. Sommation d'éléments à valeurs positives et décroissantes	
I.2.2. Sommation d'éléments à valeurs positives quelconques	
I.3. Théorèmes relatifs aux systèmes de nombres à virgule flottante	
I.4. Conclusion	
II - Calcul de la somme d'une série numérique alternée dont le terme général décroît	127
II.1. Introduction et notations	
II.2. Etude préliminaire	
II.3. Théorème	
II.4. Exemples numériques et mode opératoire.	

CHAPITRE VI UTILISATION CONJOINTE DE DEUX ARITHMETIQUES CORRECTES

I - Définitions	144
II - Suites maximisantes d'opérations sur S relatives aux algorithmes usuels	145
II.1. Somme algébrique et produit scalaire	
II.2. Evaluation numérique d'un polynôme	
II.3. Calcul de la racine carrée d'un nombre par la méthode de Newton.	
II.4. Décomposition LR et inverse d'une M - matrice	
II.4.1 Décomposition LR d'une M - Matrice	
II.4.2 Résolution du système linéaire $AX = B$, A étant une M - matrice, B un second membre positif (resp. négatif)	
II.4.3 Remarque : Algorithme classique de Gauss appliqué à la résolution de $AX = B$, A M - matrice, $B > 0$	
II.5. Inverse d'une matrice tridiagonale, symétrique, définie positive	
II.6. Inversion par la méthode des sous-matrices, d'une matrice à blocs diagonaux M - matrices, blocs extra-diagonaux positifs.	
II.7. Remarques	

III - Un exemple de suites non maximisantes d'opérations sur S générant des images encadrant la solution réelle	165
III.1. Résolution d'un système triangulaire-notations.	
III.2. Détermination de suites d'encadrement de \mathcal{G}	
III.2.1 L'erreur e_i	
III.2.2 Théorème	
III.2.3 Exemple et contre-exemple	
III.2.4 Une autre formulation des hypothèses	
III.2.5 Cas où les suites d'encadrement de \mathcal{G} sont maximisantes	
III.2.6 Exemple numérique	
IV - Théorèmes généraux	174
V - Conclusion	178

CHAPITRE VII UNE ETUDE SUR LA REPARTITION DE L'ERREUR D'ARRONDI

I - Fonction de répartition de l'erreur d'arrondi considérée comme variable aléatoire	180
I.1. Formulation du problème	
I.2. L'addition arithmétique	
I.3. L'addition de deux nombres de signes contraires	
I.3.1 Le cas de la troncature sans chiffre de garde	
I.3.2 Le cas de la troncature avec un chiffre de garde	
I.4. Conclusion relative à l'addition.	
I.5. Le cas de la multiplication	
II - Le problème d'un meilleur choix de la base	199

ANNEXE CARACTERISTIQUES DE QUELQUES ARITHMETIQUES A VIRGULE

FLOTTANTE	203
BIBLIOGRAPHIE	208
INDEX DES TERMES ET NOTATIONS	212
TABLE DES MATIERES	214