

JEAN MEINGUET

Structure et estimations de coefficients d'erreurs

Publications des séminaires de mathématiques et informatique de Rennes, 1977, fascicule S4

« Journées éléments finis », , p. 1-22

http://www.numdam.org/item?id=PSMIR_1977__S4_A9_0

© Département de mathématiques et informatique, université de Rennes, 1977, tous droits réservés.

L'accès aux archives de la série « Publications mathématiques et informatiques de Rennes » implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques
<http://www.numdam.org/>

Structure et Estimations de Coefficients d'Erreurs

Jean Meinguet

1. Introduction.

Les erreurs étant inhérentes à tout processus d'approximation, il est évidemment essentiel de pouvoir les estimer *quantitativement*. Ceci suppose en fait l'existence de méthodes générales, capables d'assurer une précision suffisante à un coût acceptable. La recherche d'un compromis à cet égard en analyse numérique (linéaire) nous a conduit à introduire récemment (cf. [4, 5]) une *méthode générale de majoration réaliste des erreurs d'approximation*. Ainsi que les quelques résultats déjà publiés (cf. [4, 5] et [1, 3]) le montrent, cette méthode permet notamment l'estimation pratique des diverses *constantes génériques* que l'on rencontre à propos des méthodes d'éléments finis (cf. en particulier [2], Théorèmes 2, 4, 5, 6).

Le présent travail poursuit l'élaboration systématique de cette méthode de majoration quantitative des erreurs dont les idées maîtresses sont dégagées ici à la fin de la Section 2, laquelle est consacrée à un résumé des éléments essentiels de l'analyse abstraite sous-jacente que nous avons développée et explicitement motivée dans [5]. On sait que la théorie des opérateurs linéaires dans les espaces vectoriels normés fournit un cadre naturel à cette analyse, qui se réfère par ailleurs à ces outils classiques que sont

Conférence faite aux *Journées Eléments Finis, Rennes,*
4-6 Mai 1977.

le *Peano kernel theorem* et sa généralisation connue sous l'appellation de *Branble-Hilbert lemma*.

Les deux dernières sections traitent de deux grandes classes d'applications particulièrement importantes en approximation linéaire multivariée, à savoir la classe des *approximations ponctuelles ou uniformes* dans des espaces de fonctions continûment dérivables (Section 3) et la classe des *approximations en moyenne quadratique* dans des espaces de Sobolev (Section 4). Notre but est ici de contribuer directement à l'obtention ultérieure de majorations quantitatives des erreurs d'approximation (et pas seulement d'interpolation) dans des situations concrètes, en fournissant un jeu de *majorations clefs*. Ainsi que nous l'avons montré dans [4, 5], de telles majorations résultent aisément de la forme intégrale du reste de la formule de Taylor (pour autant bien sûr qu'elle soit applicable !). Dans les espaces de Sobolev, il faut en fait recourir à une modification fondamentale de cette formule de Taylor. Celle que nous proposons à la Section 4 s'avère particulièrement naturelle et commode; en dépit de similitudes apparentes, elle diffère profondément de la version proposée et exploitée dans [1, 3] et échappe d'ailleurs totalement à la limitation soulignée dans [1] (cf. Remarque 2-3, p. 16).

2. Structure des coefficients d'erreurs.

Les problèmes de majoration qui nous intéressent ici relèvent d'*inégalités* de la forme très générale suivante :

$$(1) \quad \|Rf\|_Z \leq c \|Uf\|_Y, \quad \forall f \in X,$$

où c doit être une constante numérique *finie*, X , Y et Z sont des espaces vectoriels donnés (sur \mathbb{R} ou sur \mathbb{C}), $\|\cdot\|_Y$ et $\|\cdot\|_Z$ sont des normes données sur Y et sur Z , $U : X \rightarrow Y$ et

$R : X \rightarrow Z$ sont des applications *linéaires* données, U étant toujours supposée *surjective* (ce qui revient à la définition : $Y := U(X)$). On notera que (1) peut aussi bien s'écrire sous la forme :

$$(1 \text{ bis}) \quad \|Rf\|_Z \leq c q_X(f), \quad \forall f \in X,$$

où $q_X(\cdot)$ est une semi-norme sur X . Sous l'hypothèse naturelle

(H1) R continue par rapport à U ,

la majoration (1) d'une norme de Rf (considérée comme *inconnu*) en fonction d'une norme de Uf (considéré comme *connu*) n'est évidemment satisfaisante que si la majoration du *coefficient théorique d'erreur*

$$(2) \quad c_0 := \sup_{\substack{f \in X \\ Uf \neq 0}} \|Rf\|_Z / \|Uf\|_Y \equiv \min \{c \in \mathbb{R} \text{ vérifiant (1)}\}$$

par le nombre c est réaliste. Le problème que nous nous posons en définitive est dès lors celui de la détermination quantitative de "bonnes" majorations pour c_0 .

Une expression intéressante de c_0 s'obtient aisément à partir de la relation d'inclusion

$$(3) \quad K \equiv \text{Ker } U \subset \text{Ker } R$$

entre sous-espaces vectoriels de X , relation qui résulte trivialement de (1) sous l'hypothèse (H1). En vertu d'un *théorème de factorisation* élémentaire (provenant pour l'essentiel de la théorie des ensembles), il existe alors (et seulement alors) une et une seule application $Q : Y \rightarrow Z$ telle que

$$(4) \quad R = QU;$$

on sait en outre que Q est linéaire, surjective ssi R est surjective, injective ssi $\text{Ker } U = \text{Ker } R$. Etant donné l'équivalence classique (moyennant l'axiome de choix) entre les propriétés de surjectivité et d'inversibilité à droite d'une application, équivalence dont la factorisation qui précède résulte en fait directement, il existe au moins une application $V : Y \rightarrow X$ (éventuellement non-linéaire mais toujours injective) telle que l'on ait

$$(5) \quad UV = 1_y,$$

où 1_y désigne l'application identique de Y . Il s'ensuit que Q admet les représentations explicites :

$$(6a) \quad Q = RV, \quad \forall V : Y \rightarrow X \text{ vérifiant (5),}$$

d'où l'on déduit finalement le résultat annoncé, soit

$$(7) \quad c_0 := \|Q\| \equiv \|RV\|, \quad \forall V : Y \rightarrow X \text{ vérifiant (5),}$$

où

$$(8) \quad \|Q\| := \sup_{0 \neq g \in Y} \|Qg\|_Z / \|g\|_Y$$

s'interprète classiquement comme la *norme* de l'application linéaire (continue) $Q : Y \rightarrow Z$.

Une autre conséquence importante du théorème de factorisation est à rappeler ici. A savoir le fait que toute application T d'un ensemble E dans un ensemble F admet une *décomposition canonique*, laquelle peut se représenter (en notations mnémoniques) par le diagramme suivant :

$$(9) \quad E \xrightarrow{T_{cs}} E/T \xrightarrow{T_b} T(E) \xrightarrow{T_{ci}} F;$$

T désigne ici la relation d'équivalence dans E associée à T (c'est-à-dire la relation " $e_1 \in E, e_2 \in E, Te_1 = Te_2$ "), T_{cs} est la *surjection canonique* de E sur l'ensemble quotient E/T , T_b est la *bijection associée* à T , T_{ci} est l'*injection canonique* de l'image $T(E)$ dans F . Comme exemple directement utile d'application de (9), relevons que (6a) peut se récrire comme suit :

$$(6b) \quad Q = R_q U_b^{-1} \text{ avec } R_q = RVU_b,$$

$$(6c) \quad Q = R_r V_b \text{ avec } R_r = R V_{ci},$$

où $R_q : X/K \rightarrow Z$ n'est autre que l'*application quotient* issue de la factorisation de $R : X \rightarrow Z$ par $U_{cs} : X \rightarrow X/K \equiv X/U$ tandis que $R_r : V(Y) \rightarrow Z$ est simplement la *restriction* de $R : X \rightarrow Z$ à la partie $VU(X)$ de X .

Sauf dans le cas particulier (mais important !) où les espaces Y et Z sont préhilbertiens, pour lequel une approche variationnelle directe peut s'avérer préférable, il semble que la connaissance explicite de $Q : Y \rightarrow Z$ soit requise pour pouvoir évaluer c_0 . En fait, l'expression cherchée de Rf en fonction de Uf peut s'obtenir assez simplement : étant donné l'identité fondamentale (résultant de (4) et (6))

$$(10) \quad Rf \equiv RVUf, \quad \forall f \in X,$$

il suffit d'appliquer R à VUf , où $V : Y \rightarrow X$ peut même être choisi librement pour autant qu'il vérifie (5). Fort malheureusement, cette abstraction de l'idée maîtresse du *Peano kernel theorem*, si elle est évidemment susceptible

de conduire en définitive à une majoration de c_0 , ne permet qu'exceptionnellement de déterminer le coefficient d'erreur lui-même; ceci en raison de la complexité des possibilités de résonance entre "fonctions" Uf et "noyaux" RV . Le caractère foncièrement théorique de c_0 se trouve encore renforcé du fait que ce coefficient dépend non seulement de X, Y, U , *données essentiellement fixes* pour lesquelles une véritable standardisation s'est imposée (U est typiquement un opérateur différentiel et X un espace fonctionnel tel que $C^m(\Omega)$ ou $H^m(\Omega)$), mais aussi de R , *donnée essentiellement variable* (R est typiquement défini comme l'erreur associée à l'approximation d'un opérateur spécifique par un autre). Cette discussion motive l'intérêt pratique qu'il peut y avoir à se restreindre à la sous-classe suivante d'inégalités de type (1) :

$$(11) \quad \|Rf\|_Z \leq \|R\| d \|Uf\|_Y, \forall f \in X,$$

où d doit être une constante numérique *finie et indépendante* de $R : X \rightarrow Z$; ceci n'a un sens que si X est muni d'une norme $\|\cdot\|_X$ telle que

(H2a) R est continue,

(H2b) il existe au moins un inverse à droite *borné* V de U ,

cette double hypothèse devant ici s'ajouter à (H1). Le problème fondamental qui se pose dès lors est celui de l'estimation réaliste du *coefficient pratique d'erreur*

$$(12) \quad d_0 := \inf \{d \in \mathbb{R} \text{ vérifiant (11), } \forall R \in L(X; Z)\}$$

pour divers *choix standards* des données fixes; quant à l'évaluation quantitative de la norme $\|R\|$, elle constitue un problème essentiellement *spécifique* et doit être traitée en conséquence.

Il se trouve que d_0 est la solution commune d'un problème de minimisation et de son dual, ce qui suggère entre autres une technique naturelle d'encadrement. Compte tenu de (10), on peut écrire parallèlement à l'inégalité triviale

$$(13a) \quad \|Rf\|_Z \leq \|R\| \|f\|_X, \quad \forall f \in X,$$

l'inégalité généralement distincte

$$(13b) \quad \|Rf\|_Z \leq \|R\| \|VUf\|_X, \quad \forall f \in X,$$

d'où résulte une inégalité de type (11), soit

$$(13c) \quad \|Rf\|_Z \leq \|R\| \|V\| \|Uf\|_Y, \quad \forall f \in X,$$

pour tout inverse à droite *borné* de U (même si un tel V est non-linéaire, sa norme $\|V\|$ doit être interprétée ici conformément à la définition (8)). Il apparaît dès lors que d_0 est la solution du problème de *minimisation*

$$(14a) \quad d_0 = \inf \|V\|, \quad \forall V \text{ vérifiant (5),}$$

lequel est équivalent au problème de *min-max*

$$(14a \text{ bis}) \quad d_0 = \inf_V \sup_{\|Uf\|_Y=1} \|VUf\|_X,$$

où l'on peut en fait toujours supposer que V est *homogène*. Par ailleurs, la majoration fournie par (13a) peut être améliorée grâce à (3), l'optimum à cet égard étant réalisé par

$$(13d) \quad \|Rf\|_Z \leq \|R\| q(U_{cs} f), \quad \forall f \in X,$$

où

$$(15) \quad q(U_{cs} f) := \inf_{k \in K} \|f - k\|_X$$

n'est autre que la semi-norme dont est naturellement muni l'espace vectoriel quotient X/K (on sait que $q(\cdot)$ est une norme ssi K est fermé dans X , donc en particulier si K est de dimension finie ou si $U : X \rightarrow Y$ est continue). L'identité $U_{cs} = U_b^{-1} U$ conduit dès lors à une définition de d_0 comme solution du problème de *max-min*

$$(14b) \quad d_0 = \|U_b^{-1}\| = \sup_{\|Uf\|_Y=1} \inf_{k \in K} \|f - k\|_X,$$

que l'on peut en fait récrire sous la forme

$$(14b \text{ bis}) \quad d_0 = \sup_{\|Uf\|_Y=1} \inf_V \|VUf\|_X,$$

duale de (14a bis). En définitive, les inégalités (13c) peuvent être considérées comme résultant d'un affaiblissement de deux inégalités optimales et généralement distinctes, à savoir (13d) et

$$(13e) \quad \|Rf\|_Z \leq \|RV\| \|Uf\|_Y, \quad \forall f \in X.$$

Il est manifeste que diverses *équivalences conditionnelles entre normes* ont à jouer ici un rôle essentiel, notamment celles qui peuvent résulter des inégalités optimales suivantes (cf. [5], p. 314) :

- si U est bornée (ou continue), alors

$$(16a) \quad \|Uf\|_Y \leq \|U\| q(U_{cs} f),$$

$$(17a) \quad \|Uf\|_Y \leq \|V_b^{-1}\| \|VUf\|_X \leq \|U\| \|VUf\|_X;$$

- si V est un inverse à droite borné de U , alors

$$(16b) \quad q(U_{cs} f) \leq \|U_b^{-1}\| \|Uf\|_Y \leq \|V\| \|Uf\|_Y,$$

$$(17b) \quad \|VUf\|_X \leq \|V\| \|Uf\|_Y;$$

- si l'application composée VU est bornée, alors

$$(18a) \quad \|VUf\|_X \leq \|VU\| q(U_{cs} f),$$

tandis que, de toute façon,

$$(18b) \quad q(U_{cs} f) \leq \|VUf\|_X.$$

L'importance théorique et pratique de ces comparaisons entre normes est de faciliter une compréhension profonde des diverses possibilités de généralisation du "Peano kernel theorem", en particulier du célèbre "Bramble-Hilbert lemma" (dont la formulation abstraite correspond en fait à la situation où les applications U et V liées par (5) sont toutes deux linéaires et continues).

Notons encore que si l'on pose

$$(19) \quad f \equiv Pf + VUf, \quad \forall f \in X,$$

V est un inverse à droite homogène (resp. linéaire) de U ssi P est un *projecteur* homogène (resp. linéaire) de X d'image K . Les relations (16b) et (17b) impliquent alors, pour autant qu'il existe un inverse à droite borné de U , l'existence pour tout $f \in X$ d'un élément de K (noté Pf) tel que

$$(20) \quad \inf_{k \in K} \|f - k\|_X \leq \|f - Pf\|_X \leq d \|Uf\|_Y,$$

où $d (\geq d_0)$ est une constante finie; cette présentation abstraite, qui met l'accent sur l'approximant $Pf \in K$ de f plutôt que sur l'erreur associée VUf , se réduit en fait, par une particularisation appropriée des données fixes (X et Y : espaces de Sobolev, U : opérateur différentiel), à la présentation classique du "Lemme technique" dont le Lemme de Bramble-Hilbert est une conséquence directe. Si K est de dimension finie (comme c'est la règle en pratique), il existe nécessairement un projecteur de meilleure approximation P_0 (associé par (19) à un inverse à droite V_0 de U), qui peut être non-linéaire bien qu'homogène (un exemple non-linéaire familier est celui du projecteur de meilleure approximation polynomiale au sens de Tchebycheff); on obtient alors pour d_0 une expression plus concrète que (14), à savoir

$$(21) \quad d_0 := \|V_0\| \equiv \|(1_X - P_0)V_p\|, \quad \forall V_p : Y \rightarrow X \text{ vérifiant (5)}.$$

En conclusion de cette analyse, dégageons les axes de la méthode générale de majoration quantitative des erreurs que nous voulons proposer. On commence par choisir (ce qui est toujours possible si K est de dimension finie) une décomposition de l'espace normé X en une somme directe topologique de la forme

$$(22) \quad X = K \oplus \tilde{X};$$

ceci équivaut à choisir une formule de représentation de type (19), où V est un inverse à droite de U et P un projecteur de X d'image K , ces deux applications étant linéaires et continues. On décompose ensuite la restriction R_r à $\tilde{X} \equiv VU(X)$ de l'application linéaire spécifique $R : X \rightarrow Z$ (qu'il s'agit en définitive d'"estimer") en une somme finie

$$(23) \quad R_r = \sum_j S_j$$

d'applications linéaires *standards* $S_j : \tilde{X} \rightarrow Z$, de façon telle que chacune des applications composées $S_j V : Y \rightarrow Z$ soit continue (il est dans la nature des choses que de telles décompositions soient possibles dans la pratique courante). On procède à l'évaluation de constantes d_j telles que

$$(24) \quad \|S_j V U f\|_Z \leq d_j \|U f\|_Y, \quad \forall f \in X,$$

soit directement, soit en recourant à une collection élaborée par ailleurs de *majorations clefs* ayant trait notamment aux opérations courantes de dérivation, d'interpolation, d'intégration et à leurs composées. On obtient ainsi finalement la majoration quantitative annoncée

$$(25) \quad \|R f\|_Z \leq \left(\sum_j d_j \right) \|U f\|_Y, \quad \forall f \in X,$$

et cela sans avoir jamais dû tenter l'évaluation de $\|R\|$; il importe de noter que l'hypothèse (H2), en particulier la continuité de R , n'est nullement essentielle ici, à la différence de l'hypothèse (H1).

3. Exemples de majorations clefs dans $C^m(\bar{\Omega})$.

Soit Ω un *ouvert borné* (non vide) de \mathbb{R}^n tel que son adhérence $\bar{\Omega}$ soit *étoilée* par rapport à un point noté a (la dernière restriction pourrait être levée assez facilement, pour autant que Ω reste connexe, mais elle est en fait mineure). En matière d'approximations ponctuelles ou uniformes, il est à la fois naturel et commode de prendre comme *données fixes* X et U , respectivement :

- l'espace $C^m(\bar{\Omega})$ (avec m entier ≥ 1) des fonctions numériques f qui sont *uniformément continues* dans Ω ainsi que toutes leurs dérivées partielles $\partial^\alpha f$ d'ordre $|\alpha| \leq m$ (nous utilisons ici les notations multi-indicielles classiques); muni de la famille de semi-normes uniformes

$$(26a) \quad |f|_j = \max_{|\alpha|=j} \max_{x \in \bar{\Omega}} |\partial^\alpha f(x)|, \quad 0 \leq j \leq m,$$

équivalente à la norme de Tchebycheff

$$(26b) \quad \|f\|_m = \max_{0 \leq j \leq m} |f|_j,$$

cet espace est de Banach.

- l'opérateur D^m de dérivation d'ordre m au sens de Fréchet; le noyau K de U dans X est dès lors l'espace vectoriel P_{m-1} , de dimension

$$(27a) \quad N = \binom{n+m-1}{n},$$

des polynômes à n indéterminées (à coefficients réels) de degré total $\leq m-1$; quant à l'espace $Y := D^m(C^m(\bar{\Omega}))$, muni de la norme

$$(28) \quad \|D^m f\|_Y = |f|_m,$$

il peut s'identifier par l'isométrie naturelle $D^m f \mapsto \{\partial^\alpha f : |\alpha| = m\}$ à un sous-espace fermé de $(C^0(\bar{\Omega}))^{\bar{N}}$ où

$$(27b) \quad \bar{N} = \binom{n+m-1}{m}.$$

Pour le terme \tilde{X} dans la décomposition (22), il s'avère particulièrement intéressant de choisir l'espace vectoriel

$$(29) \quad C_a^m(\bar{\Omega}) := \{f \in C^m(\bar{\Omega}) : \partial^\alpha f(a) = 0, 0 \leq |\alpha| < m\};$$

ceci revient à choisir pour (19) la *formule de Taylor d'ordre m relative au point a*, le projecteur P de $C^m(\bar{\Omega})$ d'image \mathcal{P}_{m-1} étant dès lors défini par

$$(30a) \quad Pf(x) := \sum_{|\alpha|=0}^{m-1} \partial^\alpha f(a) (x-a)^\alpha / \alpha!$$

et l'inverse à droite V de D^m par

$$(30b) \quad V(D^m f)(x) := \int_0^1 \frac{(1-t)^{m-1}}{(m-1)!} D^m f(a+t(x-a)) \cdot (x-a)^m dt,$$

ainsi qu'il résulte de la forme intégrale du reste pour tout $x \in \bar{\Omega}$. On vérifiera sans peine que $\partial^\alpha (VD^m f)$, dérivée partielle quelconque d'ordre $< m$ d'un élément arbitraire de $C_a^m(\bar{\Omega})$, peut semblablement s'exprimer sous la forme intégrale

$$(31) \quad \partial^\alpha (VD^m f)(x) = \int_0^1 \frac{(1-t)^{m-1-|\alpha|}}{(m-1-|\alpha|)!} D^{m-|\alpha|} \partial^\alpha f(a+t(x-a)) \cdot (x-a)^{m-|\alpha|} dt,$$

en fonction seulement de (certaines) dérivées partielles d'ordre m de $f \in C^m(\bar{\Omega})$ ou, ce qui revient au même (puisque $f - VD^m f \in \mathcal{P}_{m-1}$), de $VD^m f \in C_a^m(\bar{\Omega})$.

On sait que toute expression du type $D^m f(y) \cdot (x-a)^m$ (où $y \in \mathbb{R}^n$) s'interprète simplement comme le produit de la matrice ligne formée avec les coordonnées du tenseur symétrique $D^m f(y)$ par la matrice colonne formée avec les coordonnées du produit de Kronecker $(x-a)^m$, ces coordonnées se rapportant à la base canonique de $\mathbb{R}^{(n^m)}$ (pour une analyse algébrique détaillée, cf. [4], Section 3). Compte

tenu de cette *interprétation matricielle*, les majorations clefs annoncées résultent directement de (31) par une simple application (sous le signe \int) d'un cas limite de l'inégalité de Hölder (pour les sommes); pour tout $f \in C^m(\bar{\Omega})$ et tout $\alpha \in N^n$ d'ordre $|\alpha| \leq m$, on obtient en définitive les inégalités (précises) suivantes :

$$(32a) \quad |\partial^\alpha \text{VD}^m f(x)| \leq \{ \|x-a\|_1^{m-|\alpha|} / (m-|\alpha|)! \} |f|_m, \quad \forall x \in \bar{\Omega},$$

$$(32b) \quad \|\partial^\alpha \text{VD}^m f\|_0 \leq |\text{VD}^m f|_{|\alpha|} \leq \{ h_1^{m-|\alpha|} / (m-|\alpha|)! \} |f|_m$$

où l'on a posé

$$(33) \quad h_1 := \max_{x \in \bar{\Omega}} \|x-a\|_1,$$

$\|\cdot\|_1$ désignant ici la 1-norme de Hölder sur \mathbb{R}^n . Il s'ensuit globalement l'inégalité de type (17b) :

$$(34) \quad \|g\|_m \leq \|V\| |g|_m \quad \text{avec} \quad \|V\| \leq d = \max_{0 \leq j \leq m} (h_1^j / j!), \quad \forall g \in C_a^m(\bar{\Omega}),$$

laquelle permet de préciser quantitativement l'équivalence sur $C_a^m(\bar{\Omega})$ des normes $|\cdot|_m$ et $\|\cdot\|_m$. Quant à l'application concrète des majorations (32), elle a été abordée dans [4] (cf. Section 5) à propos de l'intégration numérique multivariée et de l'interpolation bidimensionnelle de Lagrange sur un triangle; nous ne rappellerons pas ici les résultats obtenus à cette occasion.

4. Exemples de majorations clefs dans $H^m(\Omega)$.

Soit Ω un ouvert borné convexe (non vide) de \mathbb{R}^n (il existe ici également diverses possibilités de généralisation), de diamètre euclidien h et de mesure (de Lebesgue) S

Nous prendrons comme *données fixes* X et U, respectivement :

- l'espace de Sobolev $H^m(\Omega)$ (avec m entier ≥ 1) des (classes de) fonctions numériques $f \in L_2(\Omega)$ telles que $\partial^\alpha f \in L_2(\Omega)$ pour $|\alpha| \leq m$, les dérivées $\partial^\alpha f$ étant prises au sens des distributions sur Ω ; muni de la famille de semi-normes

$$(35a) \quad |f|_j = \left(\sum_{i_1, \dots, i_j=1}^n \int_{\Omega} \left| \frac{\partial^j f(x)}{\partial x_{i_1} \dots \partial x_{i_j}} \right|^2 dx \right)^{1/2}, \quad 0 \leq j \leq m,$$

équivalente à la norme

$$(35b) \quad \|f\|_m = \left(\sum_{j=0}^m |f|_j^2 \right)^{1/2},$$

$H^m(\Omega)$ est un espace de Hilbert dans lequel $C^m(\bar{\Omega})$ est dense.

- l'opérateur D^m de dérivation totale d'ordre m au sens des distributions sur Ω ; comme Ω est connexe, le noyau K de U dans X est à nouveau l'espace de polynômes P_{m-1} ; quant à l'espace $Y := D^m(H^m(\Omega))$, muni de la norme naturelle

$$(36) \quad \|D^m f\|_Y = |f|_m$$

définie par (35a), il peut s'identifier par l'isomorphisme (non isométrique ici !) $D^m f \mapsto \{\partial^\alpha f : |\alpha| = m\}$ à un sous-espace fermé de $(L_2(\Omega))^{\bar{N}}$.

Pour le projecteur (linéaire continu) P de $H^m(\Omega)$ d'image P_{m-1} et l'inverse à droite V de D^m , à associer dans une formule de représentation de type (19), nous pro-

posons ici les définitions suivantes :

$$(37a) \quad Pf(x) := S^{-1} \sum_{|\alpha|=0}^{m-1} \int_{\Omega} \{ \partial^{\alpha} f(a) (x-a)^{\alpha} / \alpha! \} da,$$

$$(37b) \quad V(D^m f)(x) := S^{-1} \int_0^1 \frac{(1-t)^{m-1}}{(m-1)!} \left\{ \int_{\Omega} D^m f(a+t(x-a)) \cdot (x-a)^m da \right\} dt,$$

lesquelles résultent tout simplement de l'intégration sur Ω par rapport à a des expressions respectives (30a) et (30b) où $x \in \bar{\Omega}$ et $f \in C^m(\bar{\Omega})$. Il est manifeste que la définition (37a) se prolonge par continuité à $H^m(\Omega)$; il en va de même de la définition (37b) ssi la condition

$$(37c) \quad m > n/2$$

est vérifiée (cas particulier du Théorème d'immersion de Sobolev). Comme les raisons qui justifient l'introduction de (31) à partir de (30b) sont également d'application ici, nous pouvons considérer (37b) comme un cas particulier de

$$(38a) \quad \partial^{\alpha} (VD^m f)(x) = S^{-1} \int_0^1 \frac{(1-t)^{m-1-|\alpha|}}{(m-1-|\alpha|)!} \left\{ \int_{\Omega} D^{m-|\alpha|} \partial^{\alpha} f(a+t(x-a)) \cdot (x-a)^{m-|\alpha|} da \right\} dt,$$

la condition (37c) devant alors être remplacée par

$$(38b) \quad |\alpha| < m - n/2.$$

La concrétisation par (37) de la formule générale (19) n'est pas strictement nouvelle, encore qu'on ne la rencontre guère dans la littérature (une exception notable : cf. [6], p. 79) et que jusqu'à présent elle ne semble

jamais avoir été exploitée pour obtenir des majorations quantitatives; d'autres particularisations de (19), et donc de la décomposition de $H^m(\Omega)$ en une somme directe topologique de type (22), sont apparemment plus connues (cf. entre autres [9], p. 50 et [8], p. 111), bien que leur exploitation pratique soit beaucoup plus difficile; sans doute est-ce dû à la complexité apparente de la définition

$$(39) \quad \tilde{X} := \{f \in H^m(\Omega) : \sum_{|\alpha|=0}^{m-1-|\beta|} \int_{\Omega} \{\partial^{\alpha+\beta} f(a) (0-a)^{\alpha}/\alpha!\} da, 0 \leq |\beta| \leq m-1,$$

, impliquée par (37).

Compte tenu de l'interprétation matricielle que nous avons rappelée lors de la recherche des inégalités (32a, b), il est facile de déduire de (38a), pour tout $f \in H^m(\Omega)$ et sous l'hypothèse (38b), les majorations suivantes :

$$(40) \quad |\partial^{\alpha} \nabla D^m f(x)| \leq \{S^{-1/2} h^{m-|\alpha|} / [(m-|\alpha|-1)!(m-|\alpha|-n/2)]\} |f|_m, \forall x \in \bar{\Omega};$$

il suffit en fait d'appliquer l'inégalité de Cauchy-Schwarz sous le signe \int_{Ω} et de majorer uniformément par h la norme euclidienne de $x-a$, d'opérer ensuite le changement de variables $a \mapsto y = a + t(x-a)$ de jacobien $(1-t)^n$, d'appliquer l'inégalité de Schwarz à l'intégrale ainsi transformée en notant bien qu'elle est à prendre sur l'ensemble

$$(41) \quad \Omega_{x,t} := \{y = a+t(x-a) : a \in \Omega, (x,t) \text{ fixé dans } \bar{\Omega} \times [0,1]\}$$

de mesure $(1 - t)^n S$ (qui est l'image* de Ω par une homothétie de centre x et de rapport $1 - t$), de procéder finalement à l'intégration par rapport à t . Comme le recours à l'inégalité de Schwarz dans cette méthode (déjà utilisée dans [6], cf. p. 79) n'est pas intervenu en tout dernier lieu, il est clair que le résultat (40) n'est pas optimal. Il est donc naturel de songer à intervertir les intégrations dans (38a); en exploitant alors correctement le fait (suggéré par (41)) que l'intégrale "double" (38a) sur $\mathbb{R}^n \times (0, 1)$ est à prendre sur le cône ouvert de base $\Omega \times \{0\}$ et de sommet $(x, 1)$, on obtient en définitive, pour tout $f \in H^m(\Omega)$ et sous l'hypothèse (38b), les majorations optimales suivantes :

$$(42) \quad |\partial^\alpha \text{VD}^m f(x)| \leq |D^{|\alpha|} \text{VD}^m f(x)| \leq \frac{S^{-1/2} h^{m-|\alpha|}}{(m-|\alpha|)!} \binom{m-|\alpha|}{m-|\alpha|-n/2}^{1/2} \|f\|_m, \quad \forall x \in \bar{\Omega},$$

qu'il y a lieu de substituer définitivement aux majorations (40).

Le fait que le projecteur P de $H^m(\Omega)$ d'image P_{m-1} soit toujours continu, que la condition (37c) soit satisfaite ou non, implique l'équivalence sur le sous-espace vectoriel \tilde{X} défini par (39) des normes de Sobolev $|\cdot|_m$ et $\|\cdot\|_m$ définies par (35a, b) (démonstration par l'absurde, exploitant le Lemme de compacité de Rellich, très semblable en fait à la démonstration classique du "Lemme technique" dont il a été question à propos de l'inégalité (20), cf. notamment [6], p. 86 et [8], p. 111). Il y a donc a priori un sens à vouloir majorer les normes $|\cdot|_j$ de $\text{VD}^m f \in \tilde{X}$, pour $0 \leq j \leq m - 1$, en fonction de la norme $|\cdot|_m$ de $\text{VD}^m f$ ou, ce qui revient au même (puisque $f - \text{VD}^m f \in P_{m-1}$), de $f \in H^m(\Omega)$. On obtient en définitive, à partir de (38a), les majorations suivantes :

$$(43a) \quad |VD^m f|_j \leq d(m, n, j) \{h^{m-j}/(m-j)!\} |f|_m,$$

valables pour tout $f \in H^m(\Omega)$ et $0 \leq j \leq m - 1$, où

$$(43b) \quad d(m, n, j) := (m - j) \min_p \left\{ \frac{\int_0^1 (1-t)^{2p} \min[t^{-n}, (1-t)^{-n}] dt}{2m - 2j - 2p - 1} \right\}^{1/2};$$

il suffit en fait d'appliquer l'inégalité de Cauchy-Schwarz sous le signe \int_{Ω} dans (38a) et de majorer uniformément par h la norme euclidienne de $x - a$, d'appliquer ensuite l'inégalité de Schwarz au carré de l'expression (38a) (intégrale "double" sur le cylindre $\Omega \times (0,1)$ de $\mathbb{R}^n \times (0,1)$), interprétée comme produit scalaire par rapport à la mesure $(1-t)^{2p} dt$, où $2p > -1$), d'opérer enfin à propos de $\int_{\Omega} \int_{\Omega} |D^m f(a + t(x - a))|^2 dx$ du changement de variables $a + y = a + t(x - a)$ (resp. $x + y = a + t(x - a)$) si $t \leq 1/2$ (resp. $t \geq 1/2$) de sorte que cette intégrale se trouve être majorée simplement par $\min[t^{-n}, (1-t)^{-n}] S |f|_m^2$. On peut donner de la constante $d(m, n, j)$ diverses majorations, notamment les suivantes :

$$(43c) \quad d(m, n, j) < 2(m - j) 2^{n/2} / 2^{m-j},$$

$$(43d) \quad d(m, n, j) < \left(\frac{2^{n-1} - 1}{n - 1} \right)^{1/2} \frac{m - j}{(2m - 2j - 2)^{1/2}} \text{ si } m - j \geq 2,$$

$$(43d \text{ bis}) \quad d(m, n, j) < \left(\frac{2^{n-1} - 1}{n - 1} \right)^{1/2} \text{ si } m - j = 1,$$

sans oublier celles qui résultent trivialement de (42), à savoir :

$$d(m, n, j) \leq [(m - j)/(m - j - n/2)]^{1/2} \quad \text{si } m - j > n/2.$$

Considérons pour terminer, en guise d'illustration l'application des inégalités clefs (42) et (43) en fonction de $|f|_2$ de la semi-norme $|\cdot|_1$ pour d'interpolation $f - \Pi f$ de $f \in H^2(\Omega)$, où est ici un triangle quelconque de sommets $(1, 2, 3)$, par le polynôme

$$\Pi f(x) := \sum_{i=1}^3 f(x_i) p_i(x), \quad x \in \bar{\Omega},$$

où p_i qui interpole f aux points x_i ; on sait que les $p_i(x) \in P_1$ de la base de Lagrange se réduisent en coordonnées barycentriques de x par rapport aux sommets. Comme $P_1 \subset \text{Ker}(I - \Pi)$, il vient directement à partir des lemmes généraux (23, 24, 25) et des majorations

$$\begin{aligned} |f - \Pi f|_1 &\leq |VD^2 f|_1 + \sum_{i=1}^3 |VD^2 f(x_i)| |p_i|_1 \\ &\leq 2^{1/2} h |f|_2 + (S^{-1/2} 2^{-1/2}) h^2 \left(\sum_{i=1}^3 |p_i|_1 \right) |f|_2; \end{aligned}$$

Par une discussion géométrique élémentaire, on peut obtenir

$$\sum_{i=1}^3 |p_i|_1 = 2S^{1/2}/\rho,$$

où ρ est le diamètre du cercle inscrit à Ω . D'où on obtient l'estimation cherchée :

$$(46) \quad |f - \Pi f|_1 < 2^{3/2} (h^2/\rho) |f|_2, \quad \forall f \in H^2(\Omega),$$

qui n'est que 25 % moins précise que celle obtenue par Natterer dans [7] (cf. son Théorème 1.2) pour un triangle quelconque; pour le triangle rectangle standard par contre, le coefficient de $|f|_2$ dans (46) doit être remplacé par $2^{3/2} (1 + 2^{1/2})$ et est donc nettement supérieur à celui obtenu par Natterer (soit 0.81) suite à une étude fort ingénieuse mais essentiellement spécifique, qu'il ne semble guère possible dès lors de généraliser, à la différence des résultats présentés ici.

BIBLIOGRAPHIE

1. Arcangéli, R., Gout, J.L. : Sur l'évaluation de l'erreur d'interpolation de Lagrange dans un ouvert de \mathbb{R}^n . R.A.I.R.O. Analyse Numérique 10, 5-27 (1976).
2. Ciarlet, P.G., Raviart, P.A. : General Lagrange and Hermite interpolation in \mathbb{R}^n with applications to finite element methods. Arch. Rational Mech. Anal. 46, 177-199 (1972).
3. Gout, J.L. : Sur l'estimation de l'erreur d'interpolation dans \mathbb{R}^n . Thèse 3e cycle, Université de Pau (1976).
4. Meinguet, J. : Realistic estimates for generic constants in multivariate pointwise approximation. In : Topics in numerical analysis II (J.J.H. Miller, ed.), pp. 89-107. London - New York : Academic Press 1975.
5. Meinguet, J., Descloux, J. : An operator-theoretical approach to error estimation. Numer. Math. 27, 307-326 (1977).
6. Morrey, C.B. : Multiple integrals in the calculus of variations. New York : Springer 1966.
7. Natterer, F. : Berechenbare Fehlerschranken für die Methode der finiten Elemente. In : ISNM 28, pp. 109-121. Basel : Birkhäuser 1975.
8. Nečas, J. : Les méthodes directes en théorie des équations elliptiques. Paris : Masson 1967.
9. Sobolev, S.L. : Applications of functional analysis in mathematical physics, Leningrad 1950 (traduction anglaise : Amer. Math. Soc., Transl. Math. Mono. 7, 1963).

Professeur Jean MEINGUET
Université de Louvain
Analyse Numérique et Programmation
chemin du Cyclotron, 2
B - 1348 LOUVAIN-la-NEUVE
BELGIQUE