

**NUMERICAL ANALYSIS  
FOR NONLINEAR AND BIFURCATION PROBLEMS**

by

Gabriel Caloz  
I.R.M.A.R.  
Université de Rennes I  
35042 Rennes Cedex

and

Jacques Rappaz  
D.M.A.  
École Polytechnique Fédérale de Lausanne  
1015 Lausanne

GC/JR

February 1994



## PREFACE

Computational applications generally involve nonlinear problems and often contain parameters. They may represent properties of the physical system they describe or quantities which can be varied. A basic problem in approximation consists in studying existence and convergence of approximated solutions for a given nonlinear problem, for instance when the parameters are fixed. Another problem is to represent the families or manifolds of solutions under variations of some parameters. Apart from a theoretical approach, such representations are computed and continuation methods are concerned with generating the solution manifolds. By varying one parameter, we can follow a path of solutions. Then to study the effects of change of parameters on a system, it is of prime interest to know the effects of numerical approximation on its behavior.

The goal of this article is to present a general framework in which approximations of nonlinear problems and approximations of solution manifolds can be studied. We will consider regular solutions, regular solution families, and singular solutions. Even though we will illustrate the general theory only with elementary finite element approximations of model boundary value problems, it can be applied to a much wider range of problems in connection with approximation methods. Our presentation is a remodelling of the one proposed by CROUZEIX and RAPPAZ [1989] taking its origin in DESCLOUX and RAPPAZ [1982] and in BREZZI, RAPPAZ, and RAVIART [1980, 1981a, 1981b].

The general problem we will handle and which covers a lot of applications is the following: find  $x \in X$  such that

$$F(x) = 0$$

where  $X$  and  $Z$  are Banach spaces,  $F : X \rightarrow Z$  is a smooth nonlinear mapping. Of particular interest is the case where the space  $X$  has the form  $\mathbb{R}^m \times Y$ , where  $\mathbb{R}^m$  with  $m \geq 1$  is the parameter space and the Banach space  $Y$  is the state space. We will work under the assumption that the derivative of  $F$  is a Fredholm operator of index  $n \geq 0$ . Both cases  $n = 0$  and  $n \geq 1$  with a surjective derivative are studied separately. Note that when  $n$  is positive, the family of solutions to  $F(x) = 0$  is a differentiable manifold. The singular situation with a not surjective derivative is also studied.

In the general setting, the approximation schemes are written in the form

$$F_h(x) = 0$$

where  $h$  is a parameter in  $(0, 1]$  and  $F_h : X \rightarrow Z$  is an approximation of  $F$ . The family  $\{F_h\}_{0 < h \leq 1}$  converges to  $F$  as  $h$  tends to 0. Our work then is to compare the solution set

of the problem  $F_h(x) = 0$  to the one of the problem  $F(x) = 0$ . The error estimates will generally be deduced from the inverse function theorem. Our program is the following.

In Sections 1 and 2, we present the functional analysis material we use to establish approximation results. It covers mainly variants of the inverse and of the implicit function theorems.

The sections 3, 4, and 5 present model examples which are used later to illustrate the general theory. They do not reflect the range of applications of the theory, but have been chosen because they are simple or familiar to the reader. In each case we discuss the exact problem and give some standard approximations of it.

Then we introduce a general framework to study the approximation of a nonlinear problem without parameter (or with fixed parameters). The originality of this approach is a generalization of the well-known inf-sup conditions to nonlinear problems. Furthermore we are naturally led up to discuss some a posteriori estimates in the nonlinear case, which are effective in an adaptive mesh procedure. The general case is studied in Section 6, Galerkin approximations are analyzed in Section 7 while examples are detailed in Sections 8 through 10.

Sections 11 through 17 are devoted to the study of the approximation of parametrized nonlinear problems. We consider exclusively regular solutions, that is points  $(\lambda, u) \in \mathbb{R}^m \times Y \equiv X$  where the derivative  $DF(\lambda, u)$  is surjective, as it is described in Section 11. First we study the simplest case of regular solution families, the ones parametrized by  $\lambda$ . The approximation in an abstract setting is developed in Section 12, Galerkin approximations are studied in Section 13 and an example is presented in Section 14. In the remainder of the chapter devoted to the approximation of regular solutions, simple limit points and their Galerkin approximations are also discussed. By introducing a new equation we extend the system and transform the problem to apply the results of Section 12.

In Sections 18 through 23 we are concerned with singular solutions and their approximations. Section 18 is devoted to the classical Lyapunov-Schmidt decomposition procedure, which leads to a bifurcation equation. A situation frequently encountered in the applications is presented in Section 19 where the bifurcation equation is studied carefully. Then approximations are analyzed in Section 20 and we get estimates between the bifurcation equation and the approximate bifurcation equations. In Section 21 we compare thoroughly the bifurcation equation and the approximate bifurcation equations when these are defined in  $\mathbb{R}^2$ . A simple model example illustrates in Section 22 the general theory. In Section 23 we address bibliographical complements.

As mentioned previously the emphasis here is put on the analysis of the convergence and the error estimates of the approximate problems. The methods necessary to understand the approximation of a parametrized nonlinear problem are described in a simple way with a minimum of prerequisite. Our approach is based on a detailed analysis of the error estimates and gives a global insight on the perturbation of a solution set by an approximation procedure. We notice that the solution family of the exact problem or its approximation is under mild hypotheses a differentiable manifold. An approach based on modern differential geometry is proposed in FINK and RHEINBOLDT [1983a,

1983b] or in RHEINBOLDT [1986] to study regular points. Although of great importance, we do not emphasize on the wide variety of algorithms to numerically solve the problem  $F_h(x) = 0$ , such as the predictor-corrector continuation method or the piecewise linear method, see KELLER [1977], KUEPPER, MITTELMANN, and WEBER [1984], ALLGOWER and GEORG [1990] for instance and the literature therein.



**INTRODUCTION**

Among the tools commonly used to study parameter dependent problems are the inverse and the implicit function theorems, two basic results of the differential calculus in Banach spaces. It turns out that the approach of nonlinear problems based on contraction mappings is well-suited to study numerical perturbation.

We collect in the next two sections various notations and we recall the inverse and the implicit function theorems in a form convenient for our later use. Our presentation is voluntarily brief. We refer to the literature for further details.

**1. Preliminaries**

Let  $X$  and  $Z$  be two real Banach spaces with the norms  $\|\cdot\|_X$  and  $\|\cdot\|_Z$  respectively. The set  $B(x, \delta) = \{y \in X; \|x - y\|_X < \delta\}$  is the open ball in  $X$  centered at  $x$  with radius  $\delta$ , while  $\overline{B}(x, \delta) = \{y \in X; \|y - x\|_X \leq \delta\}$  denotes the closed ball. The product space  $X \times Z$  is equipped with the norm  $\|(x, z)\|_{X \times Z} = \|x\|_X + \|z\|_Z$ .

We note  $X'$  the dual space of the real Banach space  $X$  and  $\langle \cdot, \cdot \rangle_{X', X}$  the duality pairing between  $X'$  and  $X$ . The space of continuous linear mappings from  $X$  into  $Z$  denoted by  $\mathcal{L}(X; Z)$  and endowed with the norm

$$\|A\|_{X; Z} = \sup_{\substack{x \in X \\ \|x\|_X \leq 1}} \|Ax\|_Z$$

for all  $A \in \mathcal{L}(X; Z)$ , is a Banach space. The spaces of multilinear mappings are then defined by

$$\mathcal{L}^k(X; Z) = \mathcal{L}(X; \mathcal{L}^{k-1}(X; Z)), \quad k = 2, \dots$$

Given an invertible mapping  $A \in \mathcal{L}(X; Z)$  and a mapping  $B \in \mathcal{L}(X; Z)$ , the mapping  $A + B \in \mathcal{L}(X; Z)$  is invertible if

$$(1.1) \quad \|A^{-1}B\|_{X; X} < 1.$$

Then the following bound holds

$$(1.2) \quad \|(A + B)^{-1}\|_{Z; X} \leq \frac{1}{1 - \|A^{-1}B\|_{X; X}} \|A^{-1}\|_{Z; X}.$$

Let  $\{A_n\}_{n \geq 1} \subset \mathcal{L}(X; Z)$  be a sequence of linear operators converging punctually to  $A \in \mathcal{L}(X; Z)$ , that is for all  $x \in X$ ,

$$\lim_{n \rightarrow \infty} A_n x = Ax.$$

Let  $Y$  be a Banach space and  $B \in \mathcal{L}(Y; X)$  be a compact operator. Then the sequence  $\{A_n B\}_{n \geq 1} \subset \mathcal{L}(Y; Z)$  converges to  $AB$  in the operator norm, that is

$$(1.3) \quad \lim_{n \rightarrow \infty} \|A_n B - AB\|_{Y; Z} = 0.$$

A continuous mapping  $G : X \rightarrow Z$  is Fréchet differentiable at a point  $x \in X$  if there exists an operator  $DG(x) \in \mathcal{L}(X; Z)$  such that

$$\|G(y) - G(x) - DG(x)(y - x)\|_Z = o(\|x - y\|_X).$$

The mapping  $DG : X \rightarrow \mathcal{L}(X; Z)$  is the (Fréchet) derivative of  $G$ . The mapping  $G$  is said to be of class  $C^1$  when the mapping  $DG$  is continuous on  $X$ . The mapping  $G$  is said to be of class  $C^2$  when  $DG$  is of class  $C^1$ ; the second derivative  $D^2 G(x)$  at the point  $x$  belongs to  $\mathcal{L}^2(X; Z)$ . We can repeat this procedure in order to define the class  $C^p$  with  $p \geq 2$ . Let  $G : X \rightarrow Z$  be a  $C^p$  mapping with  $p \geq 1$ . Then the Taylor expansion of  $G$  at the point  $x \in X$  reads

$$(1.4) \quad G(y) = G(x) + \sum_{k=1}^{p-1} \frac{1}{k!} D^k G(x) \underbrace{(y - x, \dots, y - x)}_{k \text{ times}} + \frac{1}{(p-1)!} \int_0^1 (1-t)^{p-1} D^p G(x + t(y-x))(y-x, \dots, y-x) dt.$$

The literature in linear and nonlinear functional analysis is rich, for further details we refer for instance to BONIC [1969], BREZIS [1983], CARTAN [1967], DIEUDONNÉ [1968], DUNFORD and SCHWARTZ [1958].

Let us finally introduce some notations useful in the applications. In the examples presented in the next chapter, the appropriate function spaces are the Sobolev spaces  $W^{m,p}(\Omega)$ , where  $\Omega$  is an open set in  $\mathbb{R}^2$  and the numbers  $m, p$  are non negative integers. These Sobolev spaces are equipped with the norms  $\|\cdot\|_{m,p,\Omega}$  and semi norms  $|\cdot|_{m,p,\Omega}$ . The space  $W_0^{m,p}(\Omega)$  is the closure of  $\mathcal{D}(\Omega)$  in  $W^{m,p}(\Omega)$  and  $W^{-m,q}(\Omega)$  is the dual space of  $W_0^{m,p}(\Omega)$  where  $\frac{1}{p} + \frac{1}{q} = 1$ . For  $p = 2$ , we shall write  $H^m(\Omega)$  and  $H_0^m(\Omega)$  instead of  $W^{m,p}(\Omega)$  and  $W_0^{m,p}(\Omega)$  respectively,  $\|\cdot\|_{m,\Omega}$  and  $|\cdot|_{m,\Omega}$  instead of  $\|\cdot\|_{m,2,\Omega}$  and  $|\cdot|_{m,2,\Omega}$ . Numerous presentations of the Sobolev spaces can be found, as for example in the books of ADAMS [1975], LIONS and MAGENES [1968], MAZ'YA [1985] or NEČAS [1967].

If  $v_1, v_2$  are two functions on  $\Omega$ ,  $v_1, v_2 : (x_1, x_2) \in \Omega \rightarrow v_1(x_1, x_2), v_2(x_1, x_2) \in \mathbb{R}$  and  $\mathbf{v} = (v_1, v_2)$ , then

$$\mathbf{grad} v_1 = \left( \frac{\partial v_1}{\partial x_1}, \frac{\partial v_1}{\partial x_2} \right), \quad \operatorname{div} \mathbf{v} = \frac{\partial v_1}{\partial x_1} + \frac{\partial v_2}{\partial x_2}$$

and

$$\mathbf{grad} \mathbf{v} = \begin{pmatrix} \frac{\partial v_1}{\partial x_1} & \frac{\partial v_1}{\partial x_2} \\ \frac{\partial v_2}{\partial x_1} & \frac{\partial v_2}{\partial x_2} \end{pmatrix}, \quad (\mathbf{v} \nabla) \mathbf{v} = \begin{pmatrix} v_1 \frac{\partial v_1}{\partial x_1} + v_2 \frac{\partial v_1}{\partial x_2} \\ v_1 \frac{\partial v_2}{\partial x_1} + v_2 \frac{\partial v_2}{\partial x_2} \end{pmatrix}.$$



## 2. Inverse and implicit function theorems

In our analysis, the inverse and implicit function theorems will play an essential role. We present here the results in a form convenient for our use. Let  $G : X \rightarrow Z$  be a  $C^1$  mapping and  $v \in X$  be such that  $DG(v) \in \mathcal{L}(X; Z)$  is an isomorphism. We introduce the notations

$$(2.1) \quad \begin{aligned} \epsilon &= \|G(v)\|_Z, \\ \gamma &= \|DG(v)^{-1}\|_{Z;X}, \\ L(\alpha) &= \sup_{x \in \overline{B}(v, \alpha)} \|DG(v) - DG(x)\|_{X;Z}, \end{aligned}$$

and we are interested in the problem to find  $x \in X$  such that

$$(2.2) \quad G(x) = 0.$$

**Theorem 2.1.** *We assume that  $DG(v) \in \mathcal{L}(X; Z)$  is an isomorphism and that  $2\gamma L(2\gamma\epsilon) \leq 1$ . Then problem (2.2) has a unique solution  $u$  in the ball  $\overline{B}(v, 2\gamma\epsilon)$  and  $DG(u) \in \mathcal{L}(X; Z)$  is invertible with*

$$(2.3) \quad \|DG(u)^{-1}\|_{Z;X} \leq 2\gamma.$$

Moreover for all  $x \in \overline{B}(v, 2\gamma\epsilon)$

$$(2.4) \quad \|x - u\|_X \leq 2\gamma \|G(x)\|_Z.$$

*Proof.* The proof follows BREZZI, RAPPAZ, and RAVIART [1980] or CROUZEIX and RAPPAZ [1989]. Let  $H : X \rightarrow X$  be the mapping defined by

$$(2.5) \quad H(x) = x - DG(v)^{-1}G(x).$$

Clearly  $x$  is a fixed point of  $H$  if and only if  $x$  is a zero of the mapping  $G$ .

For any  $x \in \overline{B}(v, 2\gamma\epsilon)$  we can write

$$H(x) - v = DG(v)^{-1} [DG(v)(x - v) - (G(x) - G(v))] - DG(v)^{-1}G(v).$$

With the Taylor expansion (1.4) with  $p = 1$  we get

$$\begin{aligned} \|H(x) - v\|_X &\leq \gamma \left[ \epsilon + \left\| \int_0^1 (DG(v) - DG(v + t(x - v)))(x - v) dt \right\|_Z \right] \\ &\leq \gamma [\epsilon + L(2\gamma\epsilon)2\gamma\epsilon] \leq 2\gamma\epsilon. \end{aligned}$$

Thus  $H$  maps the closed ball  $\overline{B}(v, 2\gamma\epsilon)$  into itself.

Let now  $x, y$  be in  $\overline{B}(v, 2\gamma\epsilon)$ , then

$$H(x) - H(y) = DG(v)^{-1} \int_0^1 [DG(v) - DG(y + t(x - y))](x - y) dt$$

and

$$\|H(x) - H(y)\|_X \leq \gamma L(2\gamma\epsilon) \|x - y\|_X \leq \frac{1}{2} \|x - y\|_X.$$

We have proved that  $H$  is a strict contraction from the ball  $\overline{B}(v, 2\gamma\epsilon)$  into itself. So  $H$  possesses a unique fixed point  $u$  in the ball  $\overline{B}(v, 2\gamma\epsilon)$ . Since

$$\|DG(v)^{-1} [DG(u) - DG(v)]\|_{X;X} \leq \gamma L(2\gamma\epsilon) \leq \frac{1}{2}$$

and

$$DG(u) = DG(v) + (DG(u) - DG(v)),$$

by (1.1) and (1.2),  $DG(u)$  is invertible and (2.3) is true, i.e.

$$\|DG(u)^{-1}\|_{Z,X} \leq 2\gamma.$$

Finally to prove the estimate (2.4), we write for  $\alpha > 0$  and  $x \in \overline{B}(v, \alpha)$

$$\begin{aligned} u - x &= H(u) - x \\ &= DG(v)^{-1} \left[ -G(x) + \int_0^1 [DG(v) - DG(u + t(x - u))](u - x) dt \right] \end{aligned}$$

and

$$\|u - x\|_X \leq \gamma [\|G(x)\|_Z + L(\alpha) \|u - x\|_X].$$

Thus we finally get

$$(1 - \gamma L(\alpha)) \|u - x\|_X \leq \gamma \|G(x)\|_Z. \quad \square$$

*Remark 2.1.* The uniqueness result in Theorem 2.1 can be improved in the following way. Under the hypotheses of the theorem and for all  $\alpha \geq 2\gamma\epsilon$  satisfying  $\gamma L(\alpha) < 1$ , problem (2.2) has a unique solution in the ball  $\overline{B}(v, \alpha)$ . Indeed if  $u_1$  and  $u_2$  are two different fixed points in  $\overline{B}(v, \alpha)$  then

$$\begin{aligned} \|u_1 - u_2\|_X &= \|H(u_1) - H(u_2)\|_X \\ &\leq \gamma L(\alpha) \|u_1 - u_2\|_X < \|u_1 - u_2\|_X \end{aligned}$$

and we have a contradiction.  $\square$

*Remark 2.2.* In Theorem 2.1 the hypothesis on the differentiability of  $G$  can be relaxed. If there exists an isomorphism  $A \in \mathcal{L}(X; Z)$  and  $v \in X$  such that  $2\gamma L(2\gamma\epsilon) \leq 1$ , with

$$\begin{aligned}\epsilon &= \|G(v)\|_Z, \\ \gamma &= \|A^{-1}\|_{Z;X}, \\ L(\alpha) &= \sup_{x,y \in \overline{B}(v,\alpha)} \frac{\|G(x) - G(y) - A(x-y)\|_Z}{\|x-y\|_X},\end{aligned}$$

then problem (2.2) has a unique solution in  $\overline{B}(v, 2\gamma\epsilon)$ . Moreover the estimate (2.4) still holds.

Note that if  $\lim_{\alpha \rightarrow 0} L(\alpha) = 0$ , then  $G$  is strongly differentiable at  $v$  and satisfies a Lipschitz condition in a neighborhood of  $v$ , see NIJENHUIS [1974] for details.  $\square$

*Remark 2.3.* Theorem 2.1 is the basic result to start an error analysis when we consider approximations of nonlinear problems. Usually  $u$  will be a solution of the approximate problem and  $x = v$  the solution of the exact problem. The estimate (2.4) is the main relation to analyze the error between exact and approximated solutions.  $\square$

The results stated in Theorem 2.1 will be sufficient for most applications. In fact under similar hypotheses we get the following variant of the inverse function theorem.

**Theorem 2.2.** *For  $v \in X$  and the function  $G : X \rightarrow Z$  of class  $C^p$ ,  $p \geq 1$ , we assume that  $DG(v) \in \mathcal{L}(X; Z)$  is an isomorphism and that  $\alpha$  satisfies  $2\gamma L(\alpha) \leq 1$ , with  $\gamma = \|DG(v)^{-1}\|_{Z;X}$ . Then there exists a  $C^p$  mapping  $F : B(G(v), \alpha/2\gamma) \rightarrow B(v, \alpha)$  such that for all  $z \in B(G(v), \alpha/2\gamma)$  we have*

$$(2.6) \quad G(F(z)) = z$$

and

$$DF(z) = [DG(F(z))]^{-1}.$$

Moreover for all  $z_1, z_2$  in  $B(G(v), \alpha/2\gamma)$

$$(2.7) \quad \|F(z_1) - F(z_2)\|_X \leq 2\gamma \|z_1 - z_2\|_Z.$$

*Proof.* We refer to BREZZI, RAPPAZ, and RAVIART [1980] for a detailed proof. We will only sketch it here. For  $z \in B(G(v), \alpha/2\gamma)$  we define the mapping  $H : X \rightarrow X$  by

$$H(x) = x + DG(v)^{-1}(z - G(x)).$$

As in the proof of Theorem 2.1, we can prove that for  $x, y \in \overline{B}(v, \alpha)$

$$\|H(x) - H(y)\|_X \leq \frac{1}{2} \|x - y\|_X$$

and for  $x \in \overline{B}(v, \alpha_1)$  with  $\alpha_1 = 2\gamma\|z - G(v)\|_Z \leq \alpha$

$$\|H(x) - v\|_X \leq \|H(x) - H(v)\|_X + \|H(v) - v\|_X \leq \alpha_1.$$

Hence the mapping  $H$  maps the closed ball  $\overline{B}(v, \alpha_1)$  into itself and is a contraction. So  $H$  has a unique fixed point in  $\overline{B}(v, \alpha_1)$ ; that is for any  $z \in B(G(v), \alpha/2\gamma)$  the equation  $G(x) = z$  has a unique solution  $x = F(z)$  in  $\overline{B}(v, \alpha_1)$ , which is unique in  $B(v, \alpha)$  (in a similar way as in Remark 2.1). Whence the mapping  $F : B(G(v), \alpha/2\gamma) \rightarrow B(v, \alpha)$  is well-defined. In a same way we proved (2.3), we get for  $x \in B(v, \alpha)$  that  $DG(x)$  is an isomorphism from  $X$  onto  $Z$  with  $\|DG(x)^{-1}\| \leq 2\gamma$ . Then in a standard way we prove that  $DF(z) = [DG(F(z))]^{-1}$  and that  $F$  is of class  $C^p$ . Finally the estimate (2.7) is a direct consequence of (2.6) and of a Taylor expansion of  $F$ .  $\square$

If we consider a  $C^1$  mapping  $G : \Lambda \times X \rightarrow Z$  where  $\Lambda$  is an other Banach space and if  $(\lambda_0, x_0)$  is an element of  $\Lambda \times X$ , then  $DG(\lambda_0, x_0) \in \mathcal{L}(\Lambda \times X; Z)$ ,  $D_\lambda G(\lambda_0, x_0) \in \mathcal{L}(\Lambda; Z)$ , and  $D_x G(\lambda_0, x_0) \in \mathcal{L}(X; Z)$  are respectively the total derivative of  $G$ , the derivative of  $G$  with respect to  $\lambda$  and the derivative of  $G$  with respect to  $x$  at the point  $(\lambda_0, x_0)$ .

Finally we recall a version of the implicit function theorem. Let  $\Lambda$ ,  $X$ , and  $Z$  be three Banach spaces and  $G : \Lambda \times X \rightarrow Z$  be a  $C^p$  mapping with  $p \geq 1$ . For a given  $(\lambda_0, x_0) \in \Lambda \times X$ , we assume that  $D_x G(\lambda_0, x_0) \in \mathcal{L}(X; Z)$  is an isomorphism and we introduce the notations

$$(2.8) \quad \begin{aligned} \epsilon &= \|G(\lambda_0, x_0)\|_Z, \\ \gamma_0 &= \|D_\lambda G(\lambda_0, x_0)\|_{\Lambda; Z}, \\ \gamma_1 &= \|D_x G(\lambda_0, x_0)^{-1}\|_{Z; X}, \\ L(\alpha) &= \sup_{(\lambda, x) \in \overline{B}((\lambda_0, x_0), \alpha)} \|DG(\lambda_0, x_0) - DG(\lambda, x)\|_{\Lambda \times X; Z}. \end{aligned}$$

**Theorem 2.3.** *Let  $\alpha$  satisfy  $2\gamma L(\alpha) \leq 1$  with  $\gamma = \max(\gamma_1, 1 + \gamma_0 \gamma_1)$ . If  $\epsilon < \alpha/4\gamma$ , then there exists a unique  $C^p$  mapping  $g : B(\lambda_0, \alpha/4\gamma) \subset \Lambda \rightarrow B(x_0, \alpha) \subset X$  satisfying for all  $\lambda \in B(\lambda_0, \alpha/4\gamma)$*

$$G(\lambda, g(\lambda)) = 0$$

and

$$(2.9) \quad \|g(\lambda) - x_0\|_X \leq 2\gamma(\epsilon + \|\lambda - \lambda_0\|_\Lambda).$$

*Proof.* We consider the mapping  $F : \Lambda \times X \rightarrow \Lambda \times Z$  given by

$$F(\lambda, x) = (\lambda, G(\lambda, x)).$$

Clearly  $F$  is of class  $C^p$  and for  $(\lambda, x) \in \Lambda \times X$  we have

$$DF(\lambda_0, x_0)(\lambda, x) = (\lambda, D_\lambda G(\lambda_0, x_0)\lambda + D_x G(\lambda_0, x_0)x).$$

Since  $D_x G(\lambda_0, x_0)$  is an isomorphism from  $X$  onto  $Z$ , so is  $DF(\lambda_0, x_0)$  from  $\Lambda \times X$  onto  $\Lambda \times Z$  and for all  $(\lambda, z) \in \Lambda \times Z$  we get

$$\|DF(\lambda_0, x_0)^{-1}(\lambda, z)\|_{\Lambda \times X} \leq \max(1 + \gamma_0 \gamma_1, \gamma_1) (\|\lambda\|_\Lambda + \|z\|_Z).$$

So we get the estimate

$$\|DF(\lambda_0, x_0)^{-1}\|_{\Lambda \times Z; \Lambda \times X} \leq \gamma.$$

Moreover for all  $(\lambda, x) \in \overline{B}((\lambda_0, x_0), \alpha)$ , we have

$$\|DF(\lambda, x) - DF(\lambda_0, x_0)\|_{\Lambda \times X; \Lambda \times Z} \leq L(\alpha).$$

We are in a position to apply Theorem 2.2 with  $G = F$  and  $v = (\lambda_0, x_0)$ . There exists a unique  $C^p$  mapping  $\mathcal{F} : B(F(\lambda_0, x_0), \alpha/2\gamma) \rightarrow B((\lambda_0, x_0), \alpha)$  such that

$$F(\mathcal{F}(\lambda, z)) = (\lambda, z).$$

Since we assume  $\epsilon < \alpha/4\gamma$ , then  $(\lambda, 0)$  is in the ball  $B((\lambda_0, G(\lambda_0, x_0)), \alpha/2\gamma)$  provided  $\|\lambda - \lambda_0\|_\Lambda \leq \alpha/4\gamma$ . Hence the mapping  $g : B(\lambda_0, \alpha/4\gamma) \rightarrow B(x_0, \alpha)$  given by

$$g(\lambda) = P\mathcal{F}(\lambda, 0)$$

with  $P$  the projection operator from  $\Lambda \times X$  onto  $X$ , is well defined on  $B(\lambda_0, \alpha/4\gamma)$ , of class  $C^p$ , and satisfies

$$G(\lambda, g(\lambda)) = 0 \quad \text{for } \lambda \in B(\lambda_0, \alpha/4\gamma).$$

Finally we can use the estimate (2.7) to get (2.9) in the following way. For  $\lambda \in B(\lambda_0, \alpha/4\gamma)$

$$\begin{aligned} \|g(\lambda) - x_0\|_X &= \|P\mathcal{F}(\lambda, 0) - P\mathcal{F}(\lambda_0, G(\lambda_0, x_0))\|_X \\ &\leq 2\gamma (\|\lambda - \lambda_0\|_\Lambda + \|G(\lambda_0, x_0)\|_Z). \quad \square \end{aligned}$$

*Remark 2.4.* The hypothesis that  $G$  is  $C^1$  can be relaxed in a same way it is done in Remark 2.2. For an operator  $A \in \mathcal{L}(\Lambda; Z)$  and an isomorphism  $B \in \mathcal{L}(X; Z)$ , we introduce the notations

$$\begin{aligned} \gamma_0 &= \|A\|_{\Lambda; Z}, \\ \gamma_1 &= \|B^{-1}\|_{Z; X}, \\ L(\alpha) &= \sup_{(\lambda, x), (\lambda^*, x^*) \in \overline{B}((\lambda_0, x_0), \alpha)} \frac{\|G(\lambda, x) - G(\lambda^*, x^*) - A(\lambda - \lambda^*) - B(x - x^*)\|_Z}{\|\lambda - \lambda^*\|_\Lambda + \|x - x^*\|_X}. \end{aligned}$$

We assume that  $\alpha$  is such that  $2\gamma L(\alpha) \leq 1$  and  $\epsilon < \alpha/4\gamma$  with  $\gamma = \max(\gamma_1, 1 + \gamma_0 \gamma_1)$ . Then Theorem 2.3 holds with  $p = 0$ . This generalization has several applications, see GIRAULT and RAVIART [1982], RAPPAZ [1984], CROUZEIX and RAPPAZ [1989], BARRETT and ELLIOTT [1989], CALOZ [1987] [1991] for instance.  $\square$



## MODEL EXAMPLES

In this chapter we present examples of nonlinear problems and introduce standard finite element approximations for all of them. It is not our purpose to present many applications where parametrized equations arise. We prefer to choose model problems and discuss thoroughly their finite element approximations. We shall consider successively the case of semilinear second order elliptic boundary value problem, the Navier-Stokes equations, and a stationary heat problem with convection.

### 3. Semilinear problems

A generic example to our theory is the case of semilinear boundary value problems. This class includes problems of the form

$$(3.1) \quad \begin{aligned} -\Delta u + g(\lambda, u) &= 0 & \text{in } \Omega, \\ u &= 0 & \text{on } \partial\Omega, \end{aligned}$$

where  $\Omega$  is a bounded open set in  $\mathbb{R}^2$  with smooth boundary  $\partial\Omega$ ,  $\Delta$  is the Laplacian operator,  $\lambda$  is a real parameter, and  $g : \mathcal{D}(g) \subset \mathbb{R} \times H_0^1(\Omega) \rightarrow L^2(\Omega)$  is a given mapping. To study problem (3.1) in a convenient functional setting, we define the mapping  $F : \mathcal{D}(g) \subset \mathbb{R} \times H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  by

$$\text{for all } (\lambda, v) \in \mathcal{D}(g), \quad F(\lambda, v) = -\Delta v + g(\lambda, v)$$

or equivalently for all  $(\lambda, v) \in \mathcal{D}(g)$  and for all  $w \in H_0^1(\Omega)$

$$(3.2) \quad \langle F(\lambda, v), w \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w \, dx + \int_{\Omega} g(\lambda, v) w \, dx.$$

Then we look for solutions  $(\lambda, u) \in \mathcal{D}(g)$  of the problem

$$(3.3) \quad F(\lambda, u) = 0.$$

The interesting point here is to study how the solutions  $u$  of (3.3) depend on  $\lambda$  and to develop numerical methods to compute them.

Given a function  $f \in L^2(\Omega)$ , we set  $w \in H_0^1(\Omega)$  to be the unique solution to

$$(3.4) \quad \text{for all } v \in H_0^1(\Omega) \quad \int_{\Omega} \mathbf{grad} w \mathbf{grad} v \, dx = \int_{\Omega} f v \, dx.$$

Thus we have defined the continuous linear operator  $T : f \in L^2(\Omega) \rightarrow Tf = w \in H_0^1(\Omega)$  which maps  $f$  onto the solution  $w$  of problem (3.4). It is useful but not always necessary for further development to assume that

$$(3.5) \quad T \in \mathcal{L}(L^2(\Omega); H_0^1(\Omega) \cap H^2(\Omega)),$$

which is a restriction on the regularity of  $\partial\Omega$ , see the standard elliptic regularity result in GILBARG and TRUDINGER [1977] or GRISVARD [1985] for instance. Then the operator  $T$  is compact from  $L^2(\Omega)$  into  $H_0^1(\Omega) \cap C^0(\bar{\Omega})$ . Introducing the operator  $T$ , problem (3.3) can be written in the equivalent way

$$(3.6) \quad u + Tg(\lambda, u) = 0.$$

With some specific choices of the function  $g$ , the solution set to problem (3.3) or (3.6) can be described precisely. Several cases will be developed later.

Given a finite-dimensional space  $V_h \subset H_0^1(\Omega)$ , a Galerkin approximation to problem (3.3) consists in finding  $(\lambda, u_h) \in (\mathbb{R} \times V_h) \cap \mathcal{D}(g)$  such that

$$(3.7) \quad \text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad} u_h \mathbf{grad} v_h \, dx + \int_{\Omega} g(\lambda, u_h) v_h \, dx = 0.$$

Usually problem (3.7) is not solved as it is but with numerical quadrature rules. We introduce a projector  $P_h$  onto  $V_h$  and consider the problem of finding  $(\lambda, u_h) \in (\mathbb{R} \times V_h) \cap \mathcal{D}(g)$  such that

$$(3.8) \quad \text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad} u_h \mathbf{grad} v_h \, dx + \int_{\Omega} P_h [g(\lambda, u_h) v_h] \, dx = 0.$$

In order to write the approximate problem (3.7) in a way similar to (3.6), we set for a given  $f \in L^2(\Omega)$ ,  $w_h \in V_h$  to be the unique solution to

$$\text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad} w_h \mathbf{grad} v_h \, dx = \int_{\Omega} f v_h \, dx.$$

So we have defined a mapping  $T_h \in \mathcal{L}(L^2(\Omega); V_h)$  and (3.7) reads

$$(3.9) \quad u_h + T_h g(\lambda, u_h) = 0.$$

*Remark 3.1.* Note that the two problems (3.6) and (3.9) have a structure identical to the general one studied by BREZZI, RAPPAZ, and RAVIART [1980, 1981a, 1981b]

$$u + Tg(\lambda, u) = 0, \quad u_h + T_h g(\lambda, u_h) = 0,$$

where  $T$  is in  $\mathcal{L}(W; V)$ ,  $W$  and  $V$  are two real Banach spaces,  $g : \mathcal{D}(g) \subset \mathbb{R} \times V \rightarrow W$  is a regular mapping,  $\{V_h\}_{0 < h \leq 1}$  is a family of finite-dimensional subspaces of  $V$  and for  $0 < h \leq 1$ ,  $T_h$  is a linear mapping in  $\mathcal{L}(W; V_h)$ . Under the main hypothesis

$$\lim_{h \rightarrow 0} \|T - T_h\|_{W; V} = 0$$



they have analyzed regular branches of solutions in BREZZI, RAPPAZ, and RAVIART [1980], limit points [1981a], and simple bifurcation points [1981b]. Bifurcations at multiple eigenvalues have been studied in RAPPAZ and RAUGEL [1982].  $\square$

Two model choices of a function  $g : \mathcal{D}(g) \subset \mathbb{R} \times H_0^1(\Omega) \rightarrow L^2(\Omega)$  are presented now. First for a given function  $f \in L^2(\Omega)$ , we consider problem (3.1) with  $g(\lambda, u) = u^3 - f$  and the problem reads : find  $u \in H_0^1(\Omega)$  such that for all  $v \in H_0^1(\Omega)$

$$(3.10) \quad \langle F(u), v \rangle_{H^{-1}(\Omega); H_0^1(\Omega)} \equiv \int_{\Omega} \mathbf{grad}u \mathbf{grad}v \, dx + \int_{\Omega} (u^3 - f)v \, dx = 0.$$

The mapping  $g$  is independent of  $\lambda$  and this example will illustrate our general approach applied to a problem not depending on a parameter. It is a classical matter to prove the existence of a solution to problem (3.10). The following result holds.

**Theorem 3.1.** *There exists a unique solution  $u \in H_0^1(\Omega)$  to the problem (3.10) and  $u$  is in  $H^2(\Omega)$ . Moreover the solution  $u$  is regular, in the sense that the derivative  $DF(u) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is an isomorphism.*

*Proof.* With the notation introduced in (3.6), our problem  $F(u) = 0$  is equivalent to

$$(3.11) \quad u + T(u^3 - f) = 0.$$

Since the mapping  $u \in H_0^1(\Omega) \rightarrow u^3 - f \in L^2(\Omega)$  is continuous and compact, the mapping  $u \in H_0^1(\Omega) \rightarrow T(u^3 - f) \in H_0^1(\Omega)$  is compact. If the real number  $\mu$  is in  $[0, 1]$  and the function  $w \in H_0^1(\Omega)$  satisfy

$$\text{for all } v \in H_0^1(\Omega) \quad \int_{\Omega} \mathbf{grad}w \mathbf{grad}v \, dx + \mu \int_{\Omega} (w^3 - f)v \, dx = 0,$$

by taking  $v = w$  we get the estimate

$$|w|_{1,\Omega}^2 + \mu \|w\|_{0,4,\Omega}^4 \leq \mu \|f\|_{0,\Omega} \|w\|_{0,\Omega}$$

which implies

$$|w|_{1,\Omega}^2 \leq \|f\|_{0,\Omega} \|w\|_{0,\Omega}.$$

With Poincaré Inequality we obtain a bound for  $|w|_{1,\Omega}$  independent of  $\mu$ . As a consequence of the Leray-Schauder homotopy theorem, see LERAY and SCHAUDER [1934] p.64, we get the existence of a solution to (3.11). The regularity assumption implies that  $u$  is in  $H^2(\Omega)$ .

Let now  $u_1$  and  $u_2$  be two solutions to (3.11), then

$$\text{for all } v \in H_0^1(\Omega) \quad \int_{\Omega} \mathbf{grad}(u_1 - u_2) \mathbf{grad}v \, dx + \int_{\Omega} (u_1^3 - u_2^3)v \, dx = 0$$

or with  $v = u_1 - u_2$

$$\int_{\Omega} |\mathbf{grad}(u_1 - u_2)|^2 \, dx + \int_{\Omega} (u_1 - u_2)^2 \left[ \frac{u_1^2 + u_2^2}{2} + \frac{(u_1 + u_2)^2}{2} \right] \, dx = 0,$$

which implies  $u_1 = u_2$ .

Let  $u \in H_0^1(\Omega)$  be the solution to (3.11). For all  $v, w \in H_0^1(\Omega)$  we have

$$\langle DF(u)v, w \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w \, dx + 3 \int_{\Omega} u^2 v w \, dx \equiv b(v, w).$$

The bilinear form  $b : H_0^1(\Omega) \times H_0^1(\Omega) \rightarrow \mathbb{R}$  is coercive on  $H_0^1(\Omega)$ . So the derivative  $DF(u) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is an isomorphism.  $\square$

We consider now finite element approximations to problem (3.10). For ease of exposition we assume that  $\Omega$  is a convex polygonal domain in  $\mathbb{R}^2$ , so it makes sense to consider  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  to be a regular family of triangulations of  $\overline{\Omega}$ , in the sense defined in CIARLET [1978], p.124 or in CIARLET [1991], p.128. The maximum diameter of the triangles  $T \in \mathcal{T}_h$  is equal to  $h$ . The approximation space and the test space are chosen to be

$$(3.12) \quad V_h = \{v \in C^0(\overline{\Omega}); v|_T \in \mathcal{P}_1(T), \text{ for all } T \in \mathcal{T}_h\} \cap H_0^1(\Omega),$$

where  $\mathcal{P}_1(T)$  represents the space of polynomials of degree  $\leq 1$  defined on  $T$ . The corresponding finite element approximation to (3.10) reads

$$(3.13) \quad \text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad} u_h \mathbf{grad} v_h \, dx + \int_{\Omega} (u_h^3 - f)v_h \, dx = 0.$$

Suppose that in the system (3.13), we would like to use numerical integration. For instance we could use the generalized trapezoidal quadrature rule which is exact for polynomials of degree up to 1 on each triangle. Then the problem reads: find  $u_h \in V_h$  such that for all  $v_h \in V_h$

$$(3.14) \quad \int_{\Omega} \mathbf{grad} u_h \mathbf{grad} v_h \, dx + \int_{\Omega} r_h [(u_h^3 - f)v_h] \, dx = 0,$$

where the  $V_h$ -interpolation operator  $r_h \in \mathcal{L}(C^0(\overline{\Omega}); V_h)$  is defined for  $v \in C^0(\overline{\Omega})$  by

$$r_h v \in V_h, \quad r_h v(x_i) = v(x_i) \quad \text{for all node } x_i \text{ of } \mathcal{T}_h.$$

Note that here the function  $f$  is supposed to be continuous in  $\overline{\Omega}$ . The main goal of the next chapter will be to develop a general formalism to study nonlinear problems. As an example we shall prove that the problem (3.13) or the (3.14) one has a unique solution  $u_h$  in some neighborhood of the solution  $u$  to problem (3.10). In particular we will show that  $u_h$  converges to  $u$  in  $H_0^1(\Omega)$  when  $h$  tends to zero and we will establish a priori and a posteriori error estimates for  $\|u - u_h\|_{1,\Omega}$ .

The second example of a function  $g : \mathcal{D}(g) \subset \mathbb{R} \times H_0^1(\Omega) \rightarrow L^2(\Omega)$  is

$$g(\lambda, u) = -\lambda f(u),$$

where  $f$  is the Nemytskii operator of some function  $\phi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ , that is for every function  $u : \Omega \rightarrow \mathbb{R}$  we have

$$\text{for all } x \in \Omega \quad f(u)(x) = \phi(x, u(x)).$$

The function  $\phi : (x, z) \in \Omega \times \mathbb{R} \rightarrow \phi(x, z) \in \mathbb{R}$  is assumed to satisfy the following conditions:

$$(3.15) \quad \left\{ \begin{array}{l} \text{the function } \phi \text{ is continuously differentiable with respect to } z \text{ uni-} \\ \text{formly in } \Omega \text{ and for all } z \in \mathbb{R}, \phi(., z) \text{ and } \frac{\partial \phi}{\partial z}(., z) \text{ are in } L^\infty(\Omega), \end{array} \right.$$

$$(3.16) \quad \phi(., 0) > 0 \quad \text{in } \Omega,$$

$$(3.17) \quad \frac{\partial \phi}{\partial z}(x, z) > 0 \quad \text{for } x \in \Omega, z \geq 0.$$

Under these regularity and monotonicity hypotheses on  $\phi$ , we shall study the problem: for  $\lambda > 0$  find  $u \in H_0^1(\Omega)$  nonnegative in  $\Omega$  solution to

$$(3.18) \quad u - \lambda T f(u) = 0.$$

The hypothesis (3.15) implies that  $f : H_0^1(\Omega) \cap C^0(\overline{\Omega}) \rightarrow L^2(\Omega)$  is well defined and  $C^1$  in a neighborhood of a solution of (3.18).

We mention that the corresponding problems have physical motivations. For instance they concern the temperature distribution in an object heated by an electric current, in which case only the positive temperature  $u$  is of interest. It is known that a limit current exists beyond which steady state solutions do not exist, see JOSEPH [1965]. The choice  $\lambda = 0$  represents the steady state problem with no current. The boundary condition is assumed to be homogeneous, only to subtract a harmonic function with a prescribed value on  $\partial\Omega$ .

To prove the existence of a value  $\lambda^*$  beyond which no solution exists, we shall use the following corollary of the maximum principle.

**Theorem 3.2.** *Let  $\rho$  be a nonnegative, not identically zero function in  $L^\infty(\Omega)$  and  $\mu_1$  be the least (or principal) eigenvalue of*

$$(3.19) \quad \begin{array}{l} -\Delta \zeta - \mu \rho \zeta = 0 \quad \text{in } \Omega, \\ \zeta = 0 \quad \text{on } \partial\Omega. \end{array}$$

*Then the eigenspace  $E_1$  corresponding to  $\mu_1$  is one-dimensional and any non zero function  $\zeta \in E_1$  has a constant sign in  $\Omega$ . Furthermore let  $\psi \in H_0^1(\Omega) \cap H^2(\Omega)$  satisfy*

$$(3.20) \quad -\Delta \psi - \lambda \rho \psi \geq 0 \quad \text{in } \Omega.$$

*Then  $\psi$  is positive in  $\Omega$  if and only if  $\lambda < \mu_1$ .*

*Proof.* The theorem will be a consequence of the classical maximum principle. A function  $v \in H_0^1(\Omega) \cap H^2(\Omega)$  with a non positive not identically zero Laplacian in  $\Omega$  is strictly positive in  $\Omega$ . The first point is a classical result.

The eigenvalue  $\mu_1$  is characterized by the Rayleigh quotient

$$\mu_1 = \min_{\substack{v \in H_0^1(\Omega) \\ v \neq 0}} \frac{|v|_{1,\Omega}^2}{\|\sqrt{\rho}v\|_{0,\Omega}^2},$$

and  $E_1 = \{v \in H_0^1(\Omega); |v|_{1,\Omega}^2 - \mu_1 \|\sqrt{\rho}v\|_{0,\Omega}^2 = 0\}$ . Now let  $\zeta$  be a non zero function in  $E_1$ . Then the function  $\eta = |\zeta|$  is in  $H_0^1(\Omega)$  and necessarily

$$|\eta|_{1,\Omega}^2 - \mu_1 \|\sqrt{\rho}\eta\|_{0,\Omega}^2 = 0,$$

which means  $\eta$  is in  $E_1$ . Since  $\zeta_+ \equiv \sup(\zeta, 0) = \frac{1}{2}(\zeta + \eta)$  and  $\zeta_- \equiv -\inf(\zeta, 0) = \frac{1}{2}(\eta - \zeta)$ , we conclude,  $E_1$  being a vector space, that  $\zeta_+$  and  $\zeta_-$  are also in  $E_1$ , that is  $\zeta_+$  and  $\zeta_-$  are in  $H_0^1(\Omega)$  and

$$-\Delta\zeta_+ = \mu_1\rho\zeta_+ \quad \text{in } \Omega, \quad -\Delta\zeta_- = \mu_1\rho\zeta_- \quad \text{in } \Omega.$$

By the maximum principle and since  $\zeta \not\equiv 0$ , either  $\zeta_+ > 0$  and then  $\zeta_- \equiv 0$ , or  $\zeta_+ \equiv 0$  and then  $\zeta_- > 0$ .

Set  $p = -\Delta\psi - \lambda\rho\psi$ . Assume first that  $\lambda < \mu_1$ . Then let  $\varphi \in H_0^1(\Omega)$  be the unique minimum of the functional

$$J(v) = \frac{1}{2} \int_{\Omega} \mathbf{grad}v \mathbf{grad}v \, dx - \frac{\lambda}{2} \int_{\Omega} \rho v^2 \, dx - \int_{\Omega} p v \, dx$$

over  $v \in H_0^1(\Omega)$ . Such a function exists, is unique, and by standard elliptic regularity  $\varphi$  belongs to  $H^2(\Omega)$ , thus  $\varphi = \psi$ . Suppose that  $\psi$  is negative somewhere. Then we define the function  $\eta = |\psi|$  which is in  $H_0^1(\Omega)$ . Since the function  $p$  is non negative, we get  $J(\eta) \leq J(\psi)$ , which contradicts the fact that  $\psi$  is the unique minimum of  $J$ . By the maximum principle and the relation (3.20),  $\psi$  cannot be zero in any open subset of  $\Omega$ .

We assume now that  $\psi$  is positive and let  $\zeta_1$  be the positive eigenvector corresponding to  $\mu_1$  with  $\|\zeta_1\|_{0,\Omega} = 1$ . Then for  $v_1, v_2 \in H_0^1(\Omega)$ , we have

$$\begin{aligned} \int_{\Omega} \mathbf{grad}\psi \mathbf{grad}v_1 \, dx - \lambda \int_{\Omega} \rho\psi v_1 \, dx &= \int_{\Omega} p v_1 \, dx, \\ \int_{\Omega} \mathbf{grad}\zeta_1 \mathbf{grad}v_2 \, dx - \mu_1 \int_{\Omega} \rho\zeta_1 v_2 \, dx &= 0, \end{aligned}$$

or setting  $v_1 = \zeta_1$ ,  $v_2 = \psi$  and subtracting

$$(-\lambda + \mu_1) \int_{\Omega} \rho\psi\zeta_1 \, dx = \int_{\Omega} p\zeta_1 \, dx,$$

so  $-\lambda + \mu_1 > 0$ .  $\square$

We first note that if  $\phi$  satisfies (3.16), (3.17) then  $\phi$  is positive for positive  $z$  and problem (3.18) can have a positive solution only for positive  $\lambda$ . The only solution of (3.18) with  $\lambda = 0$  is  $u \equiv 0$ . With the implicit function theorem 2.3, it is not difficult to check that for  $\lambda$  in a right neighborhood of 0, (3.18) has a positive solution. So let  $\lambda^* > 0$  be given by

$$\lambda^* = \sup\{\lambda \geq 0; (3.18) \text{ has a positive solution}\}.$$

**Theorem 3.3.** *We assume that  $\phi$  satisfies (3.15), (3.16), and (3.17). Then for all  $\lambda \in (0, \lambda^*)$ , problem (3.18) has a minimal positive solution  $u(\lambda)$ , that is if  $v$  is a positive solution then  $u(\lambda) \leq v$  in  $\Omega$ . Moreover  $u(\lambda)$  is an increasing function of  $\lambda$ , that is if  $0 < \lambda_1 < \lambda_2 < \lambda^*$  then*

$$0 \leq u(\lambda_1) < u(\lambda_2) \quad \text{in } \Omega.$$

*Proof.* Let  $\lambda \in (0, \lambda^*)$  be fixed. There exists  $\bar{\lambda}$ ,  $\lambda \leq \bar{\lambda} \leq \lambda^*$ , such that  $(\bar{\lambda}, \bar{u})$  is a solution of (3.18) with  $\bar{u} \geq 0$ . We define then the sequences

$$\begin{aligned} \underline{u}_0 &\equiv 0, & \underline{u}_{i+1} &= \lambda T f(\underline{u}_i), \\ \bar{u}_0 &= \bar{u}, & \bar{u}_{i+1} &= \lambda T f(\bar{u}_i). \end{aligned}$$

We prove that the sequence  $\{\underline{u}_i\}_{i \geq 0}$  is monotonous increasing, that is

$$\text{for all } x \in \Omega \quad \underline{u}_{i+1}(x) > \underline{u}_i(x).$$

With the assumption (3.16), the function  $\underline{u}_1 \in H_0^1(\Omega)$  satisfies

$$-\Delta \underline{u}_1 = \lambda \phi(\cdot, 0) > 0 \quad \text{in } \Omega,$$

and by the maximum principle  $\underline{u}_1 > 0$ . Suppose now we have  $0 < \underline{u}_1 < \underline{u}_2 < \dots < \underline{u}_i$ . This implies that the difference  $\underline{u}_{i+1} - \underline{u}_i$  is a solution to

$$\begin{aligned} -\Delta(\underline{u}_{i+1} - \underline{u}_i) &= \lambda [f(\underline{u}_i) - f(\underline{u}_{i-1})] \quad \text{in } \Omega, \\ \underline{u}_{i+1} - \underline{u}_i &= 0. \end{aligned}$$

So using (3.17) we get  $\underline{u}_i < \underline{u}_{i+1}$ . Finally with the same arguments we prove

$$0 \equiv \underline{u}_0 < \underline{u}_1 < \dots < \underline{u}_n < \dots < \bar{u}_n \dots < \bar{u}_1 < \bar{u}_0.$$

The sequence  $\{\underline{u}_i\}_{i \geq 0}$  is bounded in  $L^\infty(\Omega)$ ; since  $T$  is compact as an operator in  $\mathcal{L}(L^2(\Omega); H_0^1(\Omega) \cap C^0(\bar{\Omega}))$ , there exists a subsequence  $\{\underline{u}_{n_i}\}_{i \geq 0}$  converging to  $u$  in  $H_0^1(\Omega) \cap C^0(\bar{\Omega})$ , with  $u \leq \bar{u}$ . In fact the whole sequence converges to  $u \geq 0$ , which is a solution of (3.18). Clearly from our construction  $u$  is minimal and we set  $u = u(\lambda)$ . Moreover if  $\lambda < \bar{\lambda} < \lambda^*$ , then

$$u(\lambda)(x) < u(\bar{\lambda})(x) \quad \text{for all } x \in \Omega. \quad \square$$

The method of proof used here is known as the monotonous iteration principle.

If we add some hypotheses on  $\phi$ , then we can get more information on  $\lambda^*$ . Suppose for instance there exists a positive function  $\varphi \in L^2(\Omega)$  such that

$$(3.21) \quad \text{for all } x \in \Omega \text{ and } z > 0 \quad \phi(x, z) < \varphi(x).$$

Then we can start the monotonous iteration scheme with

$$\underline{u}_0 = 0 \quad \text{and} \quad \bar{u}_0 = \lambda T\varphi.$$

Therefore we get the existence of a positive solution for all  $\lambda > 0$  and we have  $\lambda^* = \infty$ . Suppose on the contrary that  $\phi$  satisfies

$$(3.22) \quad \phi(x, z) > \varphi(x) + \rho(x)z \quad \text{for all } x \in \Omega \text{ and } z > 0$$

with positive functions  $\varphi, \rho \in L^\infty(\Omega)$ . Then  $\lambda^* \leq \mu_1$ , where  $\mu_1$  is the least eigenvalue of (3.19). Indeed suppose that the problem (3.18) has a solution  $u > 0$  for some  $\lambda > \mu_1$ . This implies

$$-\Delta u = \lambda\phi(x, u) > \lambda\rho u \quad \text{in } \Omega.$$

We get a contradiction from Theorem 3.2. So under the hypothesis (3.22), there exists a  $\lambda^* < \infty$  beyond which problem (3.18) has no positive solution.

So far we have used essentially the monotonicity of  $\phi$ . We assume furthermore

$$(3.23) \quad \text{for all } x \in \Omega, \quad \frac{\partial \phi}{\partial z}(x, \cdot) \text{ is strictly increasing in } \mathbb{R}_+.$$

Then we can say more on  $u(\cdot)$  as a function of  $\lambda$ .

**Theorem 3.4.** *We assume that  $\phi$  satisfies (3.15), (3.16), (3.17), and (3.23). Then  $\lambda^*$  is bounded by  $\mu_1$  the principal characteristic value of*

$$\psi = \mu T Df(0)\psi$$

and the mapping  $\lambda \in [0, \lambda^*) \rightarrow u(\lambda) \in H_0^1(\Omega) \cap C^0(\bar{\Omega})$  where  $u(\lambda)$  denotes the minimal positive solution is  $C^1$  and  $(\lambda, u(\lambda))$  is a regular point of (3.18), that is  $I - \lambda T Df(u(\lambda)) \in \mathcal{L}(H_0^1(\Omega) \cap C^0(\bar{\Omega}); H_0^1(\Omega) \cap C^0(\bar{\Omega}))$  is an isomorphism.

*Proof.* Since  $\phi$  is convex in  $z$ , see (3.23), we note that

$$\phi(x, z) > \phi(x, 0) + \frac{\partial}{\partial z}\phi(x, 0)z \quad \text{for } x \in \Omega \text{ and } z > 0.$$

With the arguments developed in (3.22), we get that  $\lambda^*$  is bounded by  $\mu_1$ .

The mapping  $u : \lambda \in [0, \lambda^*) \rightarrow u(\lambda) \in H_0^1(\Omega) \cap C^0(\bar{\Omega})$  is monotone increasing, so the limits for fixed  $\lambda$ ,  $u(\lambda-)$  and  $u(\lambda+)$  both exist, are positive, and are solution of (3.18). So

$$0 \leq u(\lambda-) \leq u(\lambda) \leq u(\lambda+).$$

Thus  $u$  is left continuous by the minimality of  $u(\lambda)$ , that is  $u(\lambda-) = u(\lambda)$ .

We introduce the operator  $K(\lambda) : H_0^1(\Omega) \rightarrow H_0^1(\Omega)$  by

$$\text{for } w \in H_0^1(\Omega), \quad K(\lambda)w = \lambda T Df(u(\lambda))w,$$

and the set

$$A = \{\lambda \in [0, \lambda^*); (K(\lambda)w, w)_{1,\Omega} < (w, w)_{1,\Omega} \text{ for all } w \in H_0^1(\Omega), w \neq 0\}$$

with  $(w, v)_{1,\Omega} = \int_{\Omega} \mathbf{grad} w \mathbf{grad} v \, dx$ . Since  $Df(u(\lambda))$  is an increasing function of  $\lambda$  and

$$(K(\lambda)w, w)_{1,\Omega} = \lambda \int_{\Omega} Df(u(\lambda))w^2 \, dx,$$

the set  $A$  is an interval containing zero. For  $\lambda \in A$ ,  $I - \lambda T Df(u(\lambda))$  is an isomorphism on  $H_0^1(\Omega)$  ( or  $H_0^1(\Omega) \cap C^0(\bar{\Omega})$  ). By using the implicit function theorem, it is a simple matter to prove that  $u : A \rightarrow H_0^1(\Omega) \cap C^0(\bar{\Omega})$  is a  $C^1$  function. Suppose that  $A \neq [0, \lambda^*)$  and let  $\bar{\lambda} = \sup A$ ,  $\bar{\lambda} < \lambda^*$ . Since  $u$  is continuous on  $A$  and left continuous at  $\bar{\lambda}$ , we have  $\lim_{\lambda \rightarrow \bar{\lambda}, \lambda < \bar{\lambda}} u(\lambda) = u(\bar{\lambda})$  and  $\bar{\lambda} = \mu_1(\bar{\lambda})$  the principal characteristic value of  $T Df(u(\bar{\lambda}))$ . Since the function  $u(\cdot)$  is increasing, we have for  $\lambda > \bar{\lambda}$

$$\text{for all } x \in \Omega \quad \phi(x, u(\lambda)(x)) > \phi(x, u(\bar{\lambda})(x)) + \frac{\partial}{\partial z} \phi(x, u(\bar{\lambda})(x))(u(\lambda)(x) - u(\bar{\lambda})(x)).$$

So we obtain

$$\begin{aligned} & -\Delta(u(\lambda) - u(\bar{\lambda})) - \bar{\lambda} \frac{\partial}{\partial z} \phi(x, u(\bar{\lambda}))(u(\lambda) - u(\bar{\lambda})) \\ & = (\lambda - \bar{\lambda})f(u(\lambda)) + \bar{\lambda}(f(u(\lambda)) - f(u(\bar{\lambda}))) - \bar{\lambda} Df(u(\bar{\lambda}))(u(\lambda) - u(\bar{\lambda})) > 0. \end{aligned}$$

We know that  $u(\lambda) - u(\bar{\lambda})$  is positive and we apply Theorem 3.2 to get  $\bar{\lambda} < \mu_1(\bar{\lambda})$  which is a contradiction.  $\square$

In Theorem 3.4, we prove that the minimal positive solutions of (3.18) form a regular solution path parametrized by  $\lambda$ , see Figure 3.1.

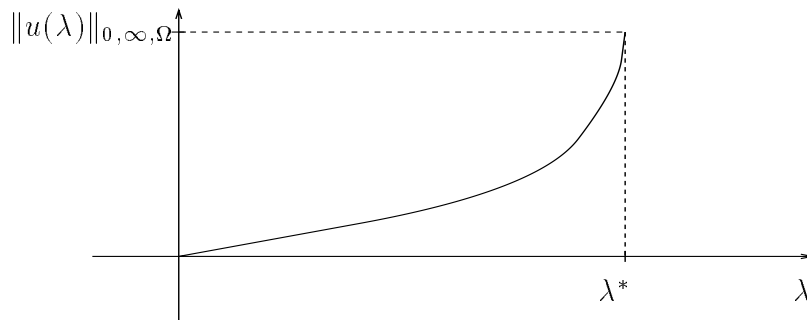


FIGURE 3.1.

The property  $\mu_1(\lambda) > \lambda$  gives a stability result of the solution branch. We shall develop a more precise study in a left neighborhood of  $\lambda^*$  when illustrating the study of simple limit points.

*Remark 3.2.* If we consider the case

$$\phi(x, z) = z^2 + 1,$$

the hypotheses (3.15), (3.16), and (3.23) are satisfied. It is a simple matter to modify Theorem 3.4 when the hypothesis (3.17) is relaxed to

$$\frac{\partial \phi}{\partial z}(x, z) > 0 \quad \text{for } x \in \Omega, z > 0.$$

We can apply the above theory and  $\lambda^*$  is bounded. Moreover one can prove that there exists  $u^* \in H_0^1(\Omega) \cap C^0(\overline{\Omega})$  satisfying (3.18) with  $\lambda = \lambda^*$ , that is

$$u^* - \lambda^* T(u^{*2} + 1) = 0$$

and

$$\lim_{\lambda \rightarrow \lambda^*} u(\lambda) = u^*;$$

furthermore  $(\lambda^*, u^*)$  is a singular point of the solution branch, in the sense that  $DF(u^*)$  is not an isomorphism, where the mapping  $F(v) = v - \lambda^* T(v^2 + 1)$ .

Suppose we can get a bound in the  $L^2(\Omega)$  norm of the right-hand side, that is

$$(3.24) \quad \|u^2(\lambda) + 1\|_{0,\Omega} \leq C \quad \text{for all } \lambda \in [0, \lambda^*),$$

this implies that the increasing function  $u^2(\lambda) + 1$  converges in  $L^2(\Omega)$  as  $\lambda$  tends to  $\lambda^*$ , so  $u(\lambda) = \lambda T(u^2(\lambda) + 1)$  converges to  $u^*$  in  $H_0^1(\Omega) \cap C^0(\overline{\Omega})$ . Clearly  $u^*$  is a solution of (3.18) with  $\lambda = \lambda^*$ , and the operator  $I - 2\lambda^* T(u^*)$  is singular. To prove (3.24) it is sufficient to prove

$$(3.25) \quad \|u(\lambda)\|_{0,4,\Omega} \leq C \quad \text{for all } \lambda \in [0, \lambda^*)$$

since  $\Omega$  is bounded.

For a given  $\lambda \in [0, \lambda^*)$ , we note  $u = u(\lambda)$ . We have

$$\text{for all } v \in H_0^1(\Omega) \quad \int_{\Omega} \mathbf{grad} u \mathbf{grad} v \, dx = \lambda \int_{\Omega} (u^2 + 1)v \, dx$$

and

$$\text{for all } w \in H_0^1(\Omega) \quad 2\lambda \int_{\Omega} uw^2 \, dx \leq \int_{\Omega} |\mathbf{grad} w|^2 \, dx.$$



We choose  $w = u^{3/2}$  and  $v = u^2$  to write successively

$$\begin{aligned} 2\lambda \int_{\Omega} u^4 dx &\leq \int_{\Omega} \frac{9}{4} u \mathbf{grad} u \mathbf{grad} u dx \\ &= \frac{9}{8} \int_{\Omega} \mathbf{grad} u^2 \mathbf{grad} u dx = \frac{9}{8} \lambda \int_{\Omega} (u^2 + 1) u^2 dx \end{aligned}$$

and so

$$(3.26) \quad \int_{\Omega} u^2 \left( u^2 - \frac{9}{7} \right) dx \leq 0.$$

Let us define the two sets

$$A = \left\{ x \in \Omega; u(x) < \frac{3\sqrt{2}}{\sqrt{7}} \right\} \quad \text{and} \quad B = \left\{ x \in \Omega; u(x) \geq \frac{3\sqrt{2}}{\sqrt{7}} \right\}.$$

Then for  $x \in B$ ,  $u^2(x) - 9/7 \geq u^2(x)/2$ , and so we have with (3.26)

$$\int_A u^2 \left( u^2 - \frac{9}{7} \right) dx + \int_B \frac{u^4}{2} dx \leq 0.$$

We are in position now to estimate

$$\begin{aligned} \int_{\Omega} u^4 dx &= \int_A u^4 dx + \int_B u^4 dx \\ &\leq \int_A \left[ -u^4 + \frac{18}{7} u^2 \right] dx \leq \int_A u^2 \left( \frac{18}{7} - u^2 \right) dx \leq 4 \left( \frac{9}{7} \right)^2 |\Omega|. \quad \square \end{aligned}$$

*Remark 3.3.* If we consider the case

$$\phi(x, z) = e^z,$$

the hypotheses (3.15), (3.16), (3.17), and (3.23) are satisfied and we can apply the above theory. Moreover one can prove that there exists  $u^* \in H_0^1(\Omega) \cap C^0(\overline{\Omega})$  satisfying (3.18), that is

$$u^* - \lambda^* T e^{u^*} = 0, \quad \lim_{\lambda \rightarrow \lambda^*} u(\lambda) = u^*;$$

furthermore  $(\lambda^*, u^*)$  is a singular point of the solution branch. The proof of this result is analogous to the one developed in Remark 3.2, see CROUZEIX and RAPPAZ [1989] for details or CRANDALL and RABINOWITZ [1973].  $\square$

*Remark 3.4.* We consider both cases when  $\phi(x, z) = z^2 + 1$  or  $\phi(x, z) = e^z$ . Then we know that the operator  $I - \lambda^* T Df(u^*)$  is singular with a kernel of dimension 1, see Theorem 3.2.

In fact since  $T$  is compact,  $I - \lambda^* T Df(u^*)$  is a self-adjoint Fredholm operator with index zero in  $H_0^1(\Omega)$  and if  $\varphi^* \in H_0^1(\Omega)$ ,  $\varphi^* \neq 0$  is such that

$$\text{Ker}(I - \lambda^* T Df(u^*)) = \text{span}\{\varphi^*\}$$

then the range of  $(I - \lambda^* T Df(u^*))$  is orthogonal to  $\text{span}\{\varphi^*\}$  in  $H_0^1(\Omega)$  with respect to the scalar product  $(\cdot, \cdot)_{1,\Omega}$  in  $H_0^1(\Omega)$ .

We define now the mapping  $\mathcal{F} : (s, \lambda, u) \in \mathbb{R} \times \mathbb{R} \times H_0^1(\Omega) \cap C^0(\overline{\Omega}) \rightarrow \mathcal{F}(s, \lambda, u) \in \mathbb{R} \times H_0^1(\Omega) \cap C^0(\overline{\Omega})$  by the relation

$$\mathcal{F}(s, \lambda, u) = \begin{pmatrix} s - (u - u^*, \varphi^*)_{1,\Omega} \\ u - \lambda T f(u) \end{pmatrix}.$$

We have  $\mathcal{F}(0, \lambda^*, u^*) = 0$  and since  $Tf(u^*)$  is a positive function and  $\varphi^*$  has a sign we can check that  $D_{\lambda u}^2 \mathcal{F}(0, \lambda^*, u^*)$  is an isomorphism on  $\mathbb{R} \times H_0^1(\Omega) \cap C^0(\overline{\Omega})$ . As a consequence of the implicit function theorem, we can parametrize  $\lambda$  and  $u$  by  $s$ . There exists a mapping  $(\lambda, u) : s \in (-\epsilon, \epsilon) \rightarrow (\lambda(s), u(s)) \in \mathbb{R} \times H_0^1(\Omega) \cap C^0(\overline{\Omega})$  satisfying

$$\lambda(0) = \lambda^*, \quad u(0) = u^*, \quad u(s) \neq u^* \quad \text{when} \quad s \neq 0,$$

and

$$u(s) - \lambda(s) T f(u(s)) = 0 \quad \text{for all } s \in (-\epsilon, \epsilon).$$

So we conclude that  $(\lambda^*, u^*)$  can be considered as a regular point of the mapping  $\mathcal{F}$ . A similar study is presented in a general setting in Section 15.  $\square$

*Remark 3.5.* There is an extensive literature concerning positive solutions of problems of the type (3.18). We mention the work of KELLER and COHEN [1967] where positive minimal solutions are studied with the monotonous iteration principle. In CRANDALL and RABINOWITZ [1973], the problem is studied by continuation and variational methods, with emphasis on the turning point  $\lambda^*$ . We can also refer to the review works of AMANN [1976], MIGNOT and PUEL [1980], and LIONS [1982].  $\square$

We consider now finite element approximations of problem (3.18). We introduce first a family of isoparametric triangulations with polynomials of degree 1 or 2, where  $h = \max_{T \in \mathcal{T}_h} h_T$  ( $h_T$  is the diameter of  $T$ ) and we denote by  $\Omega_h$  the interior of  $\overline{\Omega}_h = \cup_{T \in \mathcal{T}_h} T$ , see CIARLET [1978] p.224 or CIARLET [1991] p.230 for further details.

Actually each element  $T \in \mathcal{T}_h$  is the image  $F_T(\hat{T})$ , where  $(\hat{T}, \hat{\mathcal{P}}_k, \hat{\Sigma}_k)$  is the Lagrange reference element of type (k) ( $\hat{\mathcal{P}}_k$  is the space of polynomials of degree  $\leq k$  on  $\hat{T}$ ), and where the invertible mapping  $F_T : \hat{T} \rightarrow T$  has its components to be a polynomial of degree 1 or 2 and

$$F_T(\hat{a}_i) = a_i \quad \text{for all nodes } \hat{a}_i \in \hat{\Sigma}_k,$$

see Figure 3.2 below (for  $k = 2$ ).

The approximation space is given by

$$(3.27) \quad V_h = \{v \in C^0(\mathbb{R}^2); v(x) = 0 \text{ for all } x \notin \Omega_h, v|_T \in P_T \text{ for all } T \in \mathcal{T}_h\},$$

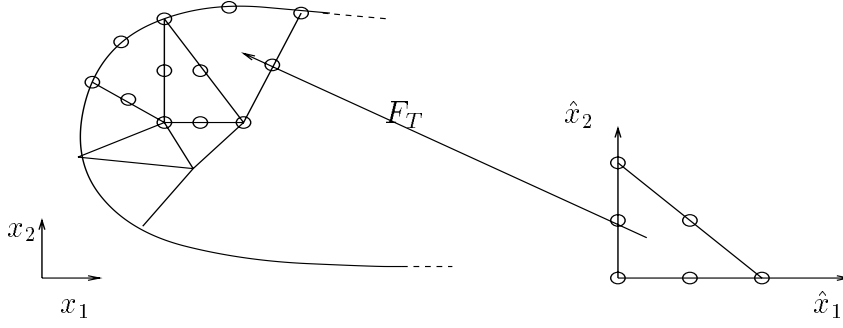


FIGURE 3.2.

with  $P_T = \{p : T \rightarrow \mathbb{R}; p \circ F_T \in \hat{\mathcal{P}}_k\}$ . The corresponding finite element approximation to (3.18) reads for all  $v_h \in V_h$

$$(3.28) \quad a_h(u_h, v_h) = \int_{\Omega \cap \Omega_h} \lambda f(u_h) v_h \, dx,$$

where

$$a_h(u_h, v_h) = \int_{\Omega_h} \mathbf{grad} u_h \mathbf{grad} v_h \, dx.$$

Generally problem (3.28) is not solved as it is, but with numerical quadrature rules, in which case we suppose that for  $v_h \in V_h$ ,  $f(v_h)$  is continuous. In the case  $k = 1$ , the elements  $F_T(\hat{T})$  are triangles and we consider the problem for all  $v_h \in V_h$

$$(3.29) \quad a_h(u_h, v_h) = \sum_{T \in \mathcal{T}_h} \left[ \frac{S_T}{3} \sum_{i=1}^3 (f(u_h) v_h)(a_{i,T}) \right]$$

where  $S_T$  is the measure of  $T$  and  $a_{1,T}$ ,  $a_{2,T}$ ,  $a_{3,T}$  the vertices of  $T \in \mathcal{T}_h$ .

To study the case  $k = 2$ , we introduce for  $u, v \in V_h$  and  $T \in \mathcal{T}_h$

$$a_T(u, v) = \sum_{i=1}^3 \frac{1}{6} \det(DF_T(\hat{m}_i)) \mathbf{grad} u(m_{i,T}) \mathbf{grad} v(m_{i,T});$$

here  $\hat{m}_i, i = 1, 2, 3$ , are the middles of the sides of  $\hat{T}$  and  $m_{i,T}$  the images  $F_T(\hat{m}_i)$ . Then problem (3.28) with numerical quadrature rules reads

$$(3.30) \quad \text{for all } v_h \in V_h \quad \sum_{T \in \mathcal{T}_h} a_T(u_h, v_h) = \sum_{T \in \mathcal{T}_h} \sum_{i=1}^3 \frac{1}{6} \det(DF_T(\hat{m}_i)) (f(u_h) v_h)(m_{i,T}).$$

As an illustration of our abstract theory developed in the next chapter, we shall analyze the problems (3.28), (3.29), and (3.30).

#### 4. A discretization of the Navier–Stokes equations

Let  $\Omega$  be a polygonal open domain in  $\mathbb{R}^2$  with boundary  $\partial\Omega$ . The stationary incompressible Navier–Stokes problem consists in finding a velocity vector  $\mathbf{u} : \Omega \rightarrow \mathbb{R}^2$  and a pressure  $p : \Omega \rightarrow \mathbb{R}$  satisfying

$$(4.1) \quad -\nu\Delta\mathbf{u} + (\mathbf{u}\nabla)\mathbf{u} + \mathbf{grad}p = \mathbf{f} \quad \text{in } \Omega,$$

$$(4.2) \quad \operatorname{div}\mathbf{u} = 0 \quad \text{in } \Omega,$$

$$(4.3) \quad \mathbf{u} = 0 \quad \text{on } \partial\Omega,$$

where the positive number  $\nu$  and the function  $\mathbf{f} \in (L^2(\Omega))^2$  are given. We introduce the spaces  $L_0^2(\Omega) = \{f \in L^2(\Omega); \int_{\Omega} f \, dx = 0\}$  and  $X = (H_0^1(\Omega))^2 \times L_0^2(\Omega)$ . A weak formulation of the problem reads: find  $(\mathbf{u}, p) \in X$  such that

$$(4.4) \quad \begin{aligned} \text{for all } \mathbf{v} \in (H_0^1(\Omega))^2 \quad & \nu \int_{\Omega} \mathbf{grad}\mathbf{u} \mathbf{grad}\mathbf{v} \, dx + \int_{\Omega} (\mathbf{u}\nabla)\mathbf{u} \mathbf{v} \, dx - \\ & \int_{\Omega} p \operatorname{div}\mathbf{v} \, dx = \int_{\Omega} \mathbf{f}\mathbf{v} \, dx, \end{aligned}$$

$$(4.5) \quad \text{for all } q \in L_0^2(\Omega) \quad \int_{\Omega} q \operatorname{div}\mathbf{u} \, dx = 0.$$

Note that if  $\mathbf{u} \in (H_0^1(\Omega))^2$ , then  $\mathbf{u} \in (L^p(\Omega))^2$  for all  $p \in [1, \infty[$  and  $(\mathbf{u}\nabla)\mathbf{u}$  is in  $(L^q(\Omega))^2$  for all  $q \in [1, 2[$ . Since  $\mathbf{v} \in (H_0^1(\Omega))^2$ , the term  $(\mathbf{u}\nabla)\mathbf{u} \mathbf{v}$  is in  $L^1(\Omega)$ .

We introduce the multilinear forms  $a : (H_0^1(\Omega))^2 \times (H_0^1(\Omega))^2 \rightarrow \mathbb{R}$  given by

$$a(\mathbf{u}, \mathbf{v}) = \nu \int_{\Omega} \mathbf{grad}\mathbf{u} \mathbf{grad}\mathbf{v} \, dx,$$

$c : (H_0^1(\Omega))^2 \times (H_0^1(\Omega))^2 \times (H_0^1(\Omega))^2 \rightarrow \mathbb{R}$  given by

$$c(\mathbf{w}; \mathbf{u}, \mathbf{v}) = \int_{\Omega} (\mathbf{w}\nabla)\mathbf{u} \mathbf{v} \, dx,$$

and  $b : (H_0^1(\Omega))^2 \times L_0^2(\Omega) \rightarrow \mathbb{R}$  given by

$$b(\mathbf{u}, p) = - \int_{\Omega} p \operatorname{div}\mathbf{u} \, dx.$$

Then the problem (4.4)–(4.5) is equivalent to find  $(\mathbf{u}, p) \in X$  such that

$$(4.6) \quad \text{for all } (\mathbf{v}, q) \in X \quad a(\mathbf{u}, \mathbf{v}) + c(\mathbf{u}; \mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) - \int_{\Omega} \mathbf{f}\mathbf{v} \, dx + b(\mathbf{u}, q) = 0.$$

The Navier–Stokes problem (4.6) has received a lot of attention and has been widely studied. It is well-known that it has at least one solution. See TEMAM [1977], GIRAULT and RAVIART [1986] and the bibliography therein for further references.

Let  $(\mathbf{u}, p) \in X$  be a solution of (4.6); we assume that  $\mathbf{u} \in (H^2(\Omega))^2$ . We also want to assume that the solution  $(\mathbf{u}, p) \in X$  is regular, in the sense that there exists  $\epsilon > 0$  such that the only solution to (4.6) in the ball  $B((\mathbf{u}, p), \epsilon)$  is  $(\mathbf{u}, p)$ . Recall that on  $X$  we introduce the norm  $\|(\mathbf{w}, r)\|_X = |\mathbf{w}|_{1,\Omega} + \|r\|_{0,\Omega}$ . Referring to Theorem 2.1 the hypothesis is stated precisely in the following terms. The bilinear form  $A : X \times X \rightarrow \mathbb{R}$  given by

$$(4.7) \quad \begin{aligned} A(\underline{w}, \underline{v}) &= A((\mathbf{w}, r), (\mathbf{v}, q)) \\ &= a(\mathbf{w}, \mathbf{v}) + b(\mathbf{w}, q) + b(\mathbf{v}, r) + c(\mathbf{u}; \mathbf{w}, \mathbf{v}) + c(\mathbf{w}; \mathbf{u}, \mathbf{v}) \end{aligned}$$

satisfies the so called inf-sup condition, see BABUŠKA and AZIZ [1972], that is

$$(4.8) \quad \begin{aligned} \inf_{\substack{\underline{w} \in X \\ \|\underline{w}\|_X=1}} \sup_{\substack{\underline{v} \in X \\ \|\underline{v}\|_X=1}} A(\underline{w}, \underline{v}) &> 0, \\ \sup_{\substack{\underline{w} \in X \\ \|\underline{w}\|_X=1}} A(\underline{w}, \underline{v}) &> 0 \quad \text{for all } \underline{v} \in X, \quad \underline{v} \neq 0. \end{aligned}$$

For simplicity we have assumed that  $\Omega$  is polygonal, so we can consider  $\mathcal{T}_h^c$  a regular triangulation of  $\bar{\Omega}$  made up of triangles  $T$  of diameter inferior to  $h$  and  $\mathcal{T}_h^f$  a finer triangulation obtained from  $\mathcal{T}_h^c$  by dividing each triangle  $T \in \mathcal{T}_h^c$  in four triangles as shown in Figure 4.1.

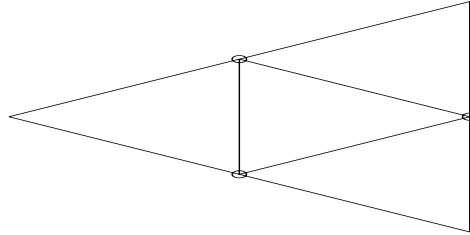


FIGURE 4.1.

The parameter  $h$  will tend to zero. We assume that the family  $\{\mathcal{T}_h^c\}_{0 < h \leq 1}$  is regular, see CIARLET [1978] p.124. To define the approximation spaces  $X_h$  to  $X$  we introduce the finite-dimensional spaces

$$V_h = \{v \in C^0(\bar{\Omega}); v|_{T'} \in \mathcal{P}_1(T') \text{ for all } T' \in \mathcal{T}_h^f, v|_{\partial\Omega} = 0\}$$

and

$$W_h = \{w \in C^0(\bar{\Omega}); w|_T \in \mathcal{P}_1(T) \text{ for all } T \in \mathcal{T}_h^c, \int_{\Omega} w \, dx = 0\},$$

where  $\mathcal{P}_1(T)$  is the space of polynomials of degree up to 1 on  $T$ . The velocity approximation space is then  $V_h \times V_h$  while the pressure approximation space is  $W_h$ , so we set  $X_h = V_h \times V_h \times W_h$ . This is a variant of the Hood-Taylor method. A finite element

approximation of (4.6), which is called the  $(4\mathcal{P}_1 - \mathcal{P}_1)$  finite element approximation, consists in finding  $(\mathbf{u}_h, p_h) \in X_h$  such that

$$(4.9) \quad \text{for all } (\mathbf{v}_h, q_h) \in X_h \quad a(\mathbf{u}_h, \mathbf{v}_h) + c(\mathbf{u}_h; \mathbf{u}_h, \mathbf{v}_h) + b(\mathbf{v}_h, p_h) - \int_{\Omega} \mathbf{f} \mathbf{v}_h \, dx + b(\mathbf{u}_h, q_h) = 0.$$

On  $\mathcal{T}_h^c$  we make the following hypothesis: there exists a set of interior nodes in  $\Omega$ ,  $P_1, \dots, P_\ell$ , such that the finite element stars  $S_1, \dots, S_\ell$

$$S_i = \bigcup_{T \ni P_i} T$$

have the property

$$(4.10) \quad \text{for all } i \neq j, \quad \overset{\circ}{S}_i \cap \overset{\circ}{S}_j = \emptyset, \quad \overline{\Omega} = \bigcup_{i=1}^{\ell} S_i,$$

where  $\overset{\circ}{S}_i$  denotes the interior of  $S_i$ . Note that we always can realize such a triangulation only to refine some elements. The study of the  $(4\mathcal{P}_1 - \mathcal{P}_1)$  finite element for the Stokes problem is well known, see GIRAULT and RAVIART [1986] p.181 for instance.

From the theory developed in the next chapter, we will be able to analyze the approximated Navier-Stokes problem (4.9) and get a priori and a posteriori error estimates. It is then essential to get the discrete analogue of the hypothesis (4.8). From the study of the  $(4\mathcal{P}_1 - \mathcal{P}_1)$  finite element for the Stokes problem, we know that under the hypothesis (4.10) we have

$$(4.11) \quad \inf_{\substack{\underline{w}_h \in X_h \\ \|\underline{w}_h\|_X = 1}} \sup_{\substack{\underline{v}_h \in X_h \\ \|\underline{v}_h\|_X = 1}} S(\underline{w}_h, \underline{v}_h) \geq \beta > 0,$$

where the bilinear form corresponding to the Stokes problem  $S : (H_0^1(\Omega))^2 \times L_0^2(\Omega) \rightarrow \mathbb{R}$  is given by

$$(4.12) \quad S(\underline{w}, \underline{v}) = S((\mathbf{w}, r), (\mathbf{v}, q)) = a(\mathbf{w}, \mathbf{v}) + b(\mathbf{w}, q) + b(\mathbf{v}, r).$$

**Theorem 4.1.** *Let  $(\mathbf{u}, p) \in X$  be a solution of (4.6) with  $\mathbf{u} \in (H^2(\Omega))^2$ , such that (4.8) is valid. Then under the assumptions (4.10) there exists  $\delta > 0$  such that*

$$(4.13) \quad \inf_{\substack{\underline{w}_h \in X_h \\ \|\underline{w}_h\|_X = 1}} \sup_{\substack{\underline{v}_h \in X_h \\ \|\underline{v}_h\|_X = 1}} A(\underline{w}_h, \underline{v}_h) \geq \delta > 0.$$

*Proof.* We define the linear operator  $\mathcal{A} : X' \rightarrow X$  by: for  $\underline{f} \in X'$ ,  $\mathcal{A}\underline{f} \in X$  is the unique solution to

$$\text{for all } \underline{v} \in X \quad A(\mathcal{A}\underline{f}, \underline{v}) = \langle \underline{f}, \underline{v} \rangle_{X'X},$$

where  $\langle \cdot, \cdot \rangle_{X'X}$  is the duality pairing between  $X'$  and  $X$ . From the hypothesis (4.8), we know that  $\mathcal{A}$  is an isomorphism from  $X'$  onto  $X$ . In the same way from the Stokes form (4.12), we define the linear operator  $\mathcal{S} : X' \rightarrow X$  by: for  $\underline{f} \in X'$ ,  $\mathcal{S}\underline{f} \in X$  is the unique solution to

$$\text{for all } \underline{v} \in X \quad S(\mathcal{S}\underline{f}, \underline{v}) = \langle \underline{f}, \underline{v} \rangle_{X'X},$$

which is an isomorphism. We also define the linear operator  $N : X \rightarrow X'$  by for all  $\underline{w} = (\mathbf{w}, r) \in X$ , for all  $\underline{v} = (\mathbf{v}, q) \in X$

$$\langle N\underline{w}, \underline{v} \rangle_{X'X} = \int_{\Omega} [(\mathbf{u}\nabla)\mathbf{w} + (\mathbf{w}\nabla)\mathbf{u}] \mathbf{v} \, dx.$$

For  $\underline{f} \in X'$  we have the relation for all  $\underline{v} \in X$

$$A(\mathcal{A}\underline{f}, \underline{v}) = S(\mathcal{A}\underline{f}, \underline{v}) + \langle N(\mathcal{A}\underline{f}), \underline{v} \rangle_{X'X} = \langle \underline{f}, \underline{v} \rangle_{X'X}$$

or

$$(4.14) \quad \langle (\mathcal{S}^{-1} + N)(\mathcal{A}\underline{f}), \underline{v} \rangle_{X'X} = \langle \underline{f}, \underline{v} \rangle_{X'X}.$$

So we have proved that

$$(4.15) \quad \mathcal{S}^{-1} + N = \mathcal{A}^{-1} \quad \text{or} \quad I + \mathcal{S}N = \mathcal{S}\mathcal{A}^{-1}$$

where  $I$  is the identity in  $X$ .

The discrete analogue of  $\mathcal{S}$  is denoted by  $\mathcal{S}_h : X' \rightarrow X_h$  and is given by: for  $\underline{f} \in X'$

$$\text{for all } \underline{v}_h \in X_h \quad S(\mathcal{S}_h\underline{f}, \underline{v}_h) = \langle \underline{f}, \underline{v}_h \rangle_{X'X}.$$

As a consequence of the hypothesis (4.10), we get

$$(4.16) \quad \lim_{h \rightarrow 0} \|\mathcal{S}_h\underline{f} - \mathcal{S}\underline{f}\|_X = 0 \quad \text{for all } \underline{f} \in X'.$$

The convergence property (4.16) and the compactness of  $N : X \rightarrow X'$  imply

$$(4.17) \quad \lim_{h \rightarrow 0} \|\mathcal{S}N - \mathcal{S}_hN\|_{X;X} = 0,$$

see (1.3).

Writing

$$I + \mathcal{S}_hN = (I + \mathcal{S}N) [I + (I + \mathcal{S}N)^{-1}(\mathcal{S}_hN - \mathcal{S}N)]$$

we deduce from (4.15), (4.17), and (1.2) that there exists an  $h_0 > 0$  such that  $I + \mathcal{S}_hN$  is an isomorphism from  $X$  onto  $X$  for  $h \leq h_0$  with a uniformly bounded inverse.

Let us prove now the discrete analogue of (4.15) or in other words

$$(4.18) \quad (I + \mathcal{S}_h N)|_{X_h} = \mathcal{S}_h \mathcal{A}^{-1}|_{X_h}.$$

We deduce from the relation (4.14) that for  $\underline{w} \in X$  and  $\underline{v}_h \in X_h$

$$S([\mathcal{S}_h \mathcal{A}^{-1} - (I + \mathcal{S}_h N)]\underline{w}, \underline{v}_h) = 0.$$

In particular with the inf-sup relation (4.11) we get

$$\text{for all } \underline{w}_h \in X_h \quad [\mathcal{S}_h \mathcal{A}^{-1} - (I + \mathcal{S}_h N)]\underline{w}_h = 0.$$

Now since  $I + \mathcal{S}_h N$  is an isomorphism from  $X$  onto  $X$  with a uniformly bounded inverse, we have

$$(4.19) \quad \text{for all } \underline{w}_h \in X_h \quad \|\mathcal{S}_h \mathcal{A}^{-1} \underline{w}_h\|_X = \|(I + \mathcal{S}_h N)\underline{w}_h\|_X \geq \gamma \|\underline{w}_h\|_X,$$

with  $\gamma > 0$ ,  $\gamma$  independent on  $h$ .

We are in a position now to prove (4.13). Let  $\underline{w}_h \in X_h$  be such that  $\|\underline{w}_h\|_X = 1$ . Then by using (4.11) and (4.19) we get

$$\begin{aligned} \sup_{\substack{\underline{v}_h \in X_h \\ \|\underline{v}_h\|_X = 1}} A(\underline{w}_h, \underline{v}_h) &= \sup_{\substack{\underline{v}_h \in X_h \\ \|\underline{v}_h\|_X = 1}} \langle \mathcal{A}^{-1} \underline{w}_h, \underline{v}_h \rangle_{X'X} = \\ &= \sup_{\substack{\underline{v}_h \in X_h \\ \|\underline{v}_h\|_X = 1}} S(\mathcal{S}_h \mathcal{A}^{-1} \underline{w}_h, \underline{v}_h) \geq \beta \|\mathcal{S}_h \mathcal{A}^{-1} \underline{w}_h\|_X \geq \beta \gamma \|\underline{w}_h\|_X = \beta \gamma, \end{aligned}$$

which proves (4.13) with  $\delta = \beta \gamma$ .  $\square$

*Remark 4.1.* We could have formulated an abstract version of the above theorem. The generalization with  $X$  a reflexive Banach space, bilinear forms  $A(.,.)$  and  $S(.,.)$ , a compact operator  $N : X \rightarrow X'$ , and  $\{X_h\}_{0 < h \leq 1}$  a family of finite-dimensional subspaces in  $X$  satisfying

$$\lim_{h \rightarrow 0} \inf_{u_h \in X_h} \|u - u_h\|_X = 0 \quad \text{for all } u \in X,$$

is straightforward.  $\square$

In the next chapter we present how to use the above result to analyze the  $(4\mathcal{P}_1 - \mathcal{P}_1)$  approximation of the Navier-Stokes problem.

## 5. A stationary heat problem with convection

The last example we develop is a convection diffusion type problem. Let  $\Omega \subset \mathbb{R}^2$  be a bounded domain with a regular boundary  $\partial\Omega$ . Given the functions  $f \in L^2(\Omega)$ ,  $k \in L^\infty(\mathbb{R})$  and the vector field  $\mathbf{c} \in (L^\infty(\Omega))^2$ , we shall consider the problem to find  $u \in H_0^1(\Omega)$  satisfying

$$(5.1) \quad -\text{div}(k(u)\mathbf{grad}u) + \mathbf{c} \mathbf{grad}u = f \quad \text{in } \Omega.$$



To study it we introduce some further assumptions on the data. The function  $k \in C^2(\mathbb{R})$  satisfies the hypotheses

$$(5.2) \quad \text{for all } s \in \mathbb{R} \quad k(s) \geq \alpha > 0,$$

$$(5.3) \quad \text{for all } s \in \mathbb{R} \quad |D^\ell k(s)| \leq \gamma_\ell, \quad \ell = 0, 1, 2,$$

with positive real numbers  $\alpha$ ,  $\gamma_0$ ,  $\gamma_1$ , and  $\gamma_2$ ; here  $D^\ell k(s)$  stands for the  $\ell$ th derivative of  $k$ . The vector field  $\mathbf{c}$  satisfies

$$(5.4) \quad \mathbf{c} \in (W^{1,\infty}(\Omega))^2 \quad \text{and} \quad \operatorname{div} \mathbf{c} = 0.$$

In the approximation schemes, we shall keep the nonlinear term in the principal part of (5.1), although to analyze the exact problem, we reduce it to a semilinear problem. Under the assumption (5.2) the primitive  $K$ ,

$$\text{for all } s \in \mathbb{R} \quad K(s) = \int_0^s k(t) dt,$$

of  $k$  is a strictly increasing function. Let  $G$  be the inverse function of  $K$  and  $g$  be the derivative of  $G$ . Then (5.1) can be written in the form

$$-\Delta K(u) + \mathbf{c} \operatorname{grad} u = f \quad \text{in } \Omega$$

or when setting  $U = K(u) \in H_0^1(\Omega)$

$$(5.5) \quad -\Delta U + g(U) \mathbf{c} \operatorname{grad} U = f \quad \text{in } \Omega,$$

and the weak form reads: find  $U \in H_0^1(\Omega)$  such that

$$(5.6) \quad \text{for all } \varphi \in H_0^1(\Omega) \quad \int_{\Omega} \operatorname{grad} U \operatorname{grad} \varphi dx + \int_{\Omega} g(U) \mathbf{c} \operatorname{grad} U \varphi dx = \int_{\Omega} f \varphi dx.$$

**Theorem 5.1.** *In addition to the hypotheses (5.2), (5.3), and (5.4) we assume that  $f \in L^\infty(\Omega)$ . Then problem (5.6) has a unique solution  $U$ ; moreover the solution  $U$  is in  $W^{2,p}(\Omega)$  for all  $1 \leq p < \infty$ .*

*Proof.* To prove the existence of a solution we introduce first the operator  $T \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$ , the inverse of the minus Laplacian operator with homogeneous boundary conditions defined in (3.4), and write (5.6) in the equivalent way

$$(5.7) \quad U - T(f - g(U) \mathbf{c} \operatorname{grad} U) = 0.$$

The mapping  $U \in H_0^1(\Omega) \rightarrow f - g(U) \mathbf{c} \operatorname{grad} U \in H^{-1}(\Omega)$  is continuous and compact. To apply the Leray-Schauder homotopy theorem, see LERAY and SCHAUDER [1934], we derive a a priori estimate for  $V \in H_0^1(\Omega)$  with

$$V - \lambda T(f - g(V) \mathbf{c} \operatorname{grad} V) = 0,$$

which is independent of  $\lambda \in [0, 1]$ . For  $\varphi \in H_0^1(\Omega)$  we have

$$\int_{\Omega} \mathbf{grad}V \mathbf{grad}\varphi \, dx - \lambda \int_{\Omega} (f - g(V)\mathbf{cgrad}V) \varphi \, dx = 0$$

or with  $\varphi = V$  and  $\mathcal{G}(s) = \int_0^s g(t)t \, dt$

$$\int_{\Omega} |\mathbf{grad}V|^2 \, dx - \lambda \int_{\Omega} fV \, dx + \lambda \int_{\Omega} \mathbf{cgrad}\mathcal{G}(V) \, dx = 0.$$

With the hypothesis (5.4) we check that  $\mathbf{cgrad}\mathcal{G}(V) = \text{div}(\mathbf{c}\mathcal{G}(V))$  and since  $\mathcal{G}(V)$  is in  $H_0^1(\Omega)$  we get

$$\int_{\Omega} |\mathbf{grad}V|^2 \, dx = \lambda \int_{\Omega} fV \, dx$$

and so

$$|V|_{1,\Omega}^2 \leq \|f\|_{0,\Omega} \|V\|_{0,\Omega} \leq C \|f\|_{0,\Omega} |V|_{1,\Omega}.$$

Thus we conclude to the existence of a solution to (5.7).

Considering the relation (5.7) and the regularity of the domain, we get from the classical elliptic regularity result that the function  $U$  is in  $H^2(\Omega)$ , from which it follows that  $f - g(U)\mathbf{cgrad}U \in L^p(\Omega)$ ,  $1 \leq p < \infty$ , and then  $U \in W^{2,p}(\Omega)$ .

Let us finally check the uniqueness result. Suppose that  $U$  and  $V$  are two solutions of (5.5). Then

$$(5.8) \quad -\Delta(U - V) + g(U)\mathbf{cgrad}U - g(V)\mathbf{cgrad}V = 0 \quad \text{in } \Omega.$$

Since  $G$  is the primitive of  $g$  with  $G(0) = 0$  and since  $\text{div } \mathbf{c} = 0$ , the relation (5.8) can be written in the form

$$(5.9) \quad -\Delta(U - V) + \text{div}(\mathbf{c}(G(U) - G(V))) = 0 \quad \text{in } \Omega$$

or with  $W = U - V$

$$(5.10) \quad -\Delta W + \text{div}\left(\mathbf{c}W \int_0^1 g(sU + (1-s)V) \, ds\right) = 0 \quad \text{in } \Omega.$$

Our goal is to prove that  $W \equiv 0$ . The weak form of (5.10) reads with  $\zeta = \int_0^1 g(sU + (1-s)V) \, ds$

$$\text{for all } \varphi \in H_0^1(\Omega) \quad \int_{\Omega} \mathbf{grad}W \mathbf{grad}\varphi \, dx - \int_{\Omega} \zeta W \mathbf{cgrad}\varphi \, dx = 0.$$

Given  $\epsilon > 0$ , we choose  $\varphi = W^+/(W^+ + \epsilon)$  where  $W^+$  stands for the positive part of  $W$ . Then we have since  $\|\zeta\|_{0,\infty} \leq 1/\alpha$ ,

$$\begin{aligned} \int_{\Omega} |\mathbf{grad} \log(1 + \frac{W^+}{\epsilon})|^2 dx &= \int_{\Omega} \frac{1}{(\epsilon + W^+)^2} |\mathbf{grad} W^+|^2 dx \\ &= \int_{\Omega} \frac{1}{\epsilon} \mathbf{grad} W \mathbf{grad} \frac{W^+}{\epsilon + W^+} dx = \int_{\Omega} \frac{1}{\epsilon} W^+ \zeta \mathbf{c} \mathbf{grad} \frac{W^+}{\epsilon + W^+} dx \\ &= \int_{\Omega} \frac{W^+}{\epsilon + W^+} \zeta \mathbf{c} \mathbf{grad} \log(1 + \frac{W^+}{\epsilon}) dx \leq \frac{\|\mathbf{c}\|_{0,\infty}}{\alpha} \int_{\Omega} |\mathbf{grad} \log(1 + \frac{W^+}{\epsilon})| dx. \end{aligned}$$

So with the Cauchy-Schwarz inequality we get

$$\int_{\Omega} |\mathbf{grad} \log(1 + \frac{W^+}{\epsilon})|^2 dx \leq C.$$

Since  $\log(1 + W^+/\epsilon)$  is in  $H_0^1(\Omega)$ , it follows from the Poincaré inequality that

$$\int_{\Omega} |\log(1 + \frac{W^+}{\epsilon})|^2 dx \leq C$$

with  $C$  independent of  $\epsilon$ . Letting  $\epsilon$  tend to 0, we get that  $W^+$  must vanish in  $\Omega$ . In the same way we prove that  $W^-$  must vanish in  $\Omega$ .  $\square$

In effective computations problem (5.1) is most often discretized as it is, but not problem (5.6). For that reason we shall focus on the approximation of (5.1). The mapping  $F : H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  is defined by: for  $v \in H_0^1(\Omega)$

$$F(v) = -\operatorname{div}(k(v)\mathbf{grad}v) + \mathbf{cgrad}v - f.$$

Let  $U$  be the unique solution of (5.6) given in Theorem 5.1. Then  $u \equiv G(U)$  is the unique zero of  $F$ ,

$$(5.11) \quad F(u) = 0,$$

and  $u \in W^{2,p}(\Omega)$ ,  $1 \leq p < \infty$ .

Generally the mapping  $F$  is not  $C^1$ , which would be a sufficient condition to apply the approximation theory developed in the next section. In fact the appropriate spaces to work with are  $W^{1,p}(\Omega)$  with  $2 < p < \infty$ , see POUSIN and RAPPAZ [1992]; the mapping  $F$  is then considered as  $F : W_0^{1,p}(\Omega) \rightarrow W^{-1,p}(\Omega)$  and has desired differentiability properties.

**Theorem 5.2.** *Under the hypotheses of Theorem 5.1, the mapping  $F : W_0^{1,p}(\Omega) \rightarrow W^{-1,p}(\Omega)$  is of class  $C^1$  when  $2 < p < \infty$  and moreover the derivative  $DF(u) \in \mathcal{L}(W_0^{1,p}(\Omega); W^{-1,p}(\Omega))$  at the solution  $u$  is an isomorphism.*

*Proof.* Let  $q$  be the conjugate of  $p$ ,  $\frac{1}{p} + \frac{1}{q} = 1$ . Then for all  $v \in W_0^{1,p}(\Omega)$  and  $\psi \in W_0^{1,q}(\Omega)$

$$\langle F(v), \psi \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} = \int_{\Omega} k(v)\mathbf{grad}v \mathbf{grad}\psi dx + \int_{\Omega} \mathbf{cgrad}v\psi dx - \int_{\Omega} f\psi dx,$$

$\langle \cdot, \cdot \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)}$  denoting the duality bracket. Then for all  $v, w$  in  $W_0^{1,p}(\Omega)$ ,  $\psi$  in  $W_0^{1,q}(\Omega)$

$$(5.12) \quad \langle DF(v)w, \psi \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} = \int_{\Omega} \mathbf{grad}(k(v)w) \mathbf{grad} \psi \, dx + \int_{\Omega} \mathbf{cgrad} w \psi \, dx$$

and with the regularity assumption on  $k$  in (5.2), (5.3), we easily verify that  $F$  is  $C^1$ .

It is known that the Laplacian operator with homogeneous boundary conditions is an isomorphism from  $W_0^{1,p}(\Omega)$  onto  $W^{-1,p}(\Omega)$ , see for instance in SIMADER [1972] or in DAUGE [1992] or the regularity result in DAUTRAY and LIONS [1987] p.538 and the characterization of  $W^{-1,p}(\Omega)$  given in GRISVARD [1985], p.17. Relation (5.12) with  $v = u$  also reads

$$(5.13) \quad TDF(u)w = k(u)w + T(\mathbf{cgrad} w)$$

where  $T$  is the inverse of the minus Laplacian operator. We introduce the mapping  $R : \varphi \in W_0^{1,p}(\Omega) \rightarrow R\varphi = \varphi/k(u) \in W_0^{1,p}(\Omega)$ , which is an isomorphism, in (5.13) with  $w = R\varphi$  to get

$$(5.14) \quad TDF(u)R\varphi = \varphi + T(\mathbf{cgrad}(R\varphi)).$$

Since the mapping  $T \in \mathcal{L}(L^p(\Omega); W_0^{1,p}(\Omega))$  is compact, so the mapping  $TDF(u)R \in \mathcal{L}(W_0^{1,p}(\Omega); W_0^{1,p}(\Omega))$  is a Fredholm operator with index 0. To prove the theorem it is then sufficient to prove that  $DF(u) \in \mathcal{L}(W_0^{1,p}(\Omega); W^{-1,p}(\Omega))$  is injective.

Let  $\varphi \in \text{Ker}(DF(u))$ , then  $\varphi \in W_0^{1,p}(\Omega)$  satisfies

$$(5.15) \quad -\Delta(k(u)\varphi) + \mathbf{cgrad}\varphi = 0 \quad \text{in } \Omega.$$

With the notations  $\psi = k(u)\varphi$ ,  $\boldsymbol{\alpha} = \mathbf{c}/k(u)$ , (5.15) takes the form

$$(5.16) \quad -\Delta\psi + \text{div}(\boldsymbol{\alpha}\psi) = 0 \quad \text{in } \Omega.$$

From the relations (5.15) and (5.16) we can conclude that  $\psi \equiv \varphi \equiv 0$  following the proof of uniqueness in Theorem 5.1.  $\square$

We introduce now a finite element approximation of (5.11). For the sake of simplicity we suppose that  $\Omega$  is a polygonal domain in  $\mathbb{R}^2$ . Given a regular family  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  of triangulations of  $\Omega$ , we associate to  $\mathcal{T}_h$  the finite element space

$$V_h = \{\varphi \in C^0(\overline{\Omega}); \varphi|_T \in \mathcal{P}_1 \text{ for all } T \in \mathcal{T}_h \text{ and } \varphi|_{\partial\Omega} = 0\},$$

where  $\mathcal{P}_1$  denotes the space of polynomials of degree up to 1. A finite element approximation of (5.11) consists in finding  $u_h \in V_h$  satisfying

$$(5.17) \quad \text{for all } v_h \in V_h \quad \int_{\Omega} k(u_h) \mathbf{grad} u_h \mathbf{grad} v_h \, dx + \int_{\Omega} \mathbf{cgrad} u_h v_h \, dx = \int_{\Omega} f v_h \, dx$$

or equivalently

$$(5.18) \quad \text{for all } v_h \in V_h \quad \langle F(u_h), v_h \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} = 0.$$

In the next chapter we will prove that (5.18) has a unique solution  $u_h \in V_h$  and get a priori and a posteriori error estimates for  $\|u - u_h\|_{1,p,\Omega}$ .

**NONLINEAR PROBLEMS  
AND THEIR NUMERICAL APPROXIMATION**

It is the purpose of this chapter to develop an approximation method for nonlinear problems not depending on a parameter. We follow here the presentation of POUSIN and RAPPAZ [1991], [1992].

A general framework is first established. Let  $X$  and  $Z$  be two real (or complex) Banach spaces,  $F : X \rightarrow Z$  be a  $C^1$  mapping, and  $\{F_h\}_{0 < h \leq 1}$  be a family of  $C^1$  mappings defined from  $X$  to  $Z$ . Under equicontinuity, consistency, and stability hypotheses we prove the convergence of a solution  $u_h$  of  $F_h(u_h) = 0$  to a solution  $u$  of  $F(u) = 0$ .

When  $Z$  is the dual space of a Banach space, a particular attention is paid to the Galerkin approximations of the problem  $F(u) = 0$ . Our approach is presented as a generalization of the linear theory. We get the analogue of the inf-sup conditions and both a priori and a posteriori error estimates.

Finally the abstract results are applied to the model problem of Section 3, to the approximation of the Navier-Stokes problem given in Section 4, and to the stationary heat problem with convection of Section 5.

**6. An abstract framework**

The two real (or complex) Banach spaces  $X$  and  $Z$  are endowed with the norms  $\|\cdot\|_X$  and  $\|\cdot\|_Z$  respectively. Given a family  $\{F_h\}_{0 < h \leq 1}$  of  $C^1$  mappings,  $F_h : X \rightarrow Z$ , we assume there exists a  $u \in X$  such that the three following hypotheses hold. The mappings  $DF_h$  are Lipschitzian at  $u$ , uniformly in a neighborhood of  $u$ , with constant  $L_h$ , that is there exists  $\eta > 0$  and for all  $h \in (0, 1]$  there exists a constant  $L_h$  such that for all  $v \in \overline{B}(u, \eta)$

$$(6.1) \quad \|DF_h(u) - DF_h(v)\|_{X;Z} \leq L_h \|u - v\|_X;$$

$$(6.2) \quad \lim_{h \rightarrow 0} (1 + L_h) \|F_h(u)\|_Z = 0;$$

$DF_h(u)$  is an isomorphism from  $X$  onto  $Z$  with an inverse uniformly bounded in  $h$ , that is there exists  $C_1$  such that

$$(6.3) \quad \text{for all } h \in (0, 1] \quad \|DF_h(u)^{-1}\|_{Z;X} \leq C_1.$$

In the examples we shall develop, we will take  $u \in X$  as a zero of a mapping  $F : X \rightarrow Z$  and  $\{F_h\}_{0 < h \leq 1}$  as an approximation family of  $F$ . Most often in the applications  $L_h$  is bounded and the assumption (6.2) can be interpreted as a consistency hypothesis while (6.3) is a stability hypothesis. Note that with our presentation we can work with  $L_h$  unbounded but with  $\|F_h(u)\|_Z$  tending to zero faster than  $L_h$  tends to  $\infty$ . Naturally the goal of this section is to prove that consistency, stability, and equicontinuity imply existence and convergence for the approximations.

**Theorem 6.1.** *Let  $\{F_h\}_{0 < h \leq 1}$  be a family of  $C^1$  mappings from  $X$  into  $Z$  and  $u \in X$  satisfy the hypotheses (6.1), (6.2), and (6.3). Then there exist  $\delta_0 > 0$  and  $h_0 > 0$  such that for all  $h \leq h_0$ , there is a unique  $u_h$  satisfying*

$$(6.4) \quad F_h(u_h) = 0 \quad \text{and} \quad \|u - u_h\|_X \leq \frac{\delta_0}{1 + L_h}.$$

Moreover for all  $h \leq h_0$

$$(6.5) \quad \|u - u_h\|_X \leq 2\|DF_h(u)^{-1}\|_{Z;X}\|F_h(u)\|_Z.$$

*Proof.* We set

$$\begin{aligned} \epsilon_h &= \|F_h(u)\|_Z, \\ \gamma_h &= \|DF_h(u)^{-1}\|_{Z;X}, \\ \tilde{L}_h(\alpha) &= \sup_{v \in \overline{B}(u, \alpha)} \|DF_h(u) - DF_h(v)\|_{X;Z}. \end{aligned}$$

From the stability (6.3), we get  $\gamma_h \leq C_1$ . If  $\eta$  is the radius of the ball  $\overline{B}(u, \eta)$  in (6.1) we set  $\delta_0 = \min(\eta, 1/(2C_1))$  and  $\delta_h = \delta_0/(1 + L_h)$ . From (6.1) we have for all  $h \in (0, 1]$

$$\tilde{L}_h(\delta_h) \leq (2C_1)^{-1}.$$

Now from the consistency (6.2), there exists  $h_0, 0 < h_0 \leq 1$  such that for all  $h \in (0, h_0]$ ,

$$2C_1\epsilon_h \leq \delta_h.$$

Thus for  $h \leq h_0$ ,

$$2\gamma_h\tilde{L}_h(2C_1\epsilon_h) \leq 2C_1\tilde{L}_h(\delta_h) \leq 1.$$

We apply Theorem 2.1 with  $G = F_h$ . So there exists a unique  $u_h \in \overline{B}(u, 2\gamma_h\epsilon_h)$  with  $F_h(u_h) = 0$ . Moreover the following estimate holds

$$\|u - u_h\|_X \leq 2\gamma_h\|F_h(u)\|_Z.$$

Applying Remark 2.1, we get the uniqueness of  $u_h$  in  $\overline{B}(u, \delta_h)$ . Finally Theorem 6.1 is proved.  $\square$

*Remark 6.1.* In Theorem 6.1, if  $L_h$  is bounded with respect to  $h$ , we have proved the existence and uniqueness of a zero of  $F_h$  in a fixed neighborhood of  $u$ .  $\square$

We consider now a  $C^1$  mapping  $F : X \rightarrow Z$ ,  $u \in X$  a zero of  $F$ , and a family  $\{F_h\}_{0 < h \leq 1}$  of mappings satisfying (6.1), (6.2), and (6.3). To relate  $u$  and  $u_h$  we have the following result.

**Theorem 6.2.** *If  $u$  is a zero of  $F$  such that  $DF(u) \in \mathcal{L}(X; Z)$  is an isomorphism from  $X$  onto  $Z$  and if the family  $\{F_h\}_{0 < h \leq 1}$  satisfies (6.1), (6.2), and (6.3), then there exists  $0 < \bar{h}_0 \leq h_0$  such that for all  $h \leq \bar{h}_0$ ,  $DF(u_h) \in \mathcal{L}(X; Z)$  is an isomorphism with a uniformly bounded inverse and*

$$(6.6) \quad \|u - u_h\|_X \leq 2\|DF(u_h)^{-1}\|_{Z;X}\|F(u_h)\|_Z,$$

where  $h_0$  and  $u_h$  are given in Theorem 6.1.

*Proof.* To apply Theorem 2.1, we set

$$\begin{aligned} \epsilon_h &= \|F(u_h)\|_Z, \\ \gamma_h &= \|DF(u_h)^{-1}\|_{Z;X}, \\ \tilde{L}_h(\alpha) &= \sup_{v \in \bar{B}(u_h, \alpha)} \|DF(v) - DF(u_h)\|_{X;Z}. \end{aligned}$$

Let  $h_0$ ,  $\delta_0$ , and  $u_h$  be given in Theorem 6.1. A consequence of the consistency (6.2) and the estimate (6.5) is

$$\lim_{h \rightarrow 0} \|u - u_h\|_X = 0.$$

Since  $DF(u)$  is an isomorphism from  $X$  onto  $Z$ , we can choose  $h_1 \leq h_0$  such that  $DF(u_h)$  is an isomorphism with

$$\gamma_h = \|DF(u_h)^{-1}\|_{Z;X} \leq 2\|DF(u)^{-1}\|_{Z;X} \equiv \gamma.$$

There exist  $0 < h_2 \leq h_1$  and  $\delta > 0$  such that

$$\tilde{L}_h(\delta) \leq \frac{1}{2\gamma} \quad \text{and} \quad \|u - u_h\|_X < \delta.$$

Now let  $h_3 \leq h_2$  be such that for all  $h \leq h_3$

$$2\gamma\epsilon_h \leq \delta.$$

We are in position to apply Theorem 2.1 with  $G = F$ ,  $v = u_h$ ,  $0 < h \leq h_3$ . We conclude that there exists a unique  $w \in \bar{B}(u_h, 2\gamma_h\epsilon_h)$  such that

$$F(w) = 0 \quad \text{and} \quad \|w - u_h\|_X \leq 2\|DF(u_h)^{-1}\|_{Z;X}\|F(u_h)\|_Z.$$

Applying Remark 2.1, we get the uniqueness of  $w$  in  $\bar{B}(u_h, \delta)$ . Since  $F(u) = 0$ , we have  $w = u$  (actually we must choose  $\bar{h}_0 < h_3$ ) and (6.6) is proved.  $\square$

The estimates (6.5) and (6.6) have a different meaning. From the first estimate we can deduce a a priori estimate bounded by the consistency error  $\|F_h(u)\|_Z$ . On the other hand we get from the second estimate a a posteriori estimate bounded by the residual  $\|F(u_h)\|_Z$ .

**Theorem 6.3.** *Let  $F : X \rightarrow Z$  be a  $C^1$  mapping and  $\{F_h\}_{0 < h \leq 1}$ ,  $F_h : X \rightarrow Z$ , be a family of  $C^2$  mappings satisfying*

$$(6.7) \quad \begin{cases} \lim_{h \rightarrow 0} \|F(v) - F_h(v)\|_Z = 0 \text{ for all } v \in X, \\ \lim_{h \rightarrow 0} \|DF(v) - DF_h(v)\|_{X;Z} = 0 \text{ for all } v \in X, \\ D^2F_h \text{ is uniformly bounded with respect to } h \text{ on all bounded domain of } X. \end{cases}$$

Let  $u \in X$  be such that

$$(6.8) \quad F(u) = 0,$$

$$(6.9) \quad DF(u) \text{ is an isomorphism from } X \text{ onto } Z.$$

Then there exist  $h_0 > 0$ ,  $\delta_0 > 0$  such that for all  $h \leq h_0$  there exists a unique  $u_h \in X$  satisfying

$$F_h(u_h) = 0 \quad \text{and} \quad \|u - u_h\|_X \leq \delta_0.$$

Moreover for all  $h \leq h_0$  we have

$$(6.10) \quad \|u - u_h\|_X \leq 2\|DF_h(u)^{-1}\|_{Z;X}\|F_h(u)\|_Z,$$

$$(6.11) \quad \|u - u_h\|_X \leq 2\|DF(u_h)^{-1}\|_{Z;X}\|F(u_h)\|_Z,$$

*Proof.* It is easy to prove that the hypotheses (6.7), (6.8), and (6.9) imply the ones of Theorems 6.1 and 6.2 with the constants  $L_h$  bounded with respect to  $h$ .  $\square$

*Remark 6.2.* For simplicity, the function  $u$  in Theorem 6.1 is not supposed to depend on  $h$ . In Section 12, we assume a dependence on  $h$ . A convenient choice of a function depending on  $h$  may lead to error estimates in different norms, see Section 14.  $\square$

In the next result we present a simple way to derive error estimates in different norms, see RAPPAZ [1983] for a variant of this technique.

We assume there exist two Banach spaces  $\mathcal{X}$  and  $\mathcal{Z}$ , with the norms  $\|\cdot\|_{\mathcal{X}}$  and  $\|\cdot\|_{\mathcal{Z}}$ , satisfying the following assumptions.

$$(6.12) \quad X \subset \mathcal{X} \text{ and } Z \subset \mathcal{Z} \text{ with continuous injection;}$$

there exist  $\epsilon > 0$  and  $\kappa > 0$  such that for all  $v \in B_X(u, \epsilon) = \{w \in X; \|w - u\|_X < \epsilon\}$ , the operators  $DF_h(v)$  and  $DF(v)$  can be extended as operators in  $\mathcal{L}(\mathcal{X}; \mathcal{Z})$  and

$$(6.13) \quad \begin{cases} \text{for all } v \in B_X(u, \epsilon) & \|DF(v) - DF(u)\|_{\mathcal{X}; \mathcal{Z}} \leq \kappa\|v - u\|_X, \\ \text{for all } v \in B_X(u, \epsilon) & \|DF_h(v) - DF_h(u)\|_{\mathcal{X}; \mathcal{Z}} \leq \kappa\|v - u\|_X; \end{cases}$$

$$(6.14) \quad \lim_{h \rightarrow 0} \|DF(u) - DF_h(u)\|_{\mathcal{X}; \mathcal{Z}} = 0;$$

$$(6.15) \quad DF(u) \text{ is an isomorphism from } \mathcal{X} \text{ onto } \mathcal{Z}.$$

The following result holds.



**Theorem 6.4.** *Under the assumptions of Theorem 6.3 and the (6.12), (6.13), (6.14), (6.15) ones, there exists  $0 < h_1 \leq h_0$  such that for all  $0 < h \leq h_1$ , we have*

$$(6.16) \quad \|u - u_h\|_{\mathcal{X}} \leq 2\|DF_h(u)^{-1}\|_{\mathcal{X}; \mathcal{X}} \|F_h(u)\|_{\mathcal{X}},$$

$$(6.17) \quad \|u - u_h\|_{\mathcal{X}} \leq 2\|DF(u_h)^{-1}\|_{\mathcal{X}; \mathcal{X}} \|F(u_h)\|_{\mathcal{X}}.$$

*Proof.* We first prove the estimate (6.16). Noticing that  $F_h(u_h) = 0$ , we use the Taylor expansion to write

$$(6.18) \quad u - u_h = DF_h(u)^{-1} \left[ \int_0^1 (DF_h(u) - DF_h(su + (1-s)u_h))(u - u_h) ds \right] \\ + DF_h(u)^{-1} F_h(u).$$

Clearly the assumptions (6.14), (6.15), the result (1.1), and the estimate (6.10) imply the existence of  $0 < \bar{h} \leq h_0$  and  $\beta > 0$  such that for all  $0 < h \leq \bar{h}$ , the mapping  $DF_h(u)$  is an isomorphism from  $\mathcal{X}$  onto  $\mathcal{Z}$  with

$$\|DF_h(u)^{-1}\|_{\mathcal{X}; \mathcal{X}} \leq \beta$$

and

$$\|u - u_h\|_X < \epsilon.$$

With (6.18), we get

$$(6.19) \quad \|u - u_h\|_{\mathcal{X}} \leq \|DF_h(u)^{-1}\|_{\mathcal{X}; \mathcal{X}} \{ \kappa \|u - u_h\|_X \|u - u_h\|_{\mathcal{X}} + \|F_h(u)\|_{\mathcal{X}} \}.$$

Let us choose now  $0 < h_1 \leq \bar{h}$  such that for all  $0 < h \leq h_1$  the following holds

$$\beta \kappa \|u - u_h\|_X \leq \frac{1}{2}.$$

Then from the estimate (6.19) we deduce the (6.16) one.

To prove the estimate (6.17), we proceed as before with the relation

$$u - u_h = DF(u)^{-1} \left[ \int_0^1 (DF(u) - DF(su + (1-s)u_h))(u - u_h) ds \right] \\ - DF(u)^{-1} F(u_h). \quad \square$$

### 7. Galerkin approximation of nonlinear problems

In order to relate the results of the previous section with the ones relative to the Galerkin approximation of linear problems, we shall focus on a particular abstract framework. We suppose from now on that we have  $Z = Y'$ ,  $Y'$  being the dual space of a real Banach space  $Y$ . We assume

$$(7.1) \quad \begin{cases} F : X \rightarrow Y' \text{ is a } C^1 \text{ mapping and } u \in X \text{ is such that } F(u) = 0 \\ \text{and } DF(u) \text{ is an isomorphism from } X \text{ onto } Y'. \end{cases}$$

Note that

$$F(u) = 0 \quad \iff \quad \text{for all } y \in Y \quad \langle F(u), y \rangle_{Y'Y} = 0.$$

Let now  $\{X_h\}_{0 < h \leq 1}$  be a family of finite-dimensional subspaces of  $X$  and  $\{Y_h\}_{0 < h \leq 1}$  be a family of finite-dimensional subspaces of  $Y$ . A Galerkin approximation of the problem to find  $u \in X$  such that

$$(7.2) \quad F(u) = 0$$

reads: find  $u_h \in X_h$  such that

$$(7.3) \quad \text{for all } y_h \in Y_h \quad \langle F(u_h), y_h \rangle_{Y'Y} = 0.$$

To introduce the analogue of the inf-sup conditions for the nonlinear case, we define the bilinear form  $b : X \times Y \rightarrow \mathbb{R}$  by

$$\text{for all } x \in X, y \in Y \quad b(x, y) = \langle DF(u)x, y \rangle_{Y'Y}.$$

We assume there exists a constant  $\beta_h > 0$  such that

$$(7.4) \quad \inf_{\substack{x \in X_h \\ \|x\|_X = 1}} \sup_{\substack{y \in Y_h \\ \|y\|_Y = 1}} b(x, y) \geq \beta_h$$

and

$$(7.5) \quad \dim X_h = \dim Y_h$$

for  $0 < h \leq 1$ . Finally we assume

$$(7.6) \quad \lim_{h \rightarrow 0} \inf_{x_h \in X_h} \beta_h^{-2} \|u - x_h\|_X = 0.$$

Remark that  $\beta_h \leq \|DF(u)\|_{X;Y'}$  and if  $\beta_h$  is bounded from below by  $\beta > 0$  for all  $h$ , then (7.4) is a standard stability condition for linear problems and (7.6) is a standard approximation property of  $X$  by  $X_h$ .

We are in a position to state the major result of the section.

**Theorem 7.1.** *Let us suppose that the hypotheses (7.1), (7.4), (7.5), and (7.6) are satisfied. Moreover we assume that  $DF$  is Lipschitzian at  $u$ , that is there exist  $\epsilon_0 > 0$  and  $L > 0$  such that for all  $v \in X$  with  $\|u - v\|_X \leq \epsilon_0$*

$$\|DF(u) - DF(v)\|_{X;Y'} \leq L\|u - v\|_X.$$

*Then there exist two constants  $h_0 > 0$ ,  $\delta_0 > 0$ , and for  $h \leq h_0$  a unique solution  $u_h$  to problem (7.3) in the closed ball  $\overline{B}(u, \delta_h)$  with  $\delta_h = \delta_0 \beta_h$ . Moreover there exists a constant  $C$  independent of  $h$  such that*

$$(7.7) \quad \|u - u_h\|_X \leq C\beta_h^{-1} \inf_{x_h \in X_h} \|u - x_h\|_X$$

and

$$(7.8) \quad \|u - u_h\|_X \leq 2\|DF(u_h)^{-1}\|_{Y';X}\|F(u_h)\|_{Y'}.$$

*Proof.* The proof is based on Theorems 6.1 and 6.2. To define the mapping  $F_h$ , we introduce the two projectors  $\Pi_{X_h} \in \mathcal{L}(X; X_h)$  and  $\Pi_{Y_h} \in \mathcal{L}(Y; Y_h)$  by setting

$$(7.9) \quad \text{for all } y \in Y_h \quad b(x - \Pi_{X_h}x, y) = 0,$$

$$(7.10) \quad \text{for all } x \in X_h \quad b(x, y - \Pi_{Y_h}y) = 0.$$

Under the hypotheses (7.4) and (7.5), both operators are well defined. Therefore we can construct a family  $\{F_h\}_{0 < h \leq 1}$  of mappings,  $F_h : X \rightarrow Y'$ , by

$$(7.11) \quad \text{for all } x \in X, y \in Y \quad \langle F_h(x), y \rangle_{Y'Y} = \langle F(x), \Pi_{Y_h}y \rangle_{Y'Y} + b(x, y - \Pi_{Y_h}y).$$

Then we can analyze the problem to find  $x \in X$  such that

$$(7.12) \quad F_h(x) = 0.$$

Problems (7.3) and (7.12) are equivalent. Clearly if  $u_h \in X_h$  is a solution to (7.3) then it is a solution to (7.12). Conversely let  $u_h \in X$  be a solution to (7.12). To prove that  $u_h$  is a solution to (7.3) it is sufficient to prove that  $u_h \in X_h$ . By taking  $y = y_1 - \Pi_{Y_h}y_1$  with  $y_1 \in Y$  in (7.11) we get

$$b(u_h, y_1 - \Pi_{Y_h}y_1) = 0$$

and so

$$\text{for all } y_1 \in Y \quad b(u_h - \Pi_{X_h}u_h, y_1) = 0.$$

Since  $DF(u)$  is an isomorphism, we conclude that

$$u_h = \Pi_{X_h}u_h \in X_h.$$

In order to conclude, we will apply Theorems 6.1 and 6.2 to the family  $\{F_h\}_{0 < h \leq 1}$  defined in (7.11). So let us check that the hypotheses (6.1), (6.2), and (6.3) are fulfilled.

The Fréchet derivative of  $F_h$  at  $v \in X$  is,

$$\text{for all } x \in X, y \in Y \quad \langle DF_h(v)x, y \rangle_{Y'Y} = \langle DF(v)x, \Pi_{Y_h} y \rangle_{Y'Y} + b(x, y - \Pi_{Y_h} y).$$

It follows that for  $x \in X, y \in Y$

$$\langle (DF_h(u) - DF_h(v))x, y \rangle_{Y'Y} = \langle (DF(u) - DF(v))x, \Pi_{Y_h} y \rangle_{Y'Y}$$

and consequently for  $v \in X$  with  $\|u - v\|_X \leq \epsilon_0$  we get

$$(7.13) \quad \|DF_h(u) - DF_h(v)\|_{X;Y'} \leq L\|u - v\|_X \|\Pi_{Y_h}\|_{Y;Y}.$$

With the notation  $\|b\| = \|DF(u)\|_{X;Y'}$ , we easily check with (7.4) that

$$(7.14) \quad \|\Pi_{X_h}\|_{X;X} \leq \frac{\|b\|}{\beta_h}.$$

In order to bound the norm of  $\Pi_{Y_h}$  we write for  $y \in Y$

$$\begin{aligned} \|\Pi_{Y_h} y\|_Y &= \sup_{\substack{\varphi \in Y' \\ \|\varphi\|_{Y'}=1}} \langle \varphi, \Pi_{Y_h} y \rangle_{Y'Y} = \sup_{\substack{\varphi \in Y' \\ \|\varphi\|_{Y'}=1}} b(DF(u)^{-1}\varphi, \Pi_{Y_h} y) \\ &= \sup_{\substack{\varphi \in Y' \\ \|\varphi\|_{Y'}=1}} b(\Pi_{X_h} DF(u)^{-1}\varphi, y). \end{aligned}$$

By using (7.14) we finally obtain

$$(7.15) \quad \|\Pi_{Y_h}\|_{Y;Y} \leq \frac{\|b\|^2}{\beta_h} \|DF(u)^{-1}\|_{Y';X}.$$

We notice that  $\beta_h \leq \|b\|$  and  $\|DF(u)^{-1}\|_{Y';X}^{-1} \leq \|b\|$ , so the inequalities (7.13) and (7.15) imply that the hypothesis (6.1) holds with

$$L_h = (1 + L) \frac{\|b\|^2}{\beta_h} \|DF(u)^{-1}\|_{Y';X} - 1.$$

Let us check the hypothesis (6.2). By the definition of  $F_h$  we have

$$(7.16) \quad \begin{aligned} \|F_h(u)\|_{Y'} &= \sup_{\substack{y \in Y \\ \|y\|_Y=1}} \langle F_h(u), y \rangle_{Y'Y} = \sup_{\substack{y \in Y \\ \|y\|_Y=1}} b(u, y - \Pi_{Y_h} y) \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y=1}} b(u - \Pi_{X_h} u, y) \leq \|b\| \|u - \Pi_{X_h} u\|_X. \end{aligned}$$

Using the assumption (7.4), we get for  $x_h \in X_h$

$$\begin{aligned} \beta_h \|x_h - \Pi_{X_h} u\|_X &\leq \sup_{\substack{y \in Y_h \\ \|y\|_Y=1}} b(x_h - \Pi_{X_h} u, y) \\ &\leq \sup_{\substack{y \in Y_h \\ \|y\|_Y=1}} b(x_h - u, y) \leq \|b\| \|u - x_h\|_X. \end{aligned}$$

So we deduce that for  $x_h \in X_h$

$$\|u - \Pi_{X_h} u\|_X \leq \left(1 + \frac{\|b\|}{\beta_h}\right) \|u - x_h\|_X$$

and using (7.16) we have

$$(7.17) \quad \|F_h(u)\|_{Y'} \leq \frac{C}{\beta_h} \inf_{x_h \in X_h} \|u - x_h\|_X$$

where  $C$  is independent of  $h$ .

When taking into account that

$$L_h = (1 + L) \|b\|^2 \|DF(u)^{-1}\|_{Y'X} / \beta_h - 1$$

and using the hypothesis (7.6), we prove that the hypothesis (6.2) is fulfilled. Going back to the definition of  $F_h$  we have  $DF_h(u) = DF(u)$ , which is an isomorphism from  $X$  onto  $Y'$  with a uniformly bounded inverse. Then we conclude that hypothesis (6.3) is satisfied.

We conclude with Theorems 6.1 and 6.2. The error estimates (7.7) is a direct consequence of (6.5), (7.17), and of the relation  $DF_h(u) = DF(u)$ .  $\square$

*Remark 7.1.* Theorem 7.1 can be viewed as a generalization of the linear theory. If the mapping  $F$  is smooth enough and the solution regular in the sense of (7.1), then we have reduced the approximation study to a verification of an inf-sup condition, see (7.4), and of an approximability result, see (7.6). Then the error estimates are deduced via an interpolation result for the a priori estimate (7.7) and via a computation of the residual for the a posteriori estimate (7.8).  $\square$

*Remark 7.2.* The approach proposed by RHEINBOLDT [1986] to study the approximation of parameter dependent problems with projection methods is in some respect close to ours.  $\square$

*Remark 7.3.* Quite often in the applications the spaces  $X$  and  $Y$  are identical Hilbert spaces and there exists a Hilbert space  $W$  satisfying  $X \subset W$  with continuous embedding and  $X$  dense in  $W$ . The space  $W$  takes the place of a pivot space in the embeddings

$$X \subset W \subset X'.$$

The scalar product in  $X$  (respectively  $W$ ) is denoted by  $(\cdot, \cdot)_X$  (respectively  $(\cdot, \cdot)_W$ ).

We assume that  $X_h = Y_h$ , the hypotheses of Theorem 7.1 hold, and the mapping  $F$  is of class  $C^2$ . Then we can derive an estimate in the norm of  $W$  from the estimate in the norm of  $X$  using a standard duality argument.

Indeed we have

$$\|u - u_h\|_W = \sup_{\substack{g \in W \\ g \neq 0}} \frac{|(g, u - u_h)_W|}{\|g\|_W}.$$

For a given  $g$  in  $W$ , let  $w$  be the unique solution to

$$\text{for all } v \in X \quad b(v, w) = (v, g)_W.$$

Then there exists a constant  $C$  independent of  $g$  such that

$$\|w\|_X \leq C\|g\|_W.$$

So for a given  $g$  in  $W$  we get

$$|(g, u - u_h)_W| = |b(u - u_h, w)| \leq |b(u - u_h, w - w_h)| + |b(u - u_h, w_h)|$$

for any function  $w_h \in X_h$ . We can estimate the term  $|b(u - u_h, w_h)|$  in the following way

$$\begin{aligned} b(u - u_h, w_h) &= \langle DF(u)(u - u_h), w_h \rangle_{X'X} = \langle -F(u) + F(u_h) - DF(u)(u_h - u), w_h \rangle_{X'X} \\ &= \left\langle \int_0^1 (1-t) D^2 F(u + t(u_h - u))(u_h - u)^2, w_h \right\rangle_{X'X} \end{aligned}$$

and so

$$|b(u - u_h, w_h)| \leq C\|u - u_h\|_X^2 \|w_h\|_X.$$

For any  $g$  in  $W$  and  $w_h \in X_h$ , we get

$$|(u - u_h, g)_W| \leq C \left\{ \|u - u_h\|_X \|w - w_h\|_X + \|u - u_h\|_X^2 \|w_h\| \right\}.$$

We have proved that for any  $g \in W$  and  $w_h \in X_h$

$$|(u - u_h, g)_W| \leq C \left\{ \|u - u_h\|_X [1 + \|u - u_h\|_X] \|w - w_h\|_X + \|u - u_h\|_X^2 \|g\|_W \right\}.$$

An illustration of the method is given in Theorem 8.2.  $\square$

*Remark 7.4.* When the mapping  $F : X \rightarrow Y'$  has the following structure

$$F(x) = Lx + G(x)$$

where  $L \in \mathcal{L}(X; Y')$  is an isomorphism and  $G : X \rightarrow Y'$  is such that  $DG(x) \in \mathcal{L}(X; Y')$  is compact, then we can use the simple bilinear form  $b : X \times Y \rightarrow \mathbb{R}$  given by

$$b(x, y) = \langle Lx, y \rangle_{Y'Y}.$$

In Section 16, this idea is used to simplify the study of the approximation of simple limit points.  $\square$

### 8. Application to a model problem

The general results of last section are applied to the semilinear problem (3.10). We shall study both approximation schemes (3.13) and (3.14).

In this particular case, the notations are the following:  $X = Y = H_0^1(\Omega)$  equipped with the norm  $\|v\|_X = \|v\|_Y = |v|_{1,\Omega} = (\int_{\Omega} |\mathbf{grad} v|^2 dx)^{1/2}$  equivalent to  $\|\cdot\|_{1,\Omega}$ ,  $Z = Y' = H^{-1}(\Omega)$ , and  $F : X \rightarrow Y'$  is given by

$$(8.1) \quad \text{for all } u, v \in H_0^1(\Omega) \quad \langle F(u), v \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} = \int_{\Omega} \mathbf{grad} u \mathbf{grad} v dx + \int_{\Omega} u^3 v dx - \int_{\Omega} f v dx,$$

where  $f$  is an  $L^2(\Omega)$  function. Theorem 3.1 implies there exists a unique solution  $u \in H_0^1(\Omega)$  to the problem

$$(8.2) \quad F(u) = 0;$$

furthermore  $u$  is in  $H^2(\Omega)$  and the derivative  $DF(u) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is an isomorphism, which means that (7.1) is satisfied; moreover  $DF$  is Lipschitzian at  $u$ .

The approximation space  $X_h$  and the test space  $Y_h$  are both  $V_h$  given in (3.12), the space of continuous piecewise polynomials of degree  $\leq 1$ , so the hypothesis (7.5) is guaranteed. The bilinear form  $b : H_0^1(\Omega) \times H_0^1(\Omega) \rightarrow \mathbb{R}$ ,

$$b(v, w) = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w dx + 3 \int_{\Omega} u^2 v w dx,$$

is coercive on  $H_0^1(\Omega)$  and this property ensures the hypothesis (7.4) with  $\beta_h = 1$ . Since the family of triangulations  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  is regular and  $\beta_h$  is equal to 1, the assumption (7.6) is satisfied.

We can apply Theorem 7.1 to study the existence, the local uniqueness and the convergence of problem (3.13), that is find  $u_h \in V_h$  such that

$$(8.3) \quad \text{for all } v_h \in V_h \quad \langle F(u_h), v_h \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} = 0.$$

A definition useful in the next theorem is the following. Let  $w$  be a function in  $W^{1,1}(\Omega)$  (so that the trace operator can be defined),  $T_i, T_j$ , with  $i < j$ , be two adjacent triangles, and let  $T'$  denote the common side. If  $w_i$  and  $w_j$  denote the restrictions of  $w$  on  $T_i$  and  $T_j$  respectively, we set

$$[w] = w_i - w_j \quad \text{on } T',$$

$[w]$  is the jump of  $w$  through the edge  $T'$ . If  $T'$  is a side of  $T_i$  on the boundary of  $\Omega$ , then  $[w] = w_i$  on  $T'$ .

**Theorem 8.1.** *There exist positive constants  $h_0, \delta_0, C$ , and for  $h \leq h_0$  a unique solution  $u_h$  to (8.3) in  $\overline{B}(u, \delta_0)$ . Moreover the following a priori and a posteriori estimates hold for  $h \leq h_0$  :*

$$(8.4) \quad |u - u_h|_{1,\Omega} \leq Ch \|f\|_{0,\Omega}$$

and

$$(8.5) \quad |u - u_h|_{1,\Omega} \leq C \left( \sum_{T \in \mathcal{T}_h} \eta(T)^2 \right)^{1/2},$$

where  $\eta(T)$  is the local estimator given by

$$\eta(T) = \left( h_T^2 \| -\Delta u_h + u_h^3 - f \|_{0,T}^2 + \sum_{i=1}^3 h_{t_i} \left\| \left[ \frac{\partial u_h}{\partial n} \right] \right\|_{0,t_i}^2 \right)^{1/2}$$

with  $h_T$  the diameter of  $T \in \mathcal{T}_h$ ,  $t_i$   $i = 1, 2, 3$  the edges of  $T$  of length  $h_{t_i}$ , and  $[\cdot]$  the jump through an edge.

*Proof.* The hypotheses of Theorem 7.1 were verified before with  $\beta_h = 1$ , so there exist positive constants  $h_0, \delta_0$ , and for  $h \leq h_0$  a unique solution  $u_h$  to (8.3) in  $\overline{B}(u, \delta_0)$ . From the estimates (7.7) and (7.8), we shall deduce (8.4) and (8.5) respectively.

Since the solution  $u$  to (8.2) is in  $H^2(\Omega)$  and  $\beta_h = 1$ , we deduce from (7.7) that

$$|u - u_h|_{1,\Omega} \leq Ch \|f\|_{0,\Omega}.$$

So we have proved (8.4).

The residual  $\|F(u_h)\|_{-1,\Omega}$  is estimated following the method presented in BARANGER and EL AMRI [1991]. For  $v \in H_0^1(\Omega)$  we have

$$\begin{aligned} \langle F(u_h), v \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} &= \int_{\Omega} \mathbf{grad} u_h \mathbf{grad} v \, dx + \int_{\Omega} (u_h^3 - f)v \, dx \\ &= \int_{\Omega} \mathbf{grad} u_h \mathbf{grad}(v - \tilde{r}_h v) \, dx + \int_{\Omega} (u_h^3 - f)(v - \tilde{r}_h v) \, dx, \end{aligned}$$

where  $\tilde{r}_h \in \mathcal{L}(H_0^1(\Omega), V_h)$  is the Clément interpolation operator onto  $V_h$ , see CLÉMENT [1975]. So for all  $v \in H_0^1(\Omega)$

$$(8.6) \quad \langle F(u_h), v \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} = \sum_{T \in \mathcal{T}_h} \left\{ \int_T \mathbf{grad} u_h \mathbf{grad}(v - \tilde{r}_h v) \, dx + \int_T (u_h^3 - f)(v - \tilde{r}_h v) \, dx \right\}.$$



For all  $T \in \mathcal{T}_h$ ,

$$\int_T \mathbf{grad} u_h \mathbf{grad}(v - \tilde{r}_h v) dx = - \int_T \Delta u_h (v - \tilde{r}_h v) dx + \int_{\partial T} \frac{\partial u_h}{\partial n_T} (v - \tilde{r}_h v) ds,$$

where  $n_T$  is the outward unit normal along the sides of  $T$ . Since  $u_h$  is in the space  $\mathcal{P}_1$  on  $T$  the above equality reads

$$\int_T \mathbf{grad} u_h \mathbf{grad}(v - \tilde{r}_h v) dx = \int_{\partial T} \frac{\partial u_h}{\partial n_T} (v - \tilde{r}_h v) ds.$$

Let  $S_h$  be the set of all sides of triangles in  $\mathcal{T}_h$ . For each  $t \in S_h$  we choose the direction of a unit normal  $n$  and denote by  $h_t$  its length, then

$$\left| \sum_{T \in \mathcal{T}_h} \int_T \mathbf{grad} u_h \mathbf{grad}(v - \tilde{r}_h v) dx \right| \leq \sum_{t \in S_h} \int_t \left| \left[ \frac{\partial u_h}{\partial n} \right] (v - \tilde{r}_h v) \right| ds;$$

if  $t$  is a side on  $\partial\Omega$ , we naturally set to zero the value outside  $\bar{\Omega}$ . Then the Cauchy Schwarz inequality gives

$$(8.7) \quad \left| \sum_{T \in \mathcal{T}_h} \int_T \mathbf{grad} u_h \mathbf{grad}(v - \tilde{r}_h v) dx \right| \leq \left( \sum_{t \in S_h} h_t \left\| \left[ \frac{\partial u_h}{\partial n} \right] \right\|_{0,t}^2 \right)^{1/2} \left( \sum_{t \in S_h} h_t^{-1} \|v - \tilde{r}_h v\|_{0,t}^2 \right)^{1/2}.$$

Using the technique of the reference element, see CIARLET [1978] or CIARLET [1991], it is a simple matter to prove the existence of a constant  $C$  such that for all  $T \in \mathcal{T}_h$  with edges  $t_1, t_2$ , and  $t_3$ , and for all  $w \in H^1(T)$  we have

$$\sum_{j=1}^3 h_{t_j}^{-1} \|w\|_{0,t_j}^2 \leq C (h_T^{-2} \|w\|_{0,T}^2 + |w|_{1,T}^2).$$

So the inequality (8.7) is transformed into

$$(8.8) \quad \left| \sum_{T \in \mathcal{T}_h} \int_T \mathbf{grad} u_h \mathbf{grad}(v - \tilde{r}_h v) dx \right| \leq C \left( \sum_{t \in S_h} h_t \left\| \left[ \frac{\partial u_h}{\partial n} \right] \right\|_{0,t}^2 \right)^{1/2} \left( \sum_{T \in \mathcal{T}_h} (h_T^{-2} \|v - \tilde{r}_h v\|_{0,T}^2 + |v - \tilde{r}_h v|_{1,T}^2) \right)^{1/2}.$$

With Clément's interpolation estimates, see CLÉMENT [1975], for  $v \in H^1(\Omega)$

$$\|v - \tilde{r}_h v\|_{0,T} \leq Ch \sum_{T' \in S_T} |v|_{1,T'} \quad \text{and} \quad |v - \tilde{r}_h v|_{1,T} \leq Ch \sum_{T' \in S_T} |v|_{1,T'},$$

where  $S_T$  is the set of triangles  $T'$  with  $T' \cap T \neq \emptyset$ , and with (8.8), the relation (8.6) is transformed by using the Cauchy Schwarz inequality into

$$(8.9) \quad \langle F(u_h), v \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} \leq C \left[ \left( \sum_{t \in S_h} h_t \left\| \left[ \frac{\partial u_h}{\partial n} \right] \right\|_{0,t}^2 \right)^{1/2} |v|_{1,\Omega} + \left( \sum_{T \in \mathcal{T}_h} h_T^2 \|u_h^3 - f\|_{0,T}^2 \right)^{1/2} |v|_{1,\Omega} \right].$$

Finally the estimate (8.5) is a consequence of (7.8), the fact that  $DF(u_h)$  for  $h \leq h_0$  is an isomorphism from  $H_0^1(\Omega)$  onto  $H^{-1}(\Omega)$  with a uniformly bounded inverse because  $\lim_{h \rightarrow 0} u_h = u$  and the estimate (8.9).  $\square$

*Remark 8.1.* Although the a posteriori estimate (8.5) is a crude one, it can justify an adaptive mesh refinement procedure, based on the minimization of the residual. The theory of a posteriori estimates for nonlinear problems has advanced considerably in recent years. Model problems in one space dimension are studied in BABUŠKA and RHEINBOLDT [1982] and in RHEINBOLDT [1981]. An other approach is presented in RHEINBOLDT [1985] which produces reliable a posteriori estimates. We refer to the review paper of VERFÜRTH [1993] for recent results on a posteriori estimates.  $\square$

It is also possible to obtain a priori estimates in other norms than the energy norm  $|\cdot|_{1,\Omega}$ . For instance with duality arguments of Aubin-Nitsche, see Remark 7.3, we shall derive an error estimate in the  $L^2(\Omega)$  norm.

**Theorem 8.2.** *Let  $h_0, \delta_0$ , and for all  $h \leq h_0$   $u_h \in \overline{B}(u, \delta_0)$  be given in Theorem 8.1. There exists a constant  $C$  such that*

$$(8.10) \quad \|u - u_h\|_{0,\Omega} \leq Ch^2.$$

*Proof.* Since the form  $b : H_0^1(\Omega) \times H_0^1(\Omega) \rightarrow \mathbb{R}$ ,  $b(v, w) = \langle DF(u)v, w \rangle_{H^{-1}(\Omega) H_0^1(\Omega)}$ , is coercive on  $H_0^1(\Omega)$ , there exists a unique solution  $w$  to the problem: find  $w \in H_0^1(\Omega)$  satisfying

$$(8.11) \quad \text{for all } v \in H_0^1(\Omega) \quad b(v, w) = \int_{\Omega} (u - u_h)v \, dx.$$

By the standard elliptic regularity, we have  $w \in H_0^1(\Omega) \cap H^2(\Omega)$  and

$$(8.12) \quad \|w\|_{2,\Omega} \leq C \|u - u_h\|_{0,\Omega};$$

for ease of exposition we have assumed that  $\Omega$  is a convex polygonal domain in  $\mathbb{R}^2$ , see the definition of  $V_h$  in (3.12). If we choose  $w_h \in V_h$  to be the solution of

$$\text{for all } v \in V_h \quad b(v, w_h) = \int_{\Omega} (u - u_h)v \, dx,$$

the following estimate holds

$$\|w - w_h\|_{1,\Omega} \leq Ch\|w\|_{2,\Omega}$$

and with (8.12)

$$(8.13) \quad \|w - w_h\|_{1,\Omega} \leq Ch\|u - u_h\|_{0,\Omega}.$$

Starting from (8.11), we deduce from the inequality (8.13)

$$(8.14) \quad \|u - u_h\|_{0,\Omega}^2 = b(u - u_h, w) \leq Ch\|u - u_h\|_{0,\Omega}\|u - u_h\|_{1,\Omega} + |b(u - u_h, w_h)|.$$

Since

$$F(u) = 0 \quad \text{and} \quad \text{for all } v_h \in V_h \quad \langle F(u_h), v_h \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} = 0,$$

we get with the  $C^2$  regularity of  $F$  and for  $w_h \in V_h$

$$\begin{aligned} b(u - u_h, w_h) &= \langle DF(u)(u - u_h), w_h \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} = \langle F(u) - F(u_h), w_h \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} \\ &\quad + \left\langle \int_0^1 (1-t)D^2F(u + t(u_h - u))(u_h - u)^2 dt, w_h \right\rangle_{H^{-1}(\Omega)H_0^1(\Omega)} \end{aligned}$$

and so

$$(8.15) \quad |b(u - u_h, w_h)| \leq C\|u_h - u\|_{1,\Omega}^2 [\|w_h - w\|_{1,\Omega} + \|w\|_{1,\Omega}].$$

With (8.4), (8.12), (8.13), (8.14), and (8.15) we finally get (8.10).  $\square$

To study the approximation scheme with numerical integration (3.14), we cannot use directly the result of Section 7 relative to Galerkin approximation. We have to adapt the formalism and use the general results of Section 6. The elliptic projection operator  $\pi_h \in \mathcal{L}(H_0^1(\Omega); V_h)$  is defined for  $v \in H_0^1(\Omega)$  by

$$\text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad}(v - \pi_h v) \mathbf{grad} v_h dx = 0.$$

We assume that  $f$  is in  $C^0(\overline{\Omega})$  and then we construct the mapping  $F_h : H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  by the formula, for  $v, w \in H_0^1(\Omega)$ ,

$$(8.16) \quad \langle F_h(v), w \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w dx + \int_{\Omega} r_h [((\pi_h v)^3 - f) \pi_h w] dx.$$

It is a simple matter to check that problem (3.14) is equivalent to find  $u_h \in H_0^1(\Omega)$  satisfying

$$(8.17) \quad \text{for all } v \in H_0^1(\Omega) \quad \langle F_h(u_h), v \rangle_{H^{-1}(\Omega)H_0^1(\Omega)} = 0.$$

Indeed let  $u_h \in V_h$  be a solution to (3.14). Then clearly  $u_h$  is also a solution to (8.17). Conversely let  $u_h \in H_0^1(\Omega)$  be a solution to (8.17). If the test function  $v$  in (8.17) has the form  $w - \pi_h w$  with  $w \in H_0^1(\Omega)$ , then the equation (8.17) reads

$$\int_{\Omega} \mathbf{grad} u_h \mathbf{grad} (w - \pi_h w) dx = 0,$$

which implies

$$u_h = \pi_h u_h \in V_h.$$

If we take  $v_h \in V_h$  to be a test function in (8.17) then we get (3.14). Thus a solution  $u_h \in H_0^1(\Omega)$  of (8.17) is a solution to (3.14).

Before studying problem (8.17), we present a technical result.

**Theorem 8.3.** *Let  $V_h$  be the space of continuous piecewise polynomials of degree  $\leq 1$  given in (3.12) and  $r_h \in \mathcal{L}(C^0(\bar{\Omega}); V_h)$  be the Lagrange interpolation operator. Then for  $M \in \mathbb{N}$ , there exists a constant  $C$  such that for any  $u_j \in V_h$   $j = 1, 2, \dots, M$ , we have*

$$\left| \int_{\Omega} \left( r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right) \right) dx \right| \leq Ch \prod_{j=1}^M \|u_j\|_{1,\Omega}.$$

*Proof.* Since

$$\left| \int_{\Omega} \left( r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right) \right) dx \right| \leq \sum_{T \in \mathcal{T}_h} \int_T \left| r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right) \right| dx,$$

we derive an estimate for each  $T \in \mathcal{T}_h$ . On an element  $T$  with the reference element technique, we get

$$\int_T \left| r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right) \right| dx = |\det B| \int_{\hat{T}} \left| \hat{r} \left( \prod_{j=1}^M \hat{u}_j \right) - \left( \prod_{j=1}^M \hat{u}_j \right) \right| d\hat{x},$$

where  $x \equiv F_T(\hat{x}) = B\hat{x} + b$  is the affine mapping from the reference element  $\hat{T}$  onto  $T$ ,  $\hat{r}$  the Lagrange  $\hat{\mathcal{P}}_1$ -interpolant on  $\hat{T}$ , and  $\hat{u}_j(\hat{x}) = u_j(x)$ ;  $\hat{\mathcal{P}}_\ell$  denotes the space of polynomials of degree  $\leq \ell$  on  $\hat{T}$ . For any  $\hat{p} \in \hat{\mathcal{P}}_0$ , we have

$$\int_T \left| r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right) \right| dx = |\det B| \int_{\hat{T}} |(I - \hat{r}) \left( \left( \prod_{j=1}^M \hat{u}_j \right) - \hat{p} \right)| d\hat{x}.$$

Clearly the mapping  $\hat{r} : \hat{\mathcal{P}}_M \rightarrow \hat{\mathcal{P}}_1$  is linear and continuous. So

$$\int_T \left| r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right) \right| dx \leq C |\det B| \left\| \left( \prod_{j=1}^M \hat{u}_j \right) - \hat{p} \right\|_{1,1,\hat{T}}$$

for any  $\hat{p} \in \hat{\mathcal{P}}_0$ . Taking the infimum over  $\hat{p}$  as in the Bramble-Hilbert Lemma and using the affine change of variables we get

$$\int_T |r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right)| dx \leq C |\det B| \left| \left( \prod_{j=1}^M \hat{u}_j \right) \right|_{1,1,\hat{T}} \leq C \|B\| \left| \left( \prod_{j=1}^M u_j \right) \right|_{1,1,T},$$

where  $\|\cdot\|$  is the matrix norm induced by the Euclidian vector norm. Noticing that

$$\|B\| \leq Ch_T$$

and that

$$\begin{aligned} \int_T |\mathbf{grad} \left( \prod_{j=1}^M u_j \right)| dx &\leq \sum_{\sigma \in \Delta} \int_T |\mathbf{grad} u_{\sigma_1}| \left| \prod_{\ell=2}^M u_{\sigma_\ell} \right| dx \\ &\leq \sum_{\sigma \in \Delta} \left( \|\mathbf{grad} u_{\sigma_1}\|_{0,T} \prod_{\ell=2}^M \|u_{\sigma_\ell}\|_{0,2(M-1),T} \right) \end{aligned}$$

with  $\Delta$  the set of circular permutations of  $\{1, 2, \dots, M\}$ , we write successively

$$\begin{aligned} &\left| \int_{\Omega} \left( r_h \left( \prod_{j=1}^M u_j \right) - \left( \prod_{j=1}^M u_j \right) \right) dx \right| \\ &\leq Ch \sum_{T \in \mathcal{T}_h} \left( \sum_{\sigma \in \Delta} |u_{\sigma_1}|_{1,T} \prod_{\ell=2}^M \|u_{\sigma_\ell}\|_{0,2(M-1),T} \right) \leq Ch \sum_{\sigma \in \Delta} |u_{\sigma_1}|_{1,\Omega} \prod_{\ell=2}^M \|u_{\sigma_\ell}\|_{0,2(M-1),\Omega} \\ &\leq \tilde{C}h \sum_{\sigma \in \Delta} \left( \prod_{j=1}^M \|u_{\sigma_j}\|_{1,\Omega} \right) = \tilde{C}hM \prod_{j=1}^M \|u_j\|_{1,\Omega}. \end{aligned}$$

and Theorem 7.3 is proved.  $\square$

As a consequence of Theorem 6.1, we get the following result.

**Theorem 8.4.** *In the framework developed above, assume that  $f$  is in  $W^{1,p}(\Omega)$  for some  $p > 2$ . Then there exist positive constants  $h_0$ ,  $\delta_0$ ,  $C$ , and for  $h \leq h_0$  a unique solution  $\tilde{u}_h$  to (8.17) (or equivalently to (3.14)) in  $\overline{B}(u, \delta_0)$ . Moreover the following a priori estimate holds*

$$(8.18) \quad \|u - \tilde{u}_h\|_{1,\Omega} \leq Ch.$$

*Proof.* The result will be a consequence of Theorem 6.1. Thus we shall check the hypotheses (6.1), (6.2), and (6.3). Note that for  $v \in H_0^1(\Omega)$ , the derivative  $DF_h(v) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is given by

$$\begin{aligned} \text{for all } w_1, w_2 \in H_0^1(\Omega) \quad &\langle DF_h(v)w_1, w_2 \rangle \\ &= \int_{\Omega} \mathbf{grad} w_1 \mathbf{grad} w_2 dx + 3 \int_{\Omega} r_h [(\pi_h v)^2 \pi_h w_1 \pi_h w_2] dx, \end{aligned}$$

with the notation  $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{H^{-1}(\Omega)H_0^1(\Omega)}$ .

Let us first check (6.1) with  $L_h = \text{constant}$ . By the definitions for  $v \in H_0^1(\Omega)$  we have

$$\|DF_h(u) - DF_h(v)\|_{H_0^1(\Omega); H^{-1}(\Omega)} = \sup_{\substack{w_1 \in H_0^1(\Omega) \\ |w_1|_{1,\Omega} \leq 1}} \sup_{\substack{w_2 \in H_0^1(\Omega) \\ |w_2|_{1,\Omega} \leq 1}} \langle [DF_h(u) - DF_h(v)] w_1, w_2 \rangle.$$

Let  $v, w_1, w_2$  in  $H_0^1(\Omega)$ , then following the technique used in the proof of Theorem 8.3, we get

$$\begin{aligned} \langle [DF_h(u) - DF_h(v)] w_1, w_2 \rangle &= 3 \int_{\Omega} r_h [((\pi_h u)^2 - (\pi_h v)^2) \pi_h w_1 \pi_h w_2] dx \\ &\leq C \|\pi_h(u - v)\|_{0,4,\Omega} \|\pi_h(u + v)\|_{0,4,\Omega} \|\pi_h w_1\|_{0,4,\Omega} \|\pi_h w_2\|_{0,4,\Omega}. \end{aligned}$$

Since the  $L^4(\Omega)$  norm is bounded by the  $H^1(\Omega)$  norm and  $\pi_h \in \mathcal{L}(H_0^1(\Omega); H_0^1(\Omega))$  is bounded with respect to  $h$ , we finally get for  $v \in \overline{B}(u, 1)$

$$(8.19) \quad \|DF_h(u) - DF_h(v)\|_{H_0^1(\Omega); H^{-1}(\Omega)} \leq C |u - v|_{1,\Omega},$$

which implies (6.1) with  $L_h = C$  and  $\eta = 1$ .

To check the hypothesis (6.3), we notice that, since  $DF(u) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is an isomorphism, it is sufficient to check that  $\lim_{h \rightarrow 0} DF_h(u) = DF(u)$ . For  $w_1, w_2 \in H_0^1(\Omega)$  we have

$$\begin{aligned} \langle [DF_h(u) - DF(u)] w_1, w_2 \rangle &= 3 \int_{\Omega} r_h [(\pi_h u)^2 \pi_h w_1 \pi_h w_2] dx \\ &\quad - 3 \int_{\Omega} (\pi_h u)^2 \pi_h w_1 \pi_h w_2 dx + 3 \int_{\Omega} (\pi_h u)^2 \pi_h w_1 \pi_h w_2 dx - 3 \int_{\Omega} u^2 w_1 w_2 dx. \end{aligned}$$

Making use of Theorem 8.3 and of the fact that  $\pi_h \in \mathcal{L}(H_0^1(\Omega); H_0^1(\Omega))$  is uniformly bounded with respect to  $h$ , we can prove that the first difference in the right hand side of the last equation tends to zero. The second difference reads

$$\begin{aligned} &\int_{\Omega} (\pi_h u)^2 \pi_h w_1 \pi_h w_2 dx - \int_{\Omega} u^2 w_1 w_2 dx = \int_{\Omega} (\pi_h u - u)(\pi_h u + u) w_1 w_2 dx \\ &\quad + \int_{\Omega} (\pi_h u)^2 (\pi_h w_1 - w_1) w_2 dx + \int_{\Omega} (\pi_h u)^2 \pi_h w_1 (\pi_h w_2 - w_2) dx \\ &\leq \|u - \pi_h u\|_{0,4,\Omega} \|u + \pi_h u\|_{0,4,\Omega} \|w_1\|_{0,4,\Omega} \|w_2\|_{0,4,\Omega} + \|\pi_h u\|_{0,4,\Omega}^2 \|w_1 - \pi_h w_1\|_{0,4,\Omega} \|w_2\|_{0,4,\Omega} \\ &\quad + \|\pi_h u\|_{0,4,\Omega}^2 \|\pi_h w_1\|_{0,4,\Omega} \|\pi_h w_2 - w_2\|_{0,4,\Omega}. \end{aligned}$$

The family of triangulations being regular, we can check that

$$\lim_{h \rightarrow 0} \sup_{\substack{w \in H_0^1(\Omega) \\ w \neq 0}} \frac{\|w - \pi_h w\|_{0,4,\Omega}}{|w|_{1,\Omega}} = 0.$$

So the hypothesis (6.3) has been verified.

Let us finally check the hypothesis (6.2). For  $v \in H_0^1(\Omega)$  with  $|v|_{1,\Omega} \leq 1$

$$\begin{aligned} \langle F_h(u), v \rangle &= \int_{\Omega} \mathbf{grad} u \mathbf{grad} v \, dx + \int_{\Omega} r_h [((\pi_h u)^3 - f) \pi_h v] \, dx \\ &= \int_{\Omega} \{r_h [(\pi_h u)^3 \pi_h v] - (\pi_h u)^3 \pi_h v\} \, dx + \int_{\Omega} [(\pi_h u)^3 - u^3] \pi_h v \, dx \\ &\quad - \int_{\Omega} \{r_h [f \pi_h v] - f \pi_h v\} \, dx + \int_{\Omega} \mathbf{grad}(u - \pi_h u) \mathbf{grad} v \, dx. \end{aligned}$$

We study successively all four terms in the last equality. We have

$$(8.20) \quad \left| \int_{\Omega} \mathbf{grad}(u - \pi_h u) \mathbf{grad} v \, dx \right| \leq |u - \pi_h u|_{1,\Omega} \leq Ch.$$

For the third difference we notice that, with  $v_h = \pi_h v$ ,

$$\begin{aligned} \int_{\Omega} [r_h(f v_h) - f v_h] \, dx &= \int_{\Omega} [r_h(r_h f v_h) - f v_h] \, dx \\ &= \int_{\Omega} [r_h(r_h f v_h) - r_h f v_h] \, dx + \int_{\Omega} (r_h f - f) v_h \, dx \end{aligned}$$

so with Theorem 8.3 we get the estimate

$$(8.21) \quad \begin{aligned} \left| \int_{\Omega} [r_h(f v_h) - f v_h] \, dx \right| &\leq Ch \|r_h f\|_{1,\Omega} \|v_h\|_{1,\Omega} + \|f - r_h f\|_{0,\Omega} \|v_h\|_{0,\Omega} \\ &\leq Ch |v_h|_{1,\Omega}. \end{aligned}$$

The second difference can be estimated in the following way

$$(8.22) \quad \left| \int_{\Omega} [(\pi_h u)^3 - u^3] \pi_h v \, dx \right| \leq Ch,$$

where  $C$  depends on  $u$  but is independent of  $h$ . As a consequence of Theorem 8.3 an estimate for the first difference reads

$$(8.23) \quad \left| \int_{\Omega} \{r_h [(\pi_h u)^3 \pi_h v] - (\pi_h u)^3 \pi_h v\} \, dx \right| \leq Ch.$$

As a consequence of (8.20) through (8.23), we get

$$(8.24) \quad \|F_h(u)\|_{-1,\Omega} \leq Ch,$$

where  $C$  is independent of  $h$  ( but actually depends on  $u$  ).

We apply Theorem 6.1 and the estimate (8.24) with (6.5) to conclude.  $\square$

For a model problem we have developed a method to get a priori and a posteriori error estimates in the cases of Galerkin approximation without and with numerical integration. The interest of the technique lies in its ability to be adapted to other cases.

### 9. Estimates for an approximation of the Navier-Stokes problem

The abstract results of Section 7 are applied now to the commonly called  $(4\mathcal{P}_1 - \mathcal{P}_1)$  approximation of the Navier-Stokes problem introduced in Section 4.

In this case we are using the following notations:  $L_0^2(\Omega) = \{f \in L^2(\Omega); \int_{\Omega} f \, dx = 0\}$ ,  $X = Y = (H_0^1(\Omega))^2 \times L_0^2(\Omega)$ , and  $F : X \rightarrow X'$  is defined for  $\underline{u} = (\mathbf{u}, p) \in X$ ,

$$(9.1) \quad \begin{aligned} \text{for all } \underline{v} = (\mathbf{v}, q) \in X \quad & \langle F(\underline{u}, p), (\mathbf{v}, q) \rangle_{X'X} = \nu \int_{\Omega} \mathbf{grad} \mathbf{u} \mathbf{grad} \mathbf{v} \, dx \\ & + \int_{\Omega} (\mathbf{u} \nabla) \mathbf{u} \mathbf{v} \, dx - \int_{\Omega} p \operatorname{div} \mathbf{v} \, dx - \int_{\Omega} \mathbf{f} \mathbf{v} \, dx - \int_{\Omega} q \operatorname{div} \mathbf{u} \, dx \\ & \equiv a(\mathbf{u}, \mathbf{v}) + c(\mathbf{u}; \mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) + b(\mathbf{u}, q) - \int_{\Omega} \mathbf{f} \mathbf{v} \, dx. \end{aligned}$$

Let  $\underline{u} = (\mathbf{u}, p)$  denote a regular solution to  $F(\underline{v}) = 0$ , in the sense that (4.8) is satisfied.

The approximation spaces and the test spaces are  $X_h \equiv Y_h = V_h^2 \times W_h$ , where the finite element spaces  $V_h$  and  $W_h$  are defined in Section 4. A Galerkin approximation of the problem  $F(\underline{v}) = 0$  reads (see (4.9)): find  $\underline{u}_h \in X_h$  such that

$$(9.2) \quad \text{for all } \underline{v}_h \in X_h \quad \langle F(\underline{u}_h), \underline{v}_h \rangle_{X'X} = 0.$$

**Theorem 9.1.** *Let  $\underline{u} = (\mathbf{u}, p) \in X$  be a solution to  $F(\underline{v}) = 0$  with  $\mathbf{u} \in (H^2(\Omega))^2$ . We assume that the hypothesis (4.8) and the assumption (4.10) on the triangulations are satisfied. Then there exist constants  $h_0 > 0$ ,  $\delta_0 > 0$ ,  $C > 0$ , and for  $h \leq h_0$  a unique solution  $\underline{u}_h = (\mathbf{u}_h, p_h) \in X_h$  to the problem (9.2) satisfying  $\|\underline{u} - \underline{u}_h\|_X \leq \delta_0$ . Furthermore the following estimates hold*

$$(9.3) \quad \|\mathbf{u} - \mathbf{u}_h\|_{1,\Omega} + \|p - p_h\|_{0,\Omega} \leq C \inf_{(\mathbf{v}_h, q_h) \in X_h} (\|\mathbf{u} - \mathbf{v}_h\|_{1,\Omega} + \|p - q_h\|_{0,\Omega})$$

and

$$(9.4) \quad \|\mathbf{u} - \mathbf{u}_h\|_{1,\Omega} + \|p - p_h\|_{0,\Omega} \leq C \left( \sum_{T \in \mathcal{T}_h^f} \eta(T)^2 \right)^{1/2},$$

where for any triangle  $T \in \mathcal{T}_h^f$  the estimator  $\eta(T)$  is expressed by

$$\eta(T)^2 = h_T^2 \left\| -\nu \Delta \mathbf{u}_h + (\mathbf{u}_h \nabla) \mathbf{u}_h + \mathbf{grad} p_h - \mathbf{f} \right\|_{0,T}^2 + \sum_{i=1}^3 h_{t_i} \left\| \left[ \frac{\partial \mathbf{u}_h}{\partial n} \right] \right\|_{0,t_i}^2 + \|\operatorname{div} \mathbf{u}_h\|_{0,T}^2$$

with  $h_T$  the diameter of  $T$ ,  $t_i$   $i = 1, 2, 3$ , the edges of  $T$  of length  $h_{t_i}$ , and  $[\cdot]$  the jump through an edge.

*Proof.* We will first check that we are in a position to apply Theorem 7.1. Since we assume that  $\underline{u} = (\mathbf{u}, p) \in X$  satisfies (4.8), so the hypothesis (7.1) is satisfied. The



assumption (7.4) is verified in Theorem 4.1 while the (7.5) and (7.6) ones are easy to check with  $\beta_h$  independent of  $h$ . We can apply Theorem 7.1 with  $\beta_h$  constant.

The estimate (9.3) is immediately deduced from (7.7) while from (7.8) we shall deduce the estimate (9.4). In our case (7.8) reads

$$(9.5) \quad |\mathbf{u} - \mathbf{u}_h|_{1,\Omega} + \|p - p_h\|_{0,\Omega} \leq C \|F(\mathbf{u}_h, p_h)\|_{X'}.$$

For  $\underline{v} = (\mathbf{v}, q) \in X$  and  $\underline{v}_h = (\mathbf{v}_h, q_h) \in X_h$  with  $\mathbf{v}_h = \tilde{r}_h \mathbf{v}$ , where  $\tilde{r}_h \in \mathcal{L}((H_0^1(\Omega))^2; V_h^2)$  is the Clément interpolation operator for discontinuous functions, see CLÉMENT [1975], we have

$$(9.6) \quad \begin{aligned} \langle F(\underline{u}_h), \underline{v} \rangle_{X'X} &= \langle F(\underline{u}_h), \underline{v} - \underline{v}_h \rangle_{X'X} \\ &= a(\mathbf{u}_h, \mathbf{v} - \mathbf{v}_h) + c(\mathbf{u}_h; \mathbf{u}_h, \mathbf{v} - \mathbf{v}_h) + b(\mathbf{v} - \mathbf{v}_h, p_h) + b(\mathbf{u}_h, q - q_h) - \int_{\Omega} \mathbf{f}(\mathbf{v} - \mathbf{v}_h) dx \\ &= \sum_{T \in \mathcal{T}_h^f} \left\{ \int_T (-\nu \Delta \mathbf{u}_h + (\mathbf{u}_h \nabla) \mathbf{u}_h + \mathbf{grad} p_h - \mathbf{f})(\mathbf{v} - \mathbf{v}_h) dx \right. \\ &\quad \left. - \int_T \operatorname{div} \mathbf{u}_h (q - q_h) dx + \nu \int_{\partial T} \frac{\partial \mathbf{u}_h}{\partial n} (\mathbf{v} - \mathbf{v}_h) ds \right\}. \end{aligned}$$

Let us bound each of these terms. Setting  $\mathbf{w} = -\nu \Delta \mathbf{u}_h + (\mathbf{u}_h \nabla) \mathbf{u}_h + \mathbf{grad} p_h - \mathbf{f}$  on  $T$ , we get successively since  $\mathbf{v}_h = \tilde{r}_h \mathbf{v}$

$$\int_T \mathbf{w}(\mathbf{v} - \mathbf{v}_h) dx \leq \|\mathbf{w}\|_{0,T} \|\mathbf{v} - \mathbf{v}_h\|_{0,T} \leq C \|\mathbf{w}\|_{0,T} h \sum_{T' \in S_T} \|\mathbf{v}\|_{1,T'}$$

where  $S_T$  is the set of triangles  $T'$  in  $\mathcal{T}_h^f$  with  $T' \cap T \neq \emptyset$ , and summing over all the elements

$$(9.7) \quad \begin{aligned} \left| \sum_{T \in \mathcal{T}_h^f} \int_T \mathbf{w}(\mathbf{v} - \mathbf{v}_h) dx \right| &\leq C \left( \sum_{T \in \mathcal{T}_h^f} h_T^2 \|\mathbf{w}\|_{0,T}^2 \right)^{1/2} \left( \sum_{T \in \mathcal{T}_h^f} \|\mathbf{v}\|_{1,T}^2 \right)^{1/2} \\ &\leq \tilde{C} |\mathbf{v}|_{1,\Omega} \left( \sum_{T \in \mathcal{T}_h^f} h_T^2 \|\mathbf{w}\|_{0,T}^2 \right)^{1/2}. \end{aligned}$$

We use the same techniques as the ones developed in the proof of Theorem 8.1 to estimate

$$(9.8) \quad \left| \sum_{T \in \mathcal{T}_h^f} \nu \int_{\partial T} \frac{\partial \mathbf{u}_h}{\partial n} (\mathbf{v} - \tilde{r}_h \mathbf{v}) ds \right| \leq C |\mathbf{v}|_{1,\Omega} \left( \sum_{t \in S_h} h_t \left\| \left[ \frac{\partial \mathbf{u}_h}{\partial n} \right] \right\|_{0,t}^2 \right)^{1/2}$$

where  $S_h$  is the set of all sides  $t$  of triangles in  $\mathcal{T}_h^f$ . Since  $q$  is only a  $L_0^2(\Omega)$  function, we take  $q_h = 0$  and so

$$(9.9) \quad \left| \sum_{T \in \mathcal{T}_h^f} \int_T \operatorname{div} \mathbf{u}_h (q - q_h) dx \right| \leq \|q\|_{0,\Omega} \left( \sum_{T \in \mathcal{T}_h^f} \|\operatorname{div} \mathbf{u}_h\|_{0,T}^2 \right)^{1/2}.$$

From the estimates (9.7), (9.8), and (9.9) we get the following bound for (9.6) noticing that  $\Delta \mathbf{u}_h = 0$  on  $T$ ,

$$(9.10) \quad |\langle F(\underline{u}_h), \underline{v} \rangle_{X'X}| \leq C \left\{ \sum_{T \in \mathcal{T}_h^f} \left[ h_T^2 \|(\mathbf{u}_h \nabla) \mathbf{u}_h + \mathbf{grad} p_h - \mathbf{f}\|_{0,T}^2 \right. \right. \\ \left. \left. + \sum_{i=1}^3 h_{t_i} \left\| \left[ \frac{\partial \mathbf{u}_h}{\partial n} \right] \right\|_{0,t_i}^2 + \|\operatorname{div} \mathbf{u}_h\|_{0,T}^2 \right] \right\}^{1/2} \|(\mathbf{v}, q)\|_X.$$

The estimates (9.10) and (9.5) imply the error estimate (9.4).  $\square$

### 10. Error estimates for the stationary heat problem

The abstract results of Section 7 are applied now to the approximation (5.18) of the stationary heat problem with convection. Our framework is well-suited to study this effective Galerkin approximation.

We are using the following notations:  $X = W_0^{1,p}(\Omega)$ ,  $Y = W_0^{1,q}(\Omega)$  with  $p, q$  in  $\mathbb{N}$ ,  $p > 2$  and  $\frac{1}{p} + \frac{1}{q} = 1$ , and  $F : X \rightarrow Y'$  defined for  $v \in W_0^{1,p}(\Omega)$  by for all  $w \in W_0^{1,q}(\Omega)$

$$(10.1) \quad \langle F(v), w \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} = \int_{\Omega} k(v) \mathbf{grad} v \mathbf{grad} w dx \\ + \int_{\Omega} \mathbf{c} \mathbf{grad} v w dx - \int_{\Omega} f w dx.$$

When the domain  $\Omega$  is regular, the hypotheses (5.2), (5.3), and (5.4) satisfied, and  $f$  in  $L^\infty(\Omega)$ , then the problem  $F(v) = 0$  has a unique solution  $u \in W_0^{1,p}(\Omega)$ , which is in  $W^{2,r}(\Omega)$  with  $1 \leq r < \infty$  and  $DF(u) \in \mathcal{L}(W_0^{1,p}(\Omega); W^{-1,p}(\Omega))$  is an isomorphism, see Theorem 5.2.

For the sake of simplicity, we assume that  $\Omega$  is a polygonal domain such that the Laplacian operator  $\Delta : W_0^{1,p}(\Omega) \rightarrow W^{-1,p}(\Omega)$  is an isomorphism. Given a family  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  of triangulations,  $h$  is the maximum of the diameter of  $T \in \mathcal{T}_h$ , the approximation spaces and the test spaces are the spaces  $V_h$  of continuous piecewise polynomials of degree 1. The family is assumed to be regular and quasi-uniform, see CIARLET [1978] p.132 and p.140. We study the Galerkin approximation: find  $u_h \in V_h$  satisfying

$$(10.2) \quad \text{for all } v_h \in V_h \quad \langle F(u_h), v_h \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} = 0.$$

*Remark 10.1.* In Theorem 5.2, we have assumed that  $\Omega$  is regular, so that the Laplacian operator with homogeneous boundary conditions is an isomorphism from  $W_0^{1,p}(\Omega)$  onto  $W^{-1,p}(\Omega)$ , see SIMADER [1972]. Here we assume  $\Omega$  polygonal so that  $\bar{\Omega} = \cup_{T \in \mathcal{T}_h} T$  and the Laplacian operator is an isomorphism, see DAUGE [1992]. Actually when  $\Omega$  is regular, we should use a family of isoparametric triangulations of  $\Omega$ .  $\square$

To apply Theorem 7.1, we need to prove a discrete inf-sup condition for the bilinear form  $b : W_0^{1,p}(\Omega) \times W_0^{1,q}(\Omega) \rightarrow \mathbb{R}$  given by: for  $v \in W_0^{1,p}(\Omega)$ ,  $w \in W_0^{1,q}(\Omega)$

$$b(v, w) = \langle DF(u)v, w \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)}.$$

**Theorem 10.1.** *Under the hypotheses of Theorem 5.1 let  $u \in W_0^{1,p}(\Omega)$  be the unique solution to  $F(v) = 0$  and  $DF(u) \in \mathcal{L}(W_0^{1,p}(\Omega); W^{-1,p}(\Omega))$  be an isomorphism. We assume that the family of triangulations is regular and quasi-uniform. Then there exists  $\zeta > 0$  independent of  $h$  such that*

$$(10.3) \quad \inf_{\substack{v_h \in V_h \\ |v_h|_{1,p,\Omega}=1}} \sup_{\substack{w_h \in V_h \\ |w_h|_{1,q,\Omega}=1}} b(v_h, w_h) \geq \zeta.$$

*Proof.* Going back to the definition of  $DF(u)$  in (5.12), we have for  $w \in W_0^{1,p}(\Omega)$ ,  $\psi \in W_0^{1,q}(\Omega)$

$$(10.4) \quad \langle DF(u)w, \psi \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} = \int_{\Omega} \mathbf{grad}(k(u)w) \mathbf{grad}\psi \, dx + \int_{\Omega} \mathbf{c} \mathbf{grad}w \psi \, dx.$$

Let  $-T \in \mathcal{L}(W^{-1,p}(\Omega); W_0^{1,p}(\Omega))$  denote the inverse of the Laplacian operator with homogeneous boundary condition, which is an isomorphism. Then (10.4) reads for  $w \in W_0^{1,p}(\Omega)$

$$(10.5) \quad TDF(u)w = k(u)w + T(\mathbf{c} \mathbf{grad}w).$$

If we introduce the mapping  $R : \varphi \in W_0^{1,p}(\Omega) \rightarrow R\varphi = \varphi/k(u) \in W_0^{1,p}(\Omega)$  which is an isomorphism, then (10.5) reads with  $w = R\varphi$

$$(10.6) \quad TDF(u)R\varphi = \varphi + T(\mathbf{c} \mathbf{grad}(R\varphi)).$$

We have proved in Theorem 5.2 that  $TDF(u)R \in \mathcal{L}(W_0^{1,p}(\Omega); W_0^{1,p}(\Omega))$  is an isomorphism.

Let  $T_h \in \mathcal{L}(W^{-1,p}(\Omega); V_h)$  be the discrete analogue of  $T$  defined for  $g \in W^{-1,p}(\Omega)$  by

$$\text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad}(T_h g) \mathbf{grad}v_h \, dx = \int_{\Omega} gv_h \, dx.$$

We easily check that for  $\varphi_h \in V_h$

$$T_h DF(u)R\varphi_h = \varphi_h + T_h(\mathbf{c} \mathbf{grad}(R\varphi_h)).$$

The family of triangulations is regular and quasi-uniform, so there exists a constant  $C > 0$  such that for all  $w \in W^{2,p}(\Omega)$

$$\|w - \pi_h w\|_{1,p,\Omega} \leq Ch \|w\|_{2,p,\Omega},$$

where  $\pi_h \in \mathcal{L}(W^{1,p}(\Omega); V_h)$  is the elliptic projector, see RANNACHER and SCOTT [1982].

Since we have  $T_h = \pi_h T$ , we get for  $\varphi_h \in V_h$

$$T_h DF(u) R\varphi_h = T DF(u) R\varphi_h + (T_h - T)(\mathbf{c} \mathbf{grad}(R\varphi_h))$$

and so there exists a constant  $\gamma_1$  independent of  $h$  such that

$$|T_h DF(u) R\varphi_h|_{1,p,\Omega} \geq \gamma_1 |\varphi_h|_{1,p,\Omega} - \epsilon_h \|\mathbf{c} \mathbf{grad}(R\varphi_h)\|_{0,p,\Omega}$$

with  $\epsilon_h$  tending to zero for  $h$  tending to zero, independent of  $\varphi_h$ . So for  $h$  small enough

$$(10.7) \quad |T_h DF(u) R\varphi_h|_{1,p,\Omega} \geq \frac{\gamma_1}{2} |\varphi_h|_{1,p,\Omega}.$$

We are working under the assumption that the Laplacian operator  $\Delta : W_0^{1,p}(\Omega) \rightarrow W^{-1,p}(\Omega)$  is an isomorphism, where for all  $v \in W_0^{1,p}(\Omega)$ , for all  $w \in W_0^{1,q}(\Omega)$

$$\langle \Delta v, w \rangle_{W^{-1,p}(\Omega) W_0^{1,q}(\Omega)} = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w \, dx.$$

So we can check there exists a  $\gamma_2 > 0$  with

$$\inf_{\substack{v \in W_0^{1,p}(\Omega) \\ |v|_{1,p,\Omega} = 1}} \sup_{\substack{w \in W_0^{1,q}(\Omega) \\ |w|_{1,q,\Omega} = 1}} \int_{\Omega} \mathbf{grad} v \mathbf{grad} w \, dx \geq \gamma_2 > 0.$$

The elliptic projector  $\pi_h \in \mathcal{L}(W_0^{1,q}(\Omega); V_h)$  is stable in  $W_0^{1,q}(\Omega)$ , which means there exists a constant  $C > 0$  such that for all  $v \in W_0^{1,q}(\Omega)$

$$|\pi_h v|_{1,q,\Omega} \leq C |v|_{1,q,\Omega},$$

see RANNACHER and SCOTT [1982]. Let  $v_h$  be in  $V_h$  with  $|v_h|_{1,p,\Omega} = 1$ . Then there exists a function  $w \in W_0^{1,q}(\Omega)$  with  $|w|_{1,q,\Omega} = 1$  such that

$$\frac{\gamma_2}{2} \leq \int_{\Omega} \mathbf{grad} v_h \mathbf{grad} w \, dx = \int_{\Omega} \mathbf{grad} v_h \mathbf{grad} \pi_h w \, dx.$$

So

$$\int_{\Omega} \mathbf{grad} v_h \mathbf{grad} \frac{\pi_h w}{|\pi_h w|_{1,q,\Omega}} \, dx \geq \frac{\gamma_2}{2 |\pi_h w|_{1,q,\Omega}}.$$

With the stability result of  $\pi_h$ , we conclude there exists a constant  $\gamma_3 > 0$  independent of  $h$  such that

$$(10.8) \quad \inf_{\substack{v_h \in V_h \\ |v_h|_{1,p,\Omega}=1}} \sup_{\substack{w_h \in V_h \\ |w_h|_{1,q,\Omega}=1}} \int_{\Omega} \mathbf{grad} v_h \mathbf{grad} w_h dx \geq \gamma_3.$$

Let the function  $\varphi_h$  be in  $V_h$ , then with the estimates (10.7) and (10.8) we get

$$\begin{aligned} \sup_{\substack{w_h \in V_h \\ |w_h|_{1,q,\Omega}=1}} b(R\varphi_h, w_h) &= \sup_{\substack{w_h \in V_h \\ |w_h|_{1,q,\Omega}=1}} \langle DF(u)R\varphi_h, w_h \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} \\ &= \sup_{\substack{w_h \in V_h \\ |w_h|_{1,q,\Omega}=1}} \int_{\Omega} \mathbf{grad}(T_h DF(u)R\varphi_h) \mathbf{grad} w_h dx \geq \gamma |\varphi_h|_{1,p,\Omega}, \end{aligned}$$

and so for all  $\varphi_h \in V_h$

$$(10.9) \quad \sup_{\substack{w_h \in V_h \\ |w_h|_{1,q,\Omega}=1}} b(R\varphi_h, w_h) \geq \gamma |\varphi_h|_{1,p,\Omega}.$$

Using the hypotheses (5.2) and (5.3) on  $k$ , the regularity of the triangulations and the regularity of the solution  $u \in W^{2,r}(\Omega)$ ,  $1 \leq r < \infty$ , we get by standard calculations on the reference element

$$(10.10) \quad \lim_{h \rightarrow 0} \max_{\substack{\chi_h \in V_h \\ |\chi_h|_{1,p,\Omega}=1}} \|R^{-1}\chi_h - r_h(R^{-1}\chi_h)\|_{1,p,\Omega} = 0,$$

where  $r_h$  is the Lagrange interpolation operator.

Finally the estimates (10.9) and (10.10) imply there exists  $\zeta > 0$  such that

$$\inf_{\substack{v_h \in V_h \\ |v_h|_{1,p,\Omega}=1}} \sup_{\substack{w_h \in V_h \\ |w_h|_{1,q,\Omega}=1}} b(v_h, w_h) \geq \zeta$$

and the proof is complete.  $\square$

We are in a position to prove the main result of the section.

**Theorem 10.2.** *We assume the hypotheses of Theorem 5.1 and denote by  $u \in W_0^{1,p}(\Omega)$ ,  $2 < p < \infty$ , the unique solution to the problem  $F(v) = 0$ . Furthermore we assume that the family  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  is regular and quasi-uniform. Then there exist positive constants  $h_0$ ,  $\delta_0$ ,  $C$ , and for  $h \leq h_0$  a unique solution  $u_h$  to problem (10.2) in the closed ball  $\overline{B}(u, \delta_0)$ . Moreover the following estimates hold*

$$(10.11) \quad |u - u_h|_{1,p,\Omega} \leq Ch$$

and

$$(10.12) \quad |u - u_h|_{1,p,\Omega} \leq C \left( \sum_{T \in \mathcal{T}_h} \eta(T)^p \right)^{1/p},$$

where for any  $T \in \mathcal{T}_h$  the estimator  $\eta(T)$  is expressed by

$$\begin{aligned} \eta(T) = h_T & \| -\operatorname{div}(k(u_h)\mathbf{grad}u_h) + \mathbf{c} \mathbf{grad}u_h - f \|_{0,p,T} \\ & + h_T^{1-2/q} \sum_{i=1}^3 h_{t_i}^{1/q} \left\| \left[ k(u_h) \frac{\partial u_h}{\partial n} \right] \right\|_{0,p,t_i}, \end{aligned}$$

where  $h_T$  is the diameter of the triangle  $T \in \mathcal{T}_h$ ,  $h_{t_i}$  is the length of the edge  $t_i$  of the triangle  $T$ ,  $i = 1, 2, 3$ , and  $[\cdot]$  the jump across the considered edge, (we adopt the convention of a zero function outside  $\bar{\Omega}$ ).

*Proof.* Let us prove that we are in a position to apply Theorem 7.1. The hypothesis (7.1) is verified in Theorems 5.1 and 5.2. The inf-sup condition (7.4) is verified in Theorem 10.1 with  $\beta_h = \beta$  independent of  $h$  while the hypothesis (7.5) is satisfied. With the classical interpolation result, we can verify (7.6). With the assumptions (5.2) and (5.3) on  $k$ , we easily check that  $DF$  is Lipschitzian at  $u$ .

The hypotheses of Theorem 7.1 are verified. The proof is complete once proved the two following error estimates

$$(10.13) \quad \inf_{v_h \in V_h} |u - v_h|_{1,p,\Omega} \leq Ch$$

and

$$(10.14) \quad \|F(u_h)\|_{-1,p,\Omega} \leq C \left( \sum_{T \in \mathcal{T}_h} \eta(T)^p \right)^{1/p}.$$

The estimate (10.13) is a consequence of the interpolation result, for  $v \in W^{2,p}(\Omega)$ ,

$$|v - r_h v|_{1,p,\Omega} \leq Ch \|v\|_{2,p,\Omega}.$$

Let us estimate now the term  $\|F(u_h)\|_{-1,p,\Omega}$  with the same method as the one used in the proof of Theorem 8.1. For  $w \in W_0^{1,q}(\Omega)$  and  $w_h \in V_h$  we have

$$\begin{aligned} \langle F(u_h), w \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} &= \langle F(u_h), w - w_h \rangle_{W^{-1,p}(\Omega)W_0^{1,q}(\Omega)} \\ &= \sum_{T \in \mathcal{T}_h} \left\{ \int_T k(u_h)\mathbf{grad}u_h \mathbf{grad}(w - w_h) dx + \int_T \mathbf{c} \mathbf{grad}u_h (w - w_h) dx \right. \\ & \quad \left. - \int_T f(w - w_h) dx \right\} \\ &= \sum_{T \in \mathcal{T}_h} \left\{ \int_T [-\operatorname{div}(k(u_h)\mathbf{grad}u_h) + \mathbf{c} \mathbf{grad}u_h - f](w - w_h) dx \right. \\ & \quad \left. + \int_{\partial T} k(u_h) \frac{\partial u_h}{\partial n_T} (w - w_h) ds \right\}, \end{aligned}$$

where  $\frac{\partial u_h}{\partial n_T}$  is the exterior normal derivative of  $u_h$  on the boundary  $\partial T$  of  $T$ . Using the Hölder inequality we get

$$(10.15) \quad \langle F(u_h), w \rangle_{W^{-1,p}(\Omega) W_0^{1,q}(\Omega)} \leq \sum_{T \in \mathcal{T}_h} \left\{ \left\| -\operatorname{div}(k(u_h) \mathbf{grad} u_h) + \mathbf{c} \mathbf{grad} u_h - f \right\|_{0,p,T} \|w - w_h\|_{0,q,T} + \sum_{i=1}^3 \left\| \left[ k(u_h) \frac{\partial u_h}{\partial n} \right] \right\|_{0,p,t_i} \|w - w_h\|_{0,q,t_i} \right\}$$

with  $\left[ \frac{\partial u_h}{\partial n} \right]$  the jump of a normal derivative of  $u_h$  on the side  $t_i$  of  $T$  with the convention  $u_h = 0$  outside  $\bar{\Omega}$ .

We choose in the inequality (10.15)  $w_h = \tilde{r}_h w$  where  $\tilde{r}_h$  is the Clément interpolation operator, see CLÉMENT [1975]. Let  $S_T = \{T' \in \mathcal{T}_h; T' \cap T \neq \emptyset\}$ . Then for all functions  $v \in W^{1,q}(\cup_{T' \in S_T} T')$  the following interpolation results hold

$$\|v - \tilde{r}_h v\|_{0,q,T} \leq Ch_T \sum_{T' \in S_T} \|v\|_{1,q,T'}$$

and

$$\|v - \tilde{r}_h v\|_{1,q,T} \leq C \sum_{T' \in S_T} \|v\|_{1,q,T'}.$$

Using the technique of the reference element and the above interpolation result, it is a simple matter to prove there exists a constant  $C$  such that for all  $v \in W^{1,q}(\cup_{T' \in S_T} T')$

$$\|v - \tilde{r}_h v\|_{0,q,t_i} \leq Ch_{t_i}^{1/q} h_T^{1-2/q} \left( \sum_{T' \in S_T} \|v\|_{1,q,T'} \right).$$

It follows that

$$\begin{aligned} \langle F(u_h), w \rangle_{W^{-1,p}(\Omega) W_0^{1,q}(\Omega)} &\leq C \sum_{T \in \mathcal{T}_h} \left\{ h_T \left\| -\operatorname{div}(k(u_h) \mathbf{grad} u_h) + \mathbf{c} \mathbf{grad} u_h - f \right\|_{0,p,T} \right. \\ &\quad \left. + \sum_{i=1}^3 h_{t_i}^{1/q} h_T^{1-2/q} \left\| \left[ k(u_h) \frac{\partial u_h}{\partial n} \right] \right\|_{0,p,t_i} \right\}^p |w|_{1,q,\Omega}, \end{aligned}$$

which completes the proof.  $\square$

*Remark 10.2.* Actually problem (10.2) is not solved as it is. The integrals are not computed exactly but with numerical quadrature rules. For the sake of simplicity we have analyzed an approximate problem without numerical integration, so that we could apply Theorem 7.1.  $\square$

*Remark 10.3.* We can obtain a priori and a posteriori error estimates in the norm  $|\cdot|_{1,\Omega}$  similar to the estimates (10.11) and (10.12) with the method developed in Theorem 6.4. Here we have  $X = W_0^{1,p}(\Omega)$  with  $p > 2$ ,  $Z = W^{-1,p}(\Omega)$ . Then let  $\mathcal{X} = H_0^1(\Omega)$  and  $\mathcal{Z} = H^{-1}(\Omega)$ . If  $q$  denotes the conjugate of  $p$ , the following inclusions hold

$$X \subset \mathcal{X} \subset W_0^{1,q}(\Omega) \subset L^2(\Omega) \subset Z \subset \mathcal{Z}.$$

Considering the expression (5.12) for  $DF(w)$  with  $w \in W_0^{1,p}(\Omega)$ ,  $2 < p < \infty$ , we easily see that  $DF(w) \in \mathcal{L}(W_0^{1,p}(\Omega); W^{-1,p}(\Omega))$  admits a continuous extension  $DF(w) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  and the mapping  $w \in W_0^{1,p}(\Omega) \rightarrow DF(w) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is continuous.

Let us check now that the mapping  $DF(u) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is an isomorphism. We can check that the operator  $TDF(u)R$  defined in the proof of Theorem 10.1, see (10.6), is a Fredholm operator with index zero from  $H_0^1(\Omega)$  into itself. The density of  $W_0^{1,p}(\Omega)$  into  $H_0^1(\Omega)$  implies that the range of  $TDF(u)R$  considered in  $H_0^1(\Omega)$  is  $H_0^1(\Omega)$  itself. It follows that  $TDF(u)R$  is an isomorphism on  $H_0^1(\Omega)$ , and since  $T \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$  is an isomorphism,  $DF(u) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is also an isomorphism.

As a consequence of (10.11), we have

$$\lim_{h \rightarrow 0} |u - u_h|_{1,p,\Omega} = 0.$$

We proceed as in the proof of Theorem 6.4 to get if  $h$  is small enough, say  $h \leq h_0$ ,

$$|u - u_h|_{1,\Omega} \leq 2 \|DF(u)^{-1}\|_{H^{-1}(\Omega); H_0^1(\Omega)} \|F(u_h)\|_{-1,\Omega}.$$

Such a posteriori estimates have been used in adaptive mesh refinement techniques applied to a two-dimensional regularized Stefan problem, see PICASSO [1992], and yield excellent results.

To derive a priori error estimates in the norm  $|\cdot|_{1,\Omega}$ , we proceed in the same way as above but with the mapping  $F_h$  defined in (7.11).  $\square$



**PARAMETRIZED NONLINEAR PROBLEMS.  
APPROXIMATION OF REGULAR SOLUTIONS**

Generally computational applications contain parameters. In Section 3 we have presented a parameter dependent problem, see problem (3.18), and we have proved the existence of a solution path. Our goal now is to study approximations of the problem

$$F(\lambda, u) = 0,$$

where  $X$  and  $Z$  are Banach spaces and  $F$  is a mapping defined on  $\mathbb{R}^m \times X$  with values in  $Z$ . For a given value of the parameter  $\lambda$  the study can be reduced to the one of Sections 6 and 7, where no parameter does occur. But then we can study how the approximation varies with respect to the parameter, which is precisely the object of this chapter, when the solutions are regular.

First we introduce the notions of coordinate spaces and state spaces. Then the case where the parameter space  $\mathbb{R}^m$  can be used as the coordinate space is analyzed both for general and Galerkin approximations. The results are applied to the example of Section 3. When the parameter space can no more be used as the coordinate space for the regular solution path, we transform the problem by adding equations. That situation is also analyzed and our results are illustrated with the example developed in Section 3.

### 11. Orientation

To present parameter dependent problems, it is adequate to work in the following general framework. Let  $X$  and  $Z$  be two real Banach spaces and  $F : X \rightarrow Z$  be a  $C^p$  mapping with  $p \geq 1$ . A point  $x \in X$  is called regular if

$$(11.1) \quad \begin{array}{ll} \text{(i)} & DF(x) \in \mathcal{L}(X; Z) \quad \text{has a finite-dimensional kernel,} \\ \text{(ii)} & DF(x) \in \mathcal{L}(X; Z) \quad \text{is surjective.} \end{array}$$

In this chapter we are interested to analyze and compute the solution set  $S$  of the problem

$$(11.2) \quad F(x) = 0$$

in a neighborhood of the regular point  $x_0 \in X$  satisfying  $F(x_0) = 0$ . Such a point  $x_0$  is called a regular solution. Here we assume that  $\dim \text{Ker}(DF(x_0)) \geq 1$ . The case  $\dim \text{Ker}(DF(x_0)) = 0$  has been studied in the previous chapter.

Since the dimension of  $\text{Ker}(DF(x_0))$  is finite this subspace is direct, that is there exists a closed subspace  $X_1 \subset X$  such that  $X = \text{Ker}(DF(x_0)) \oplus X_1$ . The point  $x_0$  is regular, so the mapping  $DF(x_0)|_{X_1} : X_1 \rightarrow Z$  is an isomorphism, the inverse of which is called  $A \in \mathcal{L}(Z; X_1)$ . A precise description of the solution set  $S$  in a neighborhood of  $x_0$  is the following.

**Theorem 11.1.** *There exist a positive number  $\epsilon$ , a neighborhood  $U$  of  $x_0 \in X$ , and a unique function  $g : B(0, \epsilon) \subset \text{Ker}(DF(x_0)) \rightarrow Z$  with*

$$g(0) = 0,$$

$$S \cap U = \{x(\lambda) = x_0 + \lambda + Ag(\lambda) \quad \text{for all } \lambda \in B(0, \epsilon)\}.$$

Moreover the mapping  $x : \lambda \in B(0, \epsilon) \rightarrow x(\lambda) = x_0 + \lambda + Ag(\lambda) \in S \cap U$  is a  $C^p$  diffeomorphism.

*Proof.* The affine mapping  $H : \text{Ker}(DF(x_0)) \times Z \rightarrow X$  is given by

$$H(\lambda, z) = x_0 + \lambda + Az.$$

Then we define the mapping  $G : \text{Ker}(DF(x_0)) \times Z \rightarrow Z$  by

$$G(\lambda, z) = F(H(\lambda, z)),$$

which is of class  $C^p$ . We notice that  $G(0, 0) = 0$  and that the derivative  $D_z G(0, 0) = DF(x_0)A \in \mathcal{L}(Z; Z)$  is an isomorphism. We apply the implicit function theorem 2.3. There exist positive constants  $\epsilon_1, \eta_1$ , and a unique mapping  $g : B(0, \epsilon_1) \subset \text{Ker}(DF(x_0)) \rightarrow B(0, \eta_1) \subset Z$ , which is of class  $C^p$  and such that

$$g(0) = 0,$$

$$\text{for all } \lambda \in B(0, \epsilon_1) \quad F(x_0 + \lambda + Ag(\lambda)) = 0.$$

We set  $U_1 = H(B(0, \epsilon_1) \times B(0, \eta_1))$  and the uniqueness implies

$$S \cap U_1 = \{x = x_0 + \lambda + Ag(\lambda), \quad \text{for all } \lambda \in B(0, \epsilon_1)\}.$$

To prove that the mapping  $x(\lambda) = x_0 + \lambda + Ag(\lambda)$  is a  $C^p$  diffeomorphism, we remark that for all  $\lambda \in B(0, \epsilon_1)$ ,  $F(x(\lambda)) = 0$  and  $DF(x_0)Dx(0) = 0$ . So  $Dx(0) \in \mathcal{L}(\text{Ker}(DF(x_0)); X)$  is an isomorphism onto  $\text{Ker}(DF(x_0))$  and we define the mapping  $K : B(0, \epsilon_1) \times X_1 \rightarrow X$  by

$$K(\lambda, x_1) = x(\lambda) + x_1.$$

Then  $K$  is of class  $C^p$  and  $DK(0, 0) \in \mathcal{L}(\text{Ker}(DF(x_0)) \times X_1, X)$  is an isomorphism. From the inverse function theorem 2.2, we deduce that  $K$  is a  $C^p$  diffeomorphism from  $B((0, 0), \epsilon) \subset \text{Ker}(DF(x_0)) \times X_1$  onto a neighborhood  $U$  of  $x_0$  with  $\epsilon \leq \epsilon_1$  and  $U \subset U_1$ . We can conclude since  $K(\lambda, 0) = x(\lambda)$ .  $\square$

With the decomposition  $X = \text{Ker}(DF(x_0)) \oplus X_1$ , we get a representation of the solution set close to  $x_0$  in terms of  $\lambda \in B(0, \epsilon) \subset \text{Ker}(DF(x_0))$ . So we say that  $\text{Ker}(DF(x_0))$  is a coordinate space while  $X_1$  is said to be a state space. Then Theorem 11.1 means that if  $m = \dim \text{Ker}(DF(x_0))$ , the mapping  $x(\lambda)$  defines a local coordinate system of the  $m$ -dimensional  $C^p$  manifold  $S$  in a neighborhood of  $x_0$ . Clearly the choice of a coordinate space is not unique. A suitable decomposition of  $X$  could be any pair of closed subspaces  $T$  and  $X_1$  such that

$$(11.3) \quad X = T \oplus X_1, \quad \dim T = m, \quad X_1 \cap \text{Ker}(DF(x_0)) = \{0\}.$$

Further developments of this analysis based on modern differential geometry can be found in FINK and RHEINBOLDT [1986], [1987], and in RHEINBOLDT [1986].

We have already noticed in the examples of Section 3 that some quantities are naturally identified as parameters. For example the space  $X$  could have the form  $X = \mathbb{R}^m \times X_1$  with  $X_1$  a Banach space and the identified parameters in  $\mathbb{R}^m$ .

In this chapter we shall focus on the following situation. Given a mapping  $F : \mathbb{R}^m \times X \rightarrow Z$  of class  $C^p$  and a regular solution  $(\lambda_0, u_0) \in \mathbb{R}^m \times X$  in the solution set  $S = \{(\lambda, u); F(\lambda, u) = 0\}$ , we study approximations of  $S$  in a neighborhood of  $(\lambda_0, u_0)$ . First we shall analyze the case when the parameter space  $\mathbb{R}^m$ , (in fact  $\mathbb{R}^m \times \{0\}$ ), can be used as the coordinate space, see (11.3). In that situation the solution set can be parametrized by  $\lambda \in \mathbb{R}^m$  and since  $(\lambda_0, u_0)$  is regular the operator  $D_u F(\lambda_0, u_0) \in \mathcal{L}(X; Z)$  is an isomorphism. The second case we shall analyze is when  $D_u F(\lambda_0, u_0)$  is no more an isomorphism from  $X$  onto  $Z$ . Then it is no more possible to describe the solution set of  $F(\lambda, u) = 0$  in a neighborhood of  $(\lambda_0, u_0)$  as a regular function of the parameter  $\lambda \in \mathbb{R}^m$ . Then one needs either to choose a coordinate space (different from  $\mathbb{R}^m$ ) satisfying (11.3) or transform the problem. We decide for the second method. By adding equations we transform the problem to apply the theory developed in the first case. Some further developments on local coordinate systems of  $S$  can be found in RHEINBOLDT [1986], Chapter 4.

## 12. Approximation of a parametrized solution set

Let  $m$  be a positive integer,  $X$  and  $Z$  be two real Banach spaces. Given a  $C^1$  mapping  $F : \mathbb{R}^m \times X \rightarrow Z$ , we consider the problem to find  $(\lambda, x) \in \mathbb{R}^m \times X$  such that

$$(12.1) \quad F(\lambda, x) = 0.$$

We assume that

$$(12.2) \quad (\lambda_0, x_0) \text{ is a solution to (12.1) and } D_x F(\lambda_0, x_0) \in \mathcal{L}(X; Z) \text{ is an isomorphism;}$$

so  $(\lambda_0, x_0)$  is a regular solution to (12.1). We can apply the implicit function theorem and get the following result. There exist a neighborhood  $I$  of  $\lambda_0$  in  $\mathbb{R}^m$ , a neighborhood  $U$  of  $x_0$  in  $X$ , and a  $C^1$  mapping  $u : I \rightarrow U$  such that  $\{(\lambda, u(\lambda)); \lambda \in I\}$  is a regular

solution set containing  $(\lambda_0, x_0)$ . Moreover if  $\{(\lambda, u(\lambda)); \lambda \in I\}$  and  $\{(\lambda, \tilde{u}(\lambda)); \lambda \in \tilde{I}\}$  are two regular solution sets containing  $(\lambda_0, x_0)$ , then  $u(\lambda)$  and  $\tilde{u}(\lambda)$  coincide on the connected component of  $I \cap \tilde{I}$  containing  $\lambda_0$ .

Note that in the case  $m = 1$ , the terminology solution branch is often used.

We introduce now a family  $\{F_h\}_{0 < h \leq 1}$  of  $C^1$  mappings

$$F_h : (\lambda, x) \in \mathbb{R}^m \times X \rightarrow F_h(\lambda, x) \in Z,$$

which are approximations of  $F$ . Our goal is to study the existence and the convergence of solutions of the problem: find  $(\lambda, x_h) \in \mathbb{R}^m \times X$  such that

$$(12.3) \quad F_h(\lambda, x_h) = 0$$

in a neighborhood of the regular solution set  $\{(\lambda, u(\lambda)); \lambda \in I\}$ .

For a fixed  $\lambda$ , we consider a solution  $(\lambda, u) \in \mathbb{R} \times X$  of (12.1) and associate an element  $\tilde{u}_h \in X$ ,  $h \in (0, 1]$ . In some cases it will be sufficient to choose  $\tilde{u}_h = u$  but in other ones an appropriate choice of  $\tilde{u}_h$  will lead to different error estimates.

Now for that fixed  $\lambda$ , the theory developed in Section 6 can be applied when  $\tilde{u}_h = u$ . Here we will work under analogous hypotheses of equicontinuity, consistency and stability. We assume that the mapping  $D_x F_h(\lambda, \cdot)$  is Lipschitzian at  $\tilde{u}_h$  with constant  $L(= L(\lambda))$ , that is there exist  $\eta > 0$  and  $L$  such that for all  $v \in \overline{B}(\tilde{u}_h, \eta)$ , for all  $h \in (0, 1]$ ,

$$(12.4) \quad \|D_x F_h(\lambda, \tilde{u}_h) - D_x F_h(\lambda, v)\|_{X;Z} \leq L \|\tilde{u}_h - v\|_X.$$

Moreover we assume that

$$(12.5) \quad \lim_{h \rightarrow 0} \|F_h(\lambda, \tilde{u}_h)\|_Z = 0$$

and for all  $h \in (0, 1]$ ,  $D_x F_h(\lambda, \tilde{u}_h)$  is an isomorphism from  $X$  onto  $Z$  with an inverse uniformly bounded in  $h$ , that is there exists  $C_1(= C_1(\lambda))$  such that for all  $h \in (0, 1]$

$$(12.6) \quad \|D_x F_h(\lambda, \tilde{u}_h)^{-1}\|_{Z;X} \leq C_1.$$

**Theorem 12.1.** *Let  $\{F_h\}_{0 < h \leq 1}$  be a family of  $C^1$  mappings,  $F_h : \mathbb{R}^m \times X \rightarrow Z$ . For a given  $\lambda \in \mathbb{R}^m$ , we associate an element  $\tilde{u}_h \in X$ ,  $h \in (0, 1]$ . We assume that the hypotheses (12.4), (12.5), and (12.6) hold. Then there exist  $h_0 > 0$ ,  $\delta_0 > 0$ , and for all  $0 < h \leq h_0$  a unique  $u_h \in X$  satisfying*

$$F_h(\lambda, u_h) = 0 \quad \text{and} \quad \|u_h - \tilde{u}_h\|_X \leq \delta_0.$$

Moreover the mapping  $D_x F_h(\lambda, u_h) \in \mathcal{L}(X; Z)$  is invertible with the bound

$$\|D_x F_h(\lambda, u_h)^{-1}\|_{Z;X} \leq 2 \|D_x F_h(\lambda, \tilde{u}_h)^{-1}\|_{Z;X},$$

and

$$(12.7) \quad \|\tilde{u}_h - u_h\|_X \leq 2\|D_x F_h(\lambda, \tilde{u}_h)^{-1}\|_{Z;X} \|F_h(\lambda, \tilde{u}_h)\|_Z.$$

*Proof.* The proof goes along the same lines as the one of Theorem 6.1. For the given  $\lambda \in I$ , we set

$$\begin{aligned} \epsilon_h &= \|F_h(\lambda, \tilde{u}_h)\|_Z, \\ \gamma_h &= \|D_x F_h(\lambda, \tilde{u}_h)^{-1}\|_{Z;X}, \\ \tilde{L}_h(\alpha) &= \sup_{v \in \overline{B}(\tilde{u}_h, \alpha)} \|D_x F_h(\lambda, \tilde{u}_h) - D_x F_h(\lambda, v)\|_{X;Z}. \end{aligned}$$

From the stability hypothesis (12.6), we get  $\gamma_h \leq C_1$ . We set  $\delta_0 = \min(\eta, (2C_1 L)^{-1})$ . Then from the Lipschitzian condition (12.4) we have

$$\tilde{L}_h(\delta_0) \leq (2C_1)^{-1}.$$

Now from the consistency (12.5), there exists  $h_0 \leq 1$  such that for all  $h \in (0, h_0]$

$$2C_1 \epsilon_h \leq \delta_0.$$

Thus for  $h \leq h_0$

$$2\gamma_h \tilde{L}_h(2\gamma_h \epsilon_h) \leq 2C_1 \tilde{L}_h(\delta_0) \leq 1.$$

We apply Theorem 2.1 with  $G = F_h(\lambda, \cdot)$ ,  $v = \tilde{u}_h \in X$ . For all  $h \leq h_0$ , there exists a unique  $u_h \in \overline{B}(\tilde{u}_h, 2\gamma_h \epsilon_h)$  satisfying  $F_h(\lambda, u_h) = 0$ , and  $D_x F_h(\lambda, u_h) \in \mathcal{L}(X; Z)$  is invertible with the bound

$$\|D_x F_h(\lambda, u_h)^{-1}\|_{Z;X} \leq 2\gamma_h.$$

The following estimate holds

$$\|\tilde{u}_h - u_h\|_X \leq 2\gamma_h \|F_h(\lambda, \tilde{u}_h)\|_Z.$$

To get the uniqueness of  $u_h$  in  $\overline{B}(\tilde{u}_h, \delta_0)$  for all  $h \in (0, h_0]$ , we use the argument developed in Remark 2.1.  $\square$

If we choose in Theorem 12.1,  $\tilde{u}_h = u$  where  $(\lambda, u)$  is a solution of (12.1), then the inequality (12.7) will lead to a priori error estimates for  $\|u - u_h\|_X$ . To get a posteriori error estimates, we use the following result.

**Theorem 12.2.** *Let  $(\lambda, u)$  be a solution of problem (12.1). We assume that the family  $\{F_h\}_{0 < h \leq 1}$  of  $C^1$  mappings satisfies the hypotheses (12.4), (12.5), and (12.6) with  $\tilde{u}_h = u$ . Then there exists  $0 < \bar{h}_0 \leq h_0$  such that for all  $h \leq \bar{h}_0$  the mapping  $D_x F(\lambda, u_h) \in \mathcal{L}(X; Z)$  is invertible with a uniformly bounded inverse and*

$$(12.8) \quad \|u - u_h\|_X \leq 2\|D_x F(\lambda, u_h)^{-1}\|_{Z;X}\|F(\lambda, u_h)\|_Z,$$

where  $h_0$  and  $u_h$  are given in Theorem 12.1.

*Proof.* The proof is identical to the one of Theorem 6.2. It is an application of Theorem 2.1 with  $G = F(\lambda, \cdot)$  and  $v = u$ .  $\square$

Actually in the theorems 12.1 and 12.2 we have worked with the parameter  $\lambda \in I$  fixed. We study now the dependence in  $\lambda$  of  $u_h$ .

The mappings  $D_x F_h(\lambda, \cdot)$  are assumed to be Lipschitzian at  $\tilde{u}_h(\lambda)$ , that is there exist  $\eta > 0$ ,  $L > 0$  such that for all  $0 < h \leq 1$ ,  $\lambda \in I$ , and  $v \in \overline{B}(\tilde{u}_h(\lambda), \eta)$

$$(12.9) \quad \|D_x F_h(\lambda, \tilde{u}_h(\lambda)) - D_x F_h(\lambda, v)\|_{X;Z} \leq L\|\tilde{u}_h(\lambda) - v\|_X;$$

moreover we assume the consistency condition

$$(12.10) \quad \lim_{h \rightarrow 0} \sup_{\lambda \in I} \|F_h(\lambda, \tilde{u}_h(\lambda))\|_Z = 0$$

and the stability condition: there exists a constant  $C$  such that for all  $0 < h \leq 1$ ,  $\lambda \in I$ ,

$$(12.11) \quad \|D_x F_h(\lambda, \tilde{u}_h(\lambda))^{-1}\|_{Z;X} \leq C.$$

**Theorem 12.3.** *Let  $\{(\lambda, u(\lambda)); \lambda \in I\}$  be a set of regular solutions of the problem (12.1). We assume that  $I$  is compact, the mappings  $\tilde{u}_h : \lambda \in I \rightarrow \tilde{u}_h(\lambda) \in X$  are continuous, and the hypotheses (12.9), (12.10), (12.11) hold. Then there exist  $h_0 > 0$ ,  $\delta_0 > 0$ , and for all  $h \leq h_0$  a  $C^1$  mapping  $u_h : \lambda \in I \rightarrow u_h(\lambda) \in X$  such that*

$$(F_h(\lambda, v) = 0 \quad \text{and} \quad v \in \overline{B}(\tilde{u}_h(\lambda), \delta_0)) \Leftrightarrow v = u_h(\lambda).$$

Moreover the mappings  $D_x F_h(\lambda, \tilde{u}_h(\lambda))$  and  $D_x F(\lambda, u_h(\lambda)) \in \mathcal{L}(X; Z)$  are isomorphisms with uniformly bounded inverse and for all  $\lambda \in I$

$$(12.12) \quad \|u_h(\lambda) - \tilde{u}_h(\lambda)\|_X \leq 2\|D_x F_h(\lambda, \tilde{u}_h(\lambda))^{-1}\|_{Z;X}\|F_h(\lambda, \tilde{u}_h(\lambda))\|_Z,$$

$$(12.13) \quad \|u_h(\lambda) - u(\lambda)\|_X \leq 2\|D_x F(\lambda, u_h(\lambda))^{-1}\|_{Z;X}\|F(\lambda, u_h(\lambda))\|_Z.$$

*Proof.* For each  $\lambda \in I$  we can apply Theorem 12.1. So for each  $\lambda \in I$ , there exist  $\delta_{0,\lambda}$ ,  $h_{0,\lambda}$ , and  $u_{h,\lambda} \in X$  for  $h \in (0, h_{0,\lambda}]$  satisfying to

$$F_h(\lambda, u_{h,\lambda}) = 0 \quad \text{and} \quad \|u_{h,\lambda} - \tilde{u}_h(\lambda)\|_X \leq \delta_{0,\lambda}$$

and

$$\|\tilde{u}_h(\lambda) - u_{h,\lambda}\|_X \leq 2\|D_x F_h(\lambda, \tilde{u}_h(\lambda))^{-1}\|_{Z;X} \|F_h(\lambda, \tilde{u}_h(\lambda))\|_Z.$$

Furthermore thanks to the uniformity in  $\lambda$  of the assumptions (12.9), (12.10), and (12.11), it is immediate to check that  $\delta_{0,\lambda}$  and  $h_{0,\lambda}$  can be chosen independently of  $\lambda \in I$ . So for  $h \leq h_0$ , we can define the function  $u_h : \lambda \in I \rightarrow u_h(\lambda) = u_{h,\lambda} \in X$ .

To prove that the function  $u_h(\cdot)$  is  $C^1$  in  $I$ , we apply the implicit function theorem 2.3 taking into account the uniqueness of  $u_h(\lambda)$  in  $\overline{B}(\tilde{u}_h(\lambda), \delta_0)$ .

The estimate (12.13) is a consequence of Theorem 12.2 only to restrict  $h_0$  if necessary.  $\square$

*Remark 12.1.* The assumption (12.11) can be replaced by

$$(12.11') \quad \lim_{h \rightarrow 0} \sup_{\lambda \in I} \|D_x F_h(\lambda, \tilde{u}_h(\lambda)) - D_x F(\lambda, u(\lambda))\|_{X;Z} = 0.$$

From the equality

$$D_x F_h(\lambda, \tilde{u}_h(\lambda)) = [D_x F_h(\lambda, \tilde{u}_h(\lambda)) - D_x F(\lambda, u(\lambda))] + D_x F(\lambda, u(\lambda))$$

and the estimate (12.11'), we deduce from (1.1) that for  $h$  small enough the mapping  $D_x F_h(\lambda, \tilde{u}_h(\lambda)) \in \mathcal{L}(X; Z)$  is an isomorphism with an inverse uniformly bounded in  $h$  and  $\lambda$ .  $\square$

If we assume that the mappings  $F_h$  are of class  $C^p$  with  $p \geq 1$ , then  $u_h(\cdot)$  is also of class  $C^p$ . Then error estimates for the derivatives can also be deduced.

**Theorem 12.4.** *If the assumptions of Theorem 12.3 hold and the mappings  $F_h$  are of class  $C^p$  with  $p \geq 1$ , then  $u_h(\cdot)$  is of class  $C^p$  in  $I$ . Moreover if the functions  $\tilde{u}_h(\cdot)$  are of class  $C^p$  then there exist constants  $C_j$ ,  $j = 0, \dots, p$ , depending only on the quantities*

$$\max_{\lambda \in I} \max_{v \in \overline{B}(\tilde{u}_h(\lambda), \delta_0)} \|D^k F_h(\lambda, v)\| \quad \text{and} \quad \max_{\lambda \in I} \|D^k \tilde{u}_h(\lambda)\|, \quad 0 \leq k \leq j$$

such that for all  $\lambda \in I$

$$(12.14) \quad \|D^j u_h(\lambda) - D^j \tilde{u}_h(\lambda)\| \leq C_j \sum_{k=0}^j \left\| \frac{d^k}{d\lambda^k} F_h(\lambda, \tilde{u}_h(\lambda)) \right\|$$

for  $j = 0, \dots, p-1$  and

$$(12.15) \quad \|D^p u_h(\lambda)\| \leq C_p.$$

For simplicity  $\|\cdot\|$  denotes the norm in the different spaces.

*Proof.* From the implicit function theorem we deduce that  $u_h$  is of class  $C^p$  in  $I$ . We have to prove yet the estimates (12.14) and (12.15). The differentiation of the relation  $F_h(\lambda, u_h(\lambda)) = 0$  for  $\lambda \in I$  gives

$$D u_h(\lambda) = -D_x F_h(\lambda, u_h(\lambda))^{-1} D_\lambda F_h(\lambda, u_h(\lambda)),$$

which leads to the estimate

$$\|Du_h(\lambda)\| \leq 2\|D_x F_h(\lambda, \tilde{u}_h(\lambda))^{-1}\| \|D_\lambda F_h(\lambda, u_h(\lambda))\|;$$

we notice that the estimate for  $\|D_x F_h(\lambda, u_h(\lambda))^{-1}\|$  is given in Theorem 12.1. This last estimate and the (12.12) one give the estimates (12.14) and (12.15) with  $p = 1$ .

Assuming the result holds for  $\ell$ ,  $1 \leq \ell < p$ , that is

$$(12.16) \quad \|D^j u_h(\lambda) - D^j \tilde{u}_h(\lambda)\| \leq C_j \sum_{k=0}^j \left\| \frac{d^k}{d\lambda^k} F_h(\lambda, \tilde{u}_h(\lambda)) \right\|$$

for  $j = 0, \dots, \ell - 1$  and

$$(12.17) \quad \|D^\ell u_h(\lambda)\| \leq C'_\ell,$$

we prove it for the index value  $\ell + 1$ . We differentiate  $\ell$  times with respect to  $\lambda$  the function  $F_h(\lambda, \tilde{u}_h(\lambda))$ , so

$$(12.18) \quad \frac{d^\ell}{d\lambda^\ell} F_h(\lambda, \tilde{u}_h(\lambda)) = D_x F_h(\lambda, \tilde{u}_h(\lambda)) D^\ell \tilde{u}_h(\lambda) + R_\ell(\lambda, \tilde{u}_h(\lambda), \dots, D^{\ell-1} \tilde{u}_h(\lambda))$$

where  $R_\ell$  has the form

$$R_1(\lambda, v) = D_\lambda F_h(\lambda, v),$$

$$R_\ell(\lambda, v, \dots, D^{\ell-1} v) = \sum_{i=2}^{\ell} \sum_{\substack{k_1+k_2+\dots+k_i=\ell \\ k_1, \dots, k_i \geq 1}} \alpha_{ik} D^i F_h(\lambda, v) \left( (\lambda, v)^{(k_1)}, \dots, (\lambda, v)^{(k_i)} \right)$$

with

$$\alpha_{ik} \in \mathbb{R} \quad \text{and} \quad (\lambda, v)^{(k_j)} = (D^{k_j} \lambda, D^{k_j} v).$$

We derive now  $\ell$  times with respect to  $\lambda$  the relation  $F_h(\lambda, u_h(\lambda)) = 0$  so

$$(12.19) \quad 0 = D_x F_h(\lambda, u_h(\lambda)) D^\ell u_h(\lambda) + R_\ell(\lambda, u_h(\lambda), \dots, D^{\ell-1} u_h(\lambda)).$$

From both expressions (12.18) and (12.19), we get

$$(12.20) \quad D^\ell \tilde{u}_h(\lambda) - D^\ell u_h(\lambda) = D_x F_h(\lambda, u_h(\lambda))^{-1} \left\{ \frac{d^\ell}{d\lambda^\ell} F_h(\lambda, \tilde{u}_h(\lambda)) \right. \\ \left. + [D_x F_h(\lambda, u_h(\lambda)) - D_x F_h(\lambda, \tilde{u}_h(\lambda))] D^\ell \tilde{u}_h(\lambda) - R_\ell(\lambda, \tilde{u}_h(\lambda), \dots, D^{\ell-1} \tilde{u}_h(\lambda)) \right. \\ \left. + R_\ell(\lambda, u_h(\lambda), \dots, D^{\ell-1} u_h(\lambda)) \right\}.$$

The mapping  $R_\ell$  is Lipschitzian, so we deduce from (12.20) and (12.16) the estimate (12.14) for  $\ell$ .

The estimate (12.15) for the index value  $\ell + 1$  is immediately deduced from the formula

$$D^{\ell+1} u_h(\lambda) = -D_x F_h(\lambda, u_h(\lambda))^{-1} R_{\ell+1}(\lambda, u_h(\lambda), \dots, D^\ell u_h(\lambda)). \quad \square$$

*Remark 12.2.* The method presented in Theorem 6.4 to derive error estimates in some other norms than the norm of  $X$  can be immediately generalized to parameter dependent problems.  $\square$



### 13. Galerkin approximation of a parametrized solution set

The goal of the present section is to transcribe the results of the previous one in the particular case of Galerkin approximations. For simplicity we will consider problem (12.1) with  $m = 1$ .

Let  $X, Y$  be two real Banach spaces, and let  $Y'$  denote the dual space of  $Y$ . Given a  $C^1$  mapping  $F : \mathbb{R} \times X \rightarrow Y'$ , we shall analyze Galerkin approximations of the problem

$$(13.1) \quad F(\lambda, x) = 0$$

in a neighborhood of the regular solution  $(\lambda_0, x_0) \in \mathbb{R} \times X$ , such that  $(\lambda_0, x_0)$  satisfies

$$(13.2) \quad F(\lambda_0, x_0) = 0 \text{ and } D_x F(\lambda_0, x_0) \in \mathcal{L}(X; Y') \text{ is an isomorphism.}$$

As a consequence of the hypothesis (13.2), we get the existence of a regular solution branch  $\{(\lambda, u(\lambda)); \lambda \in I\}$  containing  $(\lambda_0, x_0)$ , where the mapping  $u : I \rightarrow X$  is of class  $C^1$ .

Let now  $\{X_h\}_{0 < h \leq 1}$  be a family of finite-dimensional subspaces of  $X$  and  $\{Y_h\}_{0 < h \leq 1}$  be a family of finite-dimensional subspaces of  $Y$ . A Galerkin approximation of the problem (13.1) consists in finding  $\lambda \in \mathbb{R}$  and  $x_h \in X_h$  such that

$$(13.3) \quad \text{for all } y_h \in Y_h \quad \langle F(\lambda, x_h), y_h \rangle_{Y'Y} = 0.$$

For each  $\lambda \in I$ , we define the bilinear form  $b : X \times Y \rightarrow \mathbb{R}$  by

$$(13.4) \quad \text{for all } x \in X, y \in Y \quad b(x, y) = \langle D_x F(\lambda, u(\lambda))x, y \rangle_{Y'Y}.$$

Although we do not write it explicitly for notation simplicity,  $b$  is depending on  $\lambda$ . We assume that there exists a constant  $\beta > 0$  such that for all  $h \in (0, 1]$  and  $\lambda \in I$

$$(13.5) \quad \inf_{\substack{x \in X_h \\ \|x\|_X = 1}} \sup_{\substack{y \in Y_h \\ \|y\|_Y = 1}} b(x, y) \geq \beta > 0$$

and

$$(13.6) \quad \dim X_h = \dim Y_h.$$

Moreover we assume that for all  $\lambda \in I$

$$(13.7) \quad \lim_{h \rightarrow 0} \inf_{x_h \in X_h} \|u(\lambda) - x_h\|_X = 0.$$

Note that the hypotheses (13.5) and (13.6) are standard stability conditions in the linear case while the (13.7) one is a standard approximation property of  $X$  by  $X_h$ .

The following theorem presents the convergence results in the case of Galerkin approximation.

**Theorem 13.1.** *Let  $F : \mathbb{R} \times X \rightarrow Y'$  be a  $C^1$  mapping and  $\{(\lambda, u(\lambda)); \lambda \in I \subset \mathbb{R}, I \text{ compact interval}\}$  be a regular solution branch. We assume that there exist  $L > 0$  and  $\eta > 0$  such that for all  $\lambda \in I$  and  $v \in \overline{B}(u(\lambda), \eta)$*

$$\|D_x F(\lambda, u(\lambda)) - D_x F(\lambda, v)\|_{X; Y'} \leq L \|u(\lambda) - v\|_X.$$

*Then under the assumptions (13.5), (13.6), and (13.7), there exist positive constants  $h_0, \delta_0$ , and for all  $\lambda \in I$  and  $h \in (0, h_0]$  a unique solution  $(\lambda, u_h(\lambda)) \in \mathbb{R} \times X_h$  to problem (13.3) in  $I \times \overline{B}(u(\lambda), \delta_0)$ . The mapping  $u_h : \lambda \in I \rightarrow X$  is of class  $C^1$  and for  $\lambda \in I$ ,  $D_x F(\lambda, u_h(\lambda)) \in \mathcal{L}(X, Y')$  is an isomorphism. Furthermore there exists a constant  $C > 0$  (independent of  $h$  and  $\lambda$ ) such that for all  $h \in (0, h_0]$  and  $\lambda \in I$*

$$(13.8) \quad \|u(\lambda) - u_h(\lambda)\|_X \leq C \inf_{x_h \in X_h} \|u(\lambda) - x_h\|_X$$

and

$$(13.9) \quad \|u(\lambda) - u_h(\lambda)\|_X \leq 2 \|D_x F(\lambda, u_h(\lambda))^{-1}\|_{Y'; X} \|F(\lambda, u_h(\lambda))\|_{Y'}.$$

*Proof.* The proof is analogous to the one of Theorem 7.1 and will be deduced from Theorems 12.1 and 12.2. To define the mapping  $F_h : I \times X \rightarrow Y'$ , it is useful to introduce for each  $\lambda \in I$  the two projection operators  $\Pi_{X_h} \in \mathcal{L}(X; X_h)$  and  $\Pi_{Y_h} \in \mathcal{L}(Y; Y_h)$  (depending on  $\lambda$ ) given by

$$(13.10) \quad \text{for all } y_h \in Y_h \quad b(x - \Pi_{X_h} x, y_h) = 0$$

and

$$(13.11) \quad \text{for all } x_h \in X_h \quad b(x_h, y - \Pi_{Y_h} y) = 0.$$

Under the assumptions (13.5) and (13.6) both operators  $\Pi_{X_h}$  and  $\Pi_{Y_h}$  are well defined. Then we construct the family  $\{F_h\}_{0 < h \leq 1}$  of mappings  $F_h : I \times X \rightarrow Y'$  by: for  $\lambda \in I$  and for  $x \in X, y \in Y$ ,

$$\langle F_h(\lambda, x), y \rangle_{Y'Y} = \langle F(\lambda, x), \Pi_{Y_h} y \rangle_{Y'Y} + b(x, y - \Pi_{Y_h} y).$$

Then with the arguments developed in the proof of Theorem 7.1 it is a simple matter to check that problem (13.3) with  $\lambda \in I$  is equivalent to find  $(\lambda, x_h) \in I \times X_h$  such that

$$(13.12) \quad F_h(\lambda, x_h) = 0.$$

In the same way as in the proof of Theorem 7.1, there exist constants  $h_0, \delta_0$  and for  $h \in (0, h_0]$  a unique solution  $u_h \in \overline{B}(u(\lambda), \delta_0)$  to  $F_h(\lambda, x_h) = 0$ . Moreover the estimates (13.8) and (13.9) are got from (12.12) and (12.13).

Only to restrict  $h_0$  if necessary, we can apply the implicit function theorem 2.3 and the above uniqueness result to prove that  $u_h : \lambda \in I \rightarrow u_h(\lambda) \in X$  is a  $C^1$  mapping.  $\square$

*Remark 13.1.* The duality argument used in Remark 7.3 is also valid here to get error estimates in some other norms than the norm of  $X$ . The generalization to the parameter dependent problem is immediate.  $\square$

*Remark 13.2.* If the mapping  $F$  is of class  $C^p$  with  $p \geq 1$ , then under the assumptions of Theorem 13.1, the conclusions of Theorem 12.4 are true with  $\tilde{u}_h(\lambda) = u(\lambda)$  of class  $C^p$ . More precisely the function  $u_h(\cdot)$  is of class  $C^p$  in  $I$  and there exist constants  $C_j$ ,  $j = 0, \dots, p$ , depending only on the quantities

$$\max_{\lambda \in I} \max_{v \in \overline{B}(u(\lambda), \delta_0)} \|D^k F_h(\lambda, v)\| \quad \text{and} \quad \max_{\lambda \in I} \|D^k u(\lambda)\|, \quad 0 \leq k \leq j$$

such that for all  $\lambda \in I$

$$(13.13) \quad \|D^j u_h(\lambda) - D^j u(\lambda)\| \leq C_j \sum_{k=0}^j \left\| \frac{d^k}{d\lambda^k} F_h(\lambda, u(\lambda)) \right\|$$

for  $j = 0, \dots, p-1$  and

$$(13.14) \quad \|D^p u_h(\lambda)\| \leq C_p.$$

For simplicity  $\|\cdot\|$  denotes the norm in the different spaces.  $\square$

#### 14. Error estimates for a parameter dependent semilinear problem

An application of the abstract results obtained in the two previous sections to the parameter dependent problems developed in Section 3 is presented here. First we shall analyze a simple case of the general problem (3.18).

Let  $\Omega$  be a polygonal convex open bounded set in  $\mathbb{R}^2$  and let  $f : H_0^1(\Omega) \rightarrow L^2(\Omega)$  be given by: for  $v \in H_0^1(\Omega)$

$$(14.1) \quad \text{for all } x \in \Omega \quad f(v)(x) = v^2(x) + 1.$$

The mapping  $F : \mathbb{R} \times H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  is defined by: for all  $(\lambda, v) \in \mathbb{R} \times H_0^1(\Omega)$ ,  $w \in H_0^1(\Omega)$

$$(14.2) \quad \langle F(\lambda, v), w \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w \, dx - \lambda \int_{\Omega} f(v) w \, dx.$$

Since the mapping  $f$  is  $C^\infty$ , so is  $F$ .

There exist a positive number  $\lambda^*$  and a  $C^0$  mapping  $u : [0, \lambda^*] \rightarrow u(\lambda) \in H_0^1(\Omega)$ , which is  $C^\infty$  on  $[0, \lambda^*)$ , of minimal positive solutions to the problem

$$(14.3) \quad F(\lambda, v) = 0,$$

see Remark 3.2. Moreover for each  $\lambda \in [0, \lambda^*)$ , the solution  $(\lambda, u(\lambda))$  to (14.3) is regular, in fact  $D_v F(\lambda, u(\lambda)) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  given by: for all  $w_1, w_2 \in H_0^1(\Omega)$

$$\langle D_v F(\lambda, u(\lambda))w_1, w_2 \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} = \int_{\Omega} \mathbf{grad} w_1 \mathbf{grad} w_2 dx - 2\lambda \int_{\Omega} u(\lambda) w_1 w_2 dx,$$

is an isomorphism. This means that the assumption (13.2) is satisfied for  $\lambda \in [0, \lambda^*)$ .

Given a regular family  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  of triangulations of  $\bar{\Omega}$ , we consider the corresponding family  $\{V_h\}_{0 < h \leq 1}$  of finite element subspaces of degree 1,

$$(14.4) \quad V_h = \{v \in C^0(\bar{\Omega}); v|_T \in \mathcal{P}_1 \text{ for all } T \in \mathcal{T}_h\} \cap H_0^1(\Omega).$$

A Galerkin approximation to the problem (14.3) consists in finding  $\lambda \in \mathbb{R}$  and  $v_h \in V_h$  such that

$$(14.5) \quad \text{for all } w_h \in V_h \quad \langle F(\lambda, v_h), w_h \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} = 0.$$

To study the discrete problem (14.5), we shall apply the results of Section 13, with  $X = Y = H_0^1(\Omega)$ ,  $F$  given in (14.2),  $X_h = Y_h = V_h$  and for  $\lambda \in I = [0, \bar{\lambda}]$  with  $0 < \bar{\lambda} < \lambda^*$ , the bilinear form  $b : H_0^1(\Omega) \times H_0^1(\Omega) \rightarrow \mathbb{R}$  given by

$$\text{for all } w_1, w_2 \in H_0^1(\Omega) \quad b(w_1, w_2) = \int_{\Omega} \mathbf{grad} w_1 \mathbf{grad} w_2 dx - 2\lambda \int_{\Omega} u(\lambda) w_1 w_2 dx.$$

The set  $\{(\lambda, u(\lambda)); \lambda \in I\}$  is a regular solution branch. The mapping  $F$  being of class  $C^k$  for any  $k$ , the derivative  $D_v F(\lambda, \cdot)$  is clearly Lipschitzian at  $u(\lambda)$ . To check the inf-sup condition, we notice that for  $w_h \in V_h$

$$\begin{aligned} & \int_{\Omega} \mathbf{grad} w_h \mathbf{grad} w_h dx - 2\lambda \int_{\Omega} u(\lambda) w_h w_h dx \\ & \geq \left(1 - \frac{\lambda}{\lambda^*}\right) \int_{\Omega} |\mathbf{grad} w_h|^2 dx + \frac{\lambda}{\lambda^*} \left[ \int_{\Omega} |\mathbf{grad} w_h|^2 dx - 2\lambda^* \int_{\Omega} u(\lambda) w_h^2 dx \right] \\ & \geq \left(1 - \frac{\lambda}{\lambda^*}\right) \int_{\Omega} |\mathbf{grad} w_h|^2 dx \end{aligned}$$

since the expression inside the brackets is positive, see the proof of Theorem 3.4. This implies the condition (13.5) with  $\beta = (\lambda^* - \bar{\lambda})/\lambda^*$ , while the (13.6) one is immediate with our choice of subspaces  $X_h = Y_h = V_h$ . Finally the assumption (13.7) is a standard polynomial interpolation result. We can apply Theorem 13.1 and get the following result.

**Theorem 14.1.** *Let  $F : \mathbb{R} \times H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  be the mapping defined in (14.2). The set  $\{(\lambda, u(\lambda)); \lambda \in I = [0, \bar{\lambda}]\}$  with  $\bar{\lambda} < \lambda^*$  denotes a set of minimal positive regular solutions to the problem (14.3). Let the finite element subspaces  $V_h$  be given in (14.4).*

Then there exist positive constants  $h_0, \delta_0, C$ , and for all  $\lambda \in I, h \in (0, h_0]$ , a solution  $(\lambda, u_h(\lambda)) \in \mathbb{R} \times V_h$  to the Galerkin approximation (14.5); for  $\lambda \in I, u_h(\lambda)$  is unique in the closed ball  $\overline{B}(u(\lambda), \delta_0)$ . Moreover the mapping  $u_h : \lambda \in I \rightarrow H_0^1(\Omega)$  is of class  $C^1$  and the following error estimates hold: for all  $h \in (0, h_0]$  and  $\lambda \in I$

$$(14.6) \quad |u(\lambda) - u_h(\lambda)|_{1,\Omega} \leq Ch$$

and

$$(14.7) \quad |u(\lambda) - u_h(\lambda)|_{1,\Omega} \leq C \left( \sum_{T \in \mathcal{T}_h} \eta(T)^2 \right)^{1/2},$$

where the local estimator  $\eta(T)$  is

$$\eta^2(T) = h_T^2 \|\Delta u_h(\lambda) + \lambda (u_h^2(\lambda) + 1)\|_{0,T}^2 + \sum_{i=1}^3 h_{t_i} \left\| \left[ \frac{\partial u_h(\lambda)}{\partial n} \right] \right\|_{0,t_i}^2;$$

here  $h_T$  is the diameter of  $T$ ,  $t_i$   $i = 1, 2, 3$  are the edges of  $T$  of length  $h_{t_i}$ , and  $[\cdot]$  denotes the jump through an edge.

*Proof.* The hypotheses of Theorem 13.1 were checked before with  $\beta = 1 - \overline{\lambda}/\lambda^*$ . From the estimates (13.8) and (13.9), we can deduce the (14.6) and (14.7) ones in the same way as in the proof of Theorem 8.1. Notice that the constant  $C$  is depending on  $\beta$ .  $\square$

The remainder of the section is devoted to study finite element approximations of problem (3.18) including numerical quadrature rules. In the previous example we have applied the theory of Galerkin approximations, now we will apply the general results of Section 12.

Let  $\Omega$  be a regular open bounded set in  $\mathbb{R}^2$  and  $f$  be the Nemytskii operator of some function  $\phi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ , that is for a function  $u : \Omega \rightarrow \mathbb{R}$  we have

$$\text{for all } x \in \Omega \quad f(u)(x) = \phi(x, u(x)).$$

We assume that the function  $\phi$  satisfies

$$(14.8) \quad \begin{cases} \text{the function } \phi \text{ is twice continuously differentiable with respect to the} \\ \text{second argument uniformly in } \Omega \text{ and for all } z \in \mathbb{R}, \phi(\cdot, z), \frac{\partial \phi}{\partial z}(\cdot, z), \\ \text{and } \frac{\partial^2 \phi}{\partial z^2}(\cdot, z) \text{ are in } L^\infty(\Omega) \end{cases}$$

and the assumptions (3.16), (3.17), (3.23). Remark that the hypothesis (14.8) is stronger than the (3.15) one, so that  $f : H_0^1(\Omega) \cap C^0(\overline{\Omega}) \rightarrow L^2(\Omega)$  is a  $C^2$  mapping.

In Section 3, we have studied the positive solutions of

$$(14.9) \quad u - \lambda T f(u) = 0,$$

with  $-T$  the inverse of the Laplacian operator with homogeneous Dirichlet boundary conditions, see Theorem 3.4. There exist a  $\lambda^* > 0$  and a  $C^2$  mapping  $\lambda \in [0, \lambda^*) \rightarrow u(\lambda) \in H_0^1(\Omega) \cap C^0(\overline{\Omega})$ , where  $u(\lambda)$  is the minimal positive solution of (14.9) with the parameter value  $\lambda$ . The derivative  $I - \lambda T Df(u(\lambda)) \in \mathcal{L}(H_0^1(\Omega) \cap C^0(\overline{\Omega}); H_0^1(\Omega) \cap C^0(\overline{\Omega}))$  is an isomorphism.

Our goal is to study approximations of the regular solution set  $\{(\lambda, u(\lambda)); \lambda \in [0, \lambda^*)\}$  when we use the schemes (3.29) and (3.30) to numerically solve (14.9). Since we have introduced numerical quadrature rules to compute the integrals, we can no longer use the Galerkin approximation approach detailed in Section 13. We need to go back to the general results of Section 12. We are within the framework of Section 12 with  $m = 1$ ,  $X = Z = H_0^1(\Omega) \cap C^0(\overline{\Omega})$ , and  $F : \mathbb{R} \times X \rightarrow Z$  given by: for all  $\lambda \in \mathbb{R}$ ,  $v \in H_0^1(\Omega) \cap C^0(\overline{\Omega})$ ,

$$(14.10) \quad F(\lambda, v) = v - \lambda T f(v).$$

Both approximation schemes (3.29) and (3.30) can be written in the form

$$(14.11) \quad F_h(\lambda, u_h) = 0$$

where  $F_h : \mathbb{R} \times (H_0^1(\Omega) \cap C^0(\overline{\Omega})) \rightarrow H_0^1(\Omega) \cap C^0(\overline{\Omega})$  is given by: for all  $\lambda \in \mathbb{R}$ ,  $v \in H_0^1(\Omega) \cap C^0(\overline{\Omega})$

$$(14.12) \quad F_h(\lambda, v) = v - \lambda T_h f(v).$$

Here the operator  $T_h \in \mathcal{L}(C^0(\overline{\Omega}); V_h)$  is defined by: for  $g \in C^0(\overline{\Omega})$

$$(14.13) \quad \text{for all } v_h \in V_h \quad a_h(T_h g, v_h) = (g, v_h)_h,$$

where

$$(14.14) \quad a_h(u_h, v_h) = \int_{\Omega_h} \mathbf{grad}(u_h) \mathbf{grad} v_h \, dx,$$

$$(14.15) \quad (g, v_h)_h = \sum_{T \in \mathcal{T}_h} \left[ \frac{S_T}{3} \sum_{i=1}^3 (g v_h)(a_{i,T}) \right],$$

in the piecewise  $\mathcal{P}_1$  approximation (3.29) and

$$(14.16) \quad a_h(u_h, v_h) = \sum_{T \in \mathcal{T}_h} \left[ \sum_{i=1}^3 \frac{1}{6} \det(D F_T(\hat{m}_i)) (\mathbf{grad} u_h \mathbf{grad} v_h)(m_{i,T}) \right],$$

(14.17)

$$(g, v_h)_h = \sum_{T \in \mathcal{T}_h} \sum_{i=1}^3 \left[ \frac{1}{6} \det(D F_T(\hat{m}_i)) (g v_h)(m_{i,T}) \right],$$

in the piecewise  $\mathcal{P}_2$  approximation (3.30).

The main approximation properties of the operator  $T_h$  are summarized in the next theorem.

**Theorem 14.2.** *Let  $T_h$  be the operator defined in (14.13) with (14.14) and (14.15) in the case  $k = 1$  and with (14.16) and (14.17) in the case  $k = 2$ . Then there exists a constant  $C$  independent of  $h$  such that*

*i) for all  $f \in W^{1,q}(\Omega)$  with  $q > 2$ , the following estimate holds*

$$(14.18) \quad \|Tf - T_h f\|_{1,\Omega} + \|Tf - T_h f\|_{0,\infty,\Omega} \leq Ch \|f\|_{1,q,\Omega};$$

*ii) for all  $f \in W^{k+1,q}(\Omega)$  with  $q > 1$ , the following estimate holds*

$$(14.19) \quad \|Tf - T_h f\|_{0,\Omega} \leq Ch^{k+1} \|f\|_{k+1,q,\Omega};$$

*iii) for  $k = 2$  and  $f \in W^{2,q}(\Omega)$  with  $q > 1$ , the following estimate holds*

$$(14.20) \quad \|Tf - T_h f\|_{1,\Omega \cap \Omega_h} \leq Ch^2 \|f\|_{2,q,\Omega}.$$

Moreover for all  $\epsilon > 0$  there exists a constant  $C(\epsilon)$  such that for all  $f \in W^{k+1,1}(\Omega)$ ,

$$(14.21) \quad \|Tf - T_h f\|_{0,\infty,\Omega} \leq C(\epsilon) h^{k+1-\epsilon} \|f\|_{k+1,1,\Omega}.$$

If  $r_h : C^0(\overline{\Omega}) \rightarrow V_h$  denotes the corresponding Lagrange interpolation operator, then there exists a constant  $C$  such that for all  $f \in C^0(\overline{\Omega})$  and  $1 \leq p < \infty$

$$(14.22) \quad \|T_h f\|_{1,p,\Omega} \leq C \|r_h f\|_{0,\Omega}. \quad \square$$

These estimates are standard finite element estimates. For the proofs, see for instance in CIARLET [1991], in CROUZEIX and RAPPAZ [1989], and in the paper of WAHLBIN [1978].

To each solution  $(\lambda, u(\lambda))$  we associate the function  $\tilde{u}_h(\lambda) = \lambda T_h f(u(\lambda)) \in H_0^1(\Omega) \cap C^0(\overline{\Omega})$ . Since  $(\lambda, u(\lambda))$  is a solution to (14.9), the following relation holds

$$\tilde{u}_h(\lambda) = u(\lambda) - \lambda T f(u(\lambda)) + \lambda T_h f(u(\lambda)) = u(\lambda) - \lambda(T - T_h)f(u(\lambda)),$$

which implies that  $\tilde{u}_h(\lambda)$  is close to  $u(\lambda)$  in the norm of  $H_0^1(\Omega) \cap C^0(\overline{\Omega})$ . This choice is well adapted to get error estimates.

From the assumption (14.8), we deduce that  $F$  is of class  $C^2$ , so the hypothesis (12.4) is clearly satisfied.

The estimate of  $\|F_h(\lambda, \tilde{u}_h(\lambda))\|_{H_0^1(\Omega) \cap C^0(\overline{\Omega})}$  goes as follows. We set

$$\begin{aligned} \epsilon_h &= \|F_h(\lambda, \tilde{u}_h(\lambda))\|_{H_0^1(\Omega) \cap C^0(\overline{\Omega})} = \|\tilde{u}_h(\lambda) - \lambda T_h f(\tilde{u}_h(\lambda))\|_{H_0^1(\Omega) \cap C^0(\overline{\Omega})} \\ &= \lambda \|T_h(f(u(\lambda)) - f(\tilde{u}_h(\lambda)))\|_{H_0^1(\Omega) \cap C^0(\overline{\Omega})}. \end{aligned}$$

As a consequence of (14.22), we get

$$\epsilon_h \leq C \|r_h(f(u(\lambda)) - f(\tilde{u}_h(\lambda)))\|_{0,\Omega}.$$

Since for all  $a, v \in C^0(\overline{\Omega})$

$$\|r_h(av)\|_{0,\Omega} \leq C\|a\|_{0,\infty,\Omega}\|r_hv\|_{0,\Omega},$$

we have

$$\epsilon_h \leq C\|r_h(u(\lambda) - \tilde{u}_h(\lambda))\|_{0,\Omega} = C\lambda\|(r_hT - T_h)f(u(\lambda))\|_{0,\Omega}.$$

From the classical properties of the Lagrange interpolation operator  $r_h$  and from the estimate (14.19), we get the estimate

$$(14.23) \quad \epsilon_h \leq Ch^{k+1},$$

with  $k = 1$  or  $2$  corresponding to both approximations (3.29) and (3.30) respectively. So the assumption (12.5) is verified.

To check the stability (12.6), let us first notice that  $D_vF(\lambda, u(\lambda)) \in \mathcal{L}(H_0^1(\Omega) \cap C^0(\overline{\Omega}); H_0^1(\Omega) \cap C^0(\overline{\Omega}))$  is an isomorphism. Using a compactness argument on  $T$ , it is a simple matter to check that  $D_vF(\lambda, u(\lambda)) \in \mathcal{L}(W_0^{1,p}(\Omega); W_0^{1,p}(\Omega))$ , with  $p > 2$ , is an isomorphism. So there exists a constant  $C$  independent of  $h$  such that

$$(14.24) \quad \|D_vF(\lambda, u(\lambda))^{-1}\|_{H_0^1(\Omega) \cap C^0(\overline{\Omega}); H_0^1(\Omega) \cap C^0(\overline{\Omega})} \leq C,$$

$$(14.25) \quad \|D_vF(\lambda, u(\lambda))^{-1}\|_{W_0^{1,p}(\Omega); W_0^{1,p}(\Omega)} \leq C.$$

Let us notice that

$$D_vF_h(\lambda, \tilde{u}_h(\lambda)) = D_vF(\lambda, u(\lambda)) + D_vF_h(\lambda, \tilde{u}_h(\lambda)) - D_vF(\lambda, u(\lambda)).$$

So to prove the invertibility of  $D_vF_h(\lambda, \tilde{u}_h(\lambda)) \in \mathcal{L}(H_0^1(\Omega) \cap C^0(\overline{\Omega}); H_0^1(\Omega) \cap C^0(\overline{\Omega}))$ , we will bound the operator

$$A_h(\lambda) = D_vF(\lambda, u(\lambda))^{-1} [D_vF_h(\lambda, \tilde{u}_h(\lambda)) - D_vF(\lambda, u(\lambda))].$$

Let  $v$  be in  $H_0^1(\Omega) \cap C^0(\overline{\Omega})$ , then

$$\begin{aligned} A_h(\lambda)v &= D_vF(\lambda, u(\lambda))^{-1} [\lambda T Df(u(\lambda))v - \lambda T_h Df(\tilde{u}_h(\lambda))v] \\ &= D_vF(\lambda, u(\lambda))^{-1} \left[ \lambda(T - T_h)Df(u(\lambda))v + \right. \\ (14.26) \quad &\quad \left. \lambda T_h(Df(u(\lambda)) - Df(\tilde{u}_h(\lambda)))v \right]. \end{aligned}$$

Since the mapping  $T_h \in \mathcal{L}(C^0(\overline{\Omega}); W^{1,p}(\Omega))$  is uniformly bounded, see (14.22), we get from (14.26) and (14.25)

$$(14.27) \quad \|A_h(\lambda)v\|_{1,p,\Omega} \leq C\|v\|_{0,\infty,\Omega} \leq C\|v\|_{H_0^1(\Omega) \cap C^0(\overline{\Omega})}.$$



On the other hand

$$\begin{aligned} \|A_h(\lambda)v\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} &\leq C \left[ \|(T - T_h)Df(u(\lambda))v\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} \right. \\ &\quad \left. + \|[Df(u(\lambda)) - Df(\tilde{u}_h(\lambda))]v\|_{0,\infty,\Omega} \right], \end{aligned}$$

so with the regularity of  $f$ ,

$$\begin{aligned} \|A_h(\lambda)v\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} &\leq \\ &C \left[ \|(T - T_h)Df(u(\lambda))v\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} + \|(T - T_h)f(u(\lambda))\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} \|v\|_{1,p,\Omega} \right]. \end{aligned}$$

Using the convergence property (14.18), we get

$$(14.28) \quad \|A_h(\lambda)v\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} \leq Ch\|v\|_{1,p,\Omega}.$$

From the estimates (14.27) and (14.28) we deduce

$$\|A_h^2(\lambda)v\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} \leq Ch\|v\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})},$$

which implies

$$\|A_h^2(\lambda)\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega}); H_0^1(\Omega)\cap C^0(\bar{\Omega})} \leq Ch.$$

It is not difficult to check that the condition (1.1) can be relaxed to  $\|(A^{-1}B)^2\|_{X;X} \leq \ell^2$ , with  $\ell \in (0, 1)$ , see for instance CROUZEIX and RAPPAZ [1989], to get that  $A + B$  is an isomorphism. If  $h$  is chosen sufficiently small, then  $D_v F_h(\lambda, \tilde{u}_h(\lambda))$  is an isomorphism with a uniformly bounded inverse. So we have checked the stability assumption (12.6).

The mapping  $D_v F_h(\lambda, \cdot)$  is Lipschitzian since  $T_h \in \mathcal{L}(C^0(\bar{\Omega}); H_0^1(\Omega) \cap C^0(\bar{\Omega}))$  is uniformly bounded and  $f$  is  $C^2$ , that is (12.4) is satisfied.

**Theorem 14.3.** *Let  $\{(\lambda, u(\lambda)); \lambda \in [0, \bar{\lambda}]\}$  with  $\bar{\lambda} < \lambda^*$  be a regular solution set to the problem  $F(\lambda, v) = 0$  with  $F$  given in (14.10) and let  $F_h$  be the mapping defined in (14.12). There exist two constants  $h_0 > 0$ ,  $\delta_0 > 0$  and for  $h \leq h_0$  a  $C^2$  mapping  $u_h : \lambda \in [0, \bar{\lambda}] \rightarrow u_h(\lambda) \in V_h$  such that for all  $\lambda \in [0, \bar{\lambda}]$*

$$(F_h(\lambda, v_h) = 0 \text{ and } \|v_h - \tilde{u}_h(\lambda)\|_{H_0^1(\Omega)\cap C^0(\bar{\Omega})} \leq \delta_0) \iff v_h = u_h(\lambda),$$

and the set  $\{(\lambda, u_h(\lambda)); \lambda \in [0, \bar{\lambda}]\}$  is a regular solution set of problem (14.11). Moreover the following a priori estimates hold

$$(14.29) \quad \|u(\lambda) - u_h(\lambda)\|_{0,\Omega} + h\|u(\lambda) - u_h(\lambda)\|_{1,\Omega} \leq Ch^{k+1}$$

and

$$(14.30) \quad \|u(\lambda) - u_h(\lambda)\|_{0,\infty,\Omega} \leq C(\epsilon)h^{k+1-\epsilon}.$$

*Proof.* The hypotheses of Theorem 12.3 are a consequence of the above developments together with the compactness of the interval  $[0, \bar{\lambda}]$ . There exist  $h_0 > 0$ ,  $\delta_0 > 0$ , and for all  $h \leq h_0$  a  $C^2$  mapping  $u_h : \lambda \in [0, \bar{\lambda}] \rightarrow u_h(\lambda) \in H_0^1(\Omega) \cap C^0(\bar{\Omega})$  such that

$$(F_h(\lambda, v_h) = 0 \text{ and } \|v_h - \tilde{u}_h(\lambda)\|_{H_0^1(\Omega) \cap C^0(\bar{\Omega})} \leq \delta_0) \iff v_h = u_h(\lambda)$$

and the mapping  $D_v F_h(\lambda, u_h(\lambda)) \in \mathcal{L}(H_0^1(\Omega) \cap C^0(\bar{\Omega}); H_0^1(\Omega) \cap C^0(\bar{\Omega}))$  is invertible, with a uniformly bounded inverse.

The estimates (12.12) and (14.23) give

$$(14.31) \quad \|\tilde{u}_h(\lambda) - u_h(\lambda)\|_{H_0^1(\Omega) \cap C^0(\bar{\Omega})} \leq Ch^{k+1}.$$

Since

$$u(\lambda) - \tilde{u}_h(\lambda) = \lambda(T - T_h)f(u(\lambda)),$$

we deduce the estimates (14.29) and (14.30) from the (14.31) one and from Theorem 14.2.  $\square$

*Remark 14.1.* In the same way we proved the estimate (14.23), we can prove that

$$\left\| \frac{d}{d\lambda} F(\lambda, \tilde{u}_h(\lambda)) \right\|_{H_0^1(\Omega) \cap C^0(\bar{\Omega})} \leq Ch^{k+1}.$$

As a consequence of Theorem 12.4, there exist constants  $C$  and for  $\epsilon > 0$ ,  $C(\epsilon)$  such that for  $h \leq h_0$

$$\|Du(\lambda) - Du_h(\lambda)\|_{0,\Omega} + h\|Du(\lambda) - Du_h(\lambda)\|_{1,\Omega} \leq Ch^{k+1}$$

and

$$\|Du(\lambda) - Du_h(\lambda)\|_{0,\infty,\Omega} \leq C(\epsilon)h^{k+1-\epsilon}. \quad \square$$

*Remark 14.2.* In Theorem 14.3, we have got a priori estimates by using the relation (12.12). The inequality (12.13) will give a posteriori estimates. We do not want to emphasize this point here.  $\square$

## 15. Solution set containing a simple limit point

Let  $X$  and  $Z$  be two real Banach spaces and  $F : \mathbb{R} \times X \rightarrow Z$  be a nonlinear mapping of class  $C^1$ . Here for simplicity we have considered the parameter space  $\mathbb{R}^m$  with  $m = 1$ ; our approach can be immediately generalized for any  $m \in \mathbb{N}$ .

We assume that  $(\lambda_0, x_0) \in \mathbb{R} \times X$  is a regular solution to

$$(15.1) \quad F(\lambda, x) = 0,$$

but  $D_x F(\lambda_0, x_0) \in \mathcal{L}(X; Z)$  is no more an isomorphism. More precisely the following assumption holds:

$$(15.2) \quad \begin{aligned} &(\lambda_0, x_0) \text{ is a solution to (15.1), } DF(\lambda_0, x_0) \in \mathcal{L}(\mathbb{R} \times X; Z) \text{ is} \\ &\text{a Fredholm operator of index 1, } Z = \text{Range}(DF(\lambda_0, x_0)), \text{ and} \\ &D_x F(\lambda_0, x_0) \in \mathcal{L}(X; Z) \text{ is not an isomorphism.} \end{aligned}$$

The assumption (15.2) means that the solution  $(\lambda_0, x_0) \in \mathbb{R} \times X$  to (15.1) is such that  $D_x F(\lambda_0, x_0) \in \mathcal{L}(X; Z)$  is a Fredholm operator of index 0,  $\dim \text{Ker}(D_x F(\lambda_0, x_0)) = 1$ , and  $D_\lambda F(\lambda_0, x_0)$  does not belong to the range of  $D_x F(\lambda_0, x_0)$ .

Let  $\varphi \in X$ ,  $\varphi \neq 0$ , satisfy  $D_x F(\lambda_0, x_0)\varphi = 0$ ; so  $\text{Ker}(D_x F(\lambda_0, x_0)) = \text{span}\{\varphi\}$ .

Since  $D_x F(\lambda_0, x_0)$  is no more an isomorphism, it is not possible to describe the solution set of (15.1) in a neighborhood of  $(\lambda_0, x_0)$  as a solution branch parametrized by  $\lambda$ . To study it we introduce the mapping  $\Phi : \mathbb{R} \times \mathbb{R} \times X \rightarrow \mathbb{R} \times Z$  given by

$$(15.3) \quad \Phi(s, \lambda, x) = (B(x - x_0) - s, F(\lambda, x)),$$

where  $B$  is a functional in  $X'$  with  $B(\varphi) = 1$ . We check that

$$\Phi(0, \lambda_0, x_0) = 0$$

and for  $\delta \in \mathbb{R}$ ,  $v \in X$ ,

$$D_{\lambda x}^2 \Phi(0, \lambda_0, x_0)(\delta, v) = (Bv, D_\lambda F(\lambda_0, x_0)\delta + D_x F(\lambda_0, x_0)v).$$

The derivative  $D_{\lambda x}^2 \Phi(0, \lambda_0, x_0) \in \mathcal{L}(\mathbb{R} \times X; \mathbb{R} \times Z)$  is a Fredholm operator of index 0. Since  $\dim \text{Ker}(D_x F(\lambda_0, x_0))$  is finite, there exists a closed subspace  $X_1 \subset X$  such that  $X = \text{span}\{\varphi\} \oplus X_1$ . If  $(\delta, v \equiv \alpha\varphi + v_1) \in \mathbb{R} \times X$  with  $\alpha \in \mathbb{R}$  and  $v_1 \in X_1$  satisfies

$$D_{\lambda x}^2 \Phi(0, \lambda_0, x_0)(\delta, v) = 0,$$

then  $v_1 = 0$  and  $\delta = 0$  since  $DF(\lambda_0, x_0) \in \mathcal{L}(X_1; Z)$  is an isomorphism. Since  $B(\varphi) = 1$  so  $\alpha = 0$ . Finally  $D_{\lambda x}^2 \Phi(0, \lambda_0, x_0)$  is injective and therefore  $D_{\lambda x}^2 \Phi(0, \lambda_0, x_0) \in \mathcal{L}(\mathbb{R} \times X; \mathbb{R} \times Z)$  is an isomorphism.

To study the solutions of

$$(15.4) \quad \Phi(s, \lambda, x) = 0$$

in a neighborhood of  $(0, \lambda_0, x_0)$ , we can use the implicit function theorem to prove the following result.

**Theorem 15.1.** *Under the assumption (15.2), there exist a positive number  $\epsilon$ , a neighborhood  $W$  of  $(\lambda_0, x_0) \in \mathbb{R} \times X$ , and a unique mapping  $(\lambda, u) : s \in [-\epsilon, \epsilon] \rightarrow (\lambda(s), u(s)) \in \mathbb{R} \times X$  of class  $C^1$  such that for all  $s \in [-\epsilon, \epsilon]$*

$$(15.5) \quad F(\lambda(s), u(s)) = 0 \quad \text{and} \quad s = B(u(s) - x_0),$$

$$(15.6) \quad \lambda(0) = \lambda_0 \quad \text{and} \quad u(0) = x_0,$$

and if  $(\lambda, u) \in W$  is a solution to (15.1), then necessarily  $\lambda = \lambda(s)$ ,  $u = u(s)$  with  $s = B(u - x_0)$ . Moreover for  $s \in [-\epsilon, \epsilon]$ , the mapping  $D_{\lambda x}^2 \Phi(s, \lambda(s), u(s)) \in \mathcal{L}(\mathbb{R} \times X; \mathbb{R} \times Z)$  is an isomorphism.

Clearly  $(\lambda(s), u(s)) \neq (\lambda_0, x_0)$  for all  $s \neq 0$  since  $B(u(s) - x_0) = s$ . So the mapping  $s \in [-\epsilon, \epsilon] \rightarrow (\lambda(s), u(s)) \in \mathbb{R} \times X$  cannot be the constant function equal to  $(\lambda_0, x_0)$

and the set of solutions of (15.1) is not isolated. We can get more information on the variation of  $\lambda$  and  $u$  as a function of  $s$ . When differentiating with respect to  $s$  the identity  $\Phi(s, \lambda(s), u(s)) = 0$ , we get

$$B\left(\frac{d}{ds}u(0)\right) = 1,$$

$$D_\lambda F(\lambda_0, x_0)\frac{d}{ds}\lambda(0) + D_x F(\lambda_0, x_0)\frac{d}{ds}u(0) = 0,$$

and so  $\frac{d}{ds}\lambda(0) = 0$ ,  $\frac{d}{ds}u(0) = \varphi$ .

A pair  $(\lambda_0, x_0) \in \mathbb{R} \times X$  satisfying to the condition (15.2) is called a simple limit point of (15.1).

If we assume that  $F$  is a mapping of class  $C^p$  with  $p \geq 2$ , then we can get a more precise description of the solution set  $\{(\lambda(s), u(s)); s \in [-\epsilon, \epsilon]\}$ . We differentiate twice with respect to  $s$  the identity  $\Phi(s, \lambda(s), u(s)) = 0$  to get

$$D_\lambda F(\lambda_0, x_0)\frac{d^2}{ds^2}\lambda(0) + D_{xx}^2 F(\lambda_0, x_0)(\varphi, \varphi) + D_x F(\lambda_0, x_0)\frac{d^2}{ds^2}u(0) = 0,$$

since  $\frac{d}{ds}\lambda(0) = 0$  and  $\frac{d}{ds}u(0) = \varphi$ , and

$$B\left(\frac{d^2}{ds^2}u(0)\right) = 0.$$

Consequently  $\frac{d^2}{ds^2}\lambda(0)$  is different from 0 if and only if

$$(15.7) \quad D_{xx}^2 F(\lambda_0, x_0)(\varphi, \varphi) \notin \text{Range}(D_x F(\lambda_0, x_0)).$$

If  $(\lambda_0, x_0)$  is a simple limit point of (15.1) satisfying to (15.7), then for  $s$  in a neighborhood of 0 the following development holds

$$\lambda(s) = \lambda_0 + \frac{\frac{d^2}{ds^2}\lambda(0)}{2}s^2 + o(s^2).$$

In the figure 15.1, we have represented the branch  $\{(\lambda(s), B(u(s)))\}$  in a neighborhood of 0.

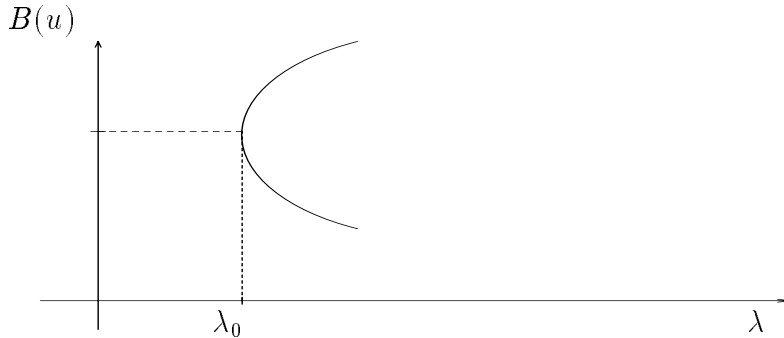


FIGURE 15.1:  $D^2\lambda(0) > 0$ .

A pair  $(\lambda_0, x_0) \in \mathbb{R} \times X$  satisfying to the conditions (15.2) and (15.7) is a nondegenerate simple limit point.

In the degenerate case  $\frac{d^2}{ds^2}\lambda(0) = 0$ , different cases can occur, see for instance the figure 15.2.

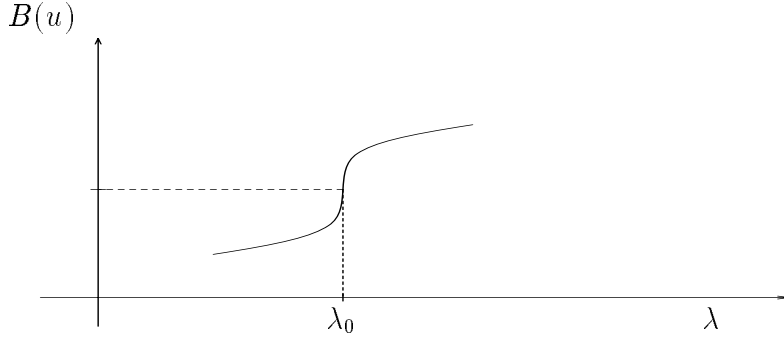


FIGURE 15.2:  $D^2\lambda(0) = 0$  AND  $D^3\lambda(0) \neq 0$ .

*Remark 15.1.* A nondegenerate simple limit point is in fact a junction point of two regular solution branches parametrized by  $\lambda$ . Let  $(\lambda_0, x_0)$  be a nondegenerate simple limit point and let  $s \in [-\epsilon, \epsilon] \rightarrow (\lambda(s), u(s)) \in \mathbb{R} \times X$  be the solution set given in Theorem 15.1. We assume here that  $F$  is of class  $C^p$  with  $p \geq 2$ . Then we will prove that there exist positive numbers  $0 < \epsilon_1 \leq \epsilon$  and  $C$  such that for all  $s \in [-\epsilon_1, \epsilon_1]$ ,  $s \neq 0$ ,

$$(15.8) \quad D_x F(\lambda(s), u(s)) \in \mathcal{L}(X; Z) \quad \text{is an isomorphism}$$

and

$$(15.9) \quad \|D_x F(\lambda(s), u(s))^{-1}\|_{Z; X} \leq \frac{C}{|s|}.$$

Since  $\frac{d^2}{ds^2}\lambda(0) \neq 0$ , there exist  $\epsilon_1$ ,  $0 < \epsilon_1 \leq \epsilon$  and a constant  $C_1$  such that

$$\text{for all } s \in [-\epsilon_1, \epsilon_1] \quad \left| \frac{d}{ds}\lambda(s) \right| \geq C_1 |s|.$$

The derivative  $D_{\lambda x}^2 \Phi(s, \lambda(s), u(s)) \in \mathcal{L}(\mathbb{R} \times X; \mathbb{R} \times Z)$  is an isomorphism with an inverse uniformly bounded with respect to  $s \in [-\epsilon_1, \epsilon_1]$ . For all  $s \in [-\epsilon_1, \epsilon_1]$  and  $f \in Z$ , there exists a unique  $(\delta, v) \in \mathbb{R} \times X$  with

$$D_{\lambda x}^2 \Phi(s, \lambda(s), u(s))(\delta, v) = (0, f)$$

and

$$|\delta| + \|v\|_X \leq C_2 \|f\|_Z,$$

where  $C_2$  is independent of  $s$ . Given  $f \in Z$  and  $s \neq 0$ , we use both relations

$$\begin{aligned} D_\lambda F(\lambda(s), u(s))\delta + D_x F(\lambda(s), u(s))v &= f, \\ D_\lambda F(\lambda(s), u(s))\frac{d}{ds}\lambda(s) + D_x F(\lambda(s), u(s))\frac{d}{ds}u(s) &= 0, \end{aligned}$$

and set  $w = v - \frac{\delta}{\frac{d}{ds}\lambda(s)}\frac{d}{ds}u(s)$ , to obtain

$$D_x F(\lambda(s), u(s))w = f$$

and

$$\|w\|_X \leq C_2 \left( 1 + \frac{1}{C_1|s|} \left\| \frac{d}{ds}u(s) \right\|_X \right) \|f\|_Z.$$

We have to check now that  $D_x F(\lambda(s), u(s))$  is injective with  $s \in [-\epsilon_1, \epsilon_1]$ ,  $s \neq 0$ . Let  $\psi \in X$  be such that

$$D_x F(\lambda(s), u(s))\psi = 0.$$

Then since  $B(\frac{d}{ds}u(s)) = 1$ , we have

$$D_{\lambda_x}^2 \Phi(s, \lambda(s), u(s))(-B(\psi)\frac{d}{ds}\lambda(s), \psi - B(\psi)\frac{d}{ds}u(s)) = 0,$$

and consequently  $\psi = B(\psi)\frac{d}{ds}u(s)$  and  $B(\psi)\frac{d}{ds}\lambda(s) = 0$ . Therefore  $\psi = 0$ .  $\square$

## 16. Galerkin approximation of simple limit points

Let  $X, Y$  be two real Banach spaces, and let  $Y'$  denote the dual space of  $Y$ . Given a  $C^p$  mapping  $F : \mathbb{R} \times X \rightarrow Y'$ , with  $p \geq 2$ , we shall analyze Galerkin approximations of the problem

$$(16.1) \quad F(\lambda, x) = 0$$

in a neighborhood of a regular solution  $(\lambda_0, x_0) \in \mathbb{R} \times X$  satisfying to (15.2). The study of the exact problem has been done in Section 15.

To analyze the approximation of the simple limit point  $(\lambda_0, x_0)$ , we consider the case where the mapping  $F$  has the following structure

$$(16.2) \quad F(\lambda, x) = Lx + G(\lambda, x)$$

where

$$(16.3) \quad L \in \mathcal{L}(X; Y') \quad \text{is an isomorphism}$$

and the  $C^p$  mapping  $G : \mathbb{R} \times X \rightarrow Y'$  satisfies:

$$(16.4) \quad \text{for all } (\lambda, x) \in \mathbb{R} \times X \quad D_x G(\lambda, x) \in \mathcal{L}(X; Y') \quad \text{is compact.}$$

Then problem (16.1) reads: find  $\lambda \in \mathbb{R}$ ,  $x \in X$  such that

$$(16.5) \quad \text{for all } y \in Y \quad \langle Lx, y \rangle_{Y'Y} + \langle G(\lambda, x), y \rangle_{Y'Y} = 0.$$

Given a family  $\{X_h\}_{0 < h \leq 1}$  of finite-dimensional subspaces of  $X$  and a family  $\{Y_h\}_{0 < h \leq 1}$  of finite-dimensional subspaces of  $Y$ , a Galerkin approximation of (16.5) consists in finding  $\lambda \in \mathbb{R}$  and  $x_h \in X_h$  such that

$$(16.6) \quad \text{for all } y_h \in Y_h \quad \langle Lx_h, y_h \rangle_{Y'Y} + \langle G(\lambda, x_h), y_h \rangle_{Y'Y} = 0.$$

Let  $b : X \times Y \rightarrow \mathbb{R}$  be the bilinear form defined by

$$(16.7) \quad \text{for all } (x, y) \in X \times Y \quad b(x, y) = \langle Lx, y \rangle_{Y'Y}.$$

Then the stability assumption reads: for all  $h \in (0, 1]$ ,

$$(16.8) \quad \inf_{\substack{x \in X_h \\ \|x\|_X = 1}} \sup_{\substack{y \in Y_h \\ \|y\|_Y = 1}} b(x, y) \geq \beta > 0$$

and

$$(16.9) \quad \dim X_h = \dim Y_h.$$

The consistency assumption reads

$$(16.10) \quad \text{for all } x \in X \quad \lim_{h \rightarrow 0} \inf_{x_h \in X_h} \|x - x_h\|_X = 0.$$

To study the approximation (16.6) we proceed in a way similar to the one in Theorem 7.1. Let  $\Pi_{X_h} \in \mathcal{L}(X; X_h)$  and  $\Pi_{Y_h} \in \mathcal{L}(Y; Y_h)$  be the two projectors defined by

$$(16.11) \quad \text{for all } y_h \in Y_h \quad b(x - \Pi_{X_h} x, y_h) = 0$$

and

$$(16.12) \quad \text{for all } x_h \in X_h \quad b(x_h, y - \Pi_{Y_h} y) = 0.$$

Then we construct the family  $\{F_h\}_{0 < h \leq 1}$  of mappings  $F_h : \mathbb{R} \times X \rightarrow Y'$  by setting: for  $\lambda \in \mathbb{R}$  and for  $x \in X$ ,  $y \in Y$ ,

$$\langle F_h(\lambda, x), y \rangle_{Y'Y} = \langle F(\lambda, x), \Pi_{Y_h} y \rangle_{Y'Y} + b(x, y - \Pi_{Y_h} y).$$

It is a simple matter (see the proof of Theorem 7.1) to check that problem (16.6) is equivalent to find  $\lambda \in \mathbb{R}$  and  $x_h \in X$  such that

$$(16.13) \quad F_h(\lambda, x_h) = 0.$$

To study problem (16.13) in some neighborhood of the simple limit point  $(\lambda_0, x_0) \in \mathbb{R} \times X$ , we define the mapping  $\Phi_h : \mathbb{R} \times \mathbb{R} \times X \rightarrow \mathbb{R} \times Y'$  by

$$\Phi_h(s, \lambda, x) = (B(x - u(s)), F_h(\lambda, x)),$$

the discrete analogue of  $\Phi$  given in (15.3). Here  $(\lambda, u) : [-\epsilon, \epsilon] \rightarrow \mathbb{R} \times X$  is the exact solution set in a neighborhood of  $(\lambda_0, x_0)$ . Clearly if  $(s, \lambda_h, x_h) \in \mathbb{R} \times \mathbb{R} \times X$  satisfies

$$(16.14) \quad \Phi_h(s, \lambda_h, x_h) = 0$$

then  $(\lambda_h, x_h)$  is a solution to (16.13) and consequently to (16.6). Conversely if  $(\lambda_h, x_h) \in \mathbb{R} \times X$  be a solution to (16.6) or (16.13), then  $(s \equiv B(x_h) - B(x_0), \lambda_h, x_h)$  is a solution to (16.14).

The study of problem (16.14) can be done as an application of the results of Section 12 with the mapping  $\Phi_h$  and the solution set  $\{(s, \lambda(s), u(s)); s \in [-\epsilon, \epsilon]\}$ . Let us check that  $\Phi_h$  satisfies successively the consistency (12.10), the stability (12.11), and the Lipschitzian condition (12.9) at the solutions  $(\lambda(s), u(s))$ .

For a given  $s \in [-\epsilon, \epsilon]$ , the following estimate holds

$$(16.15) \quad \begin{aligned} & \|\Phi_h(s, \lambda(s), u(s))\|_{\mathbb{R} \times Y'} \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y=1}} \langle F_h(\lambda(s), u(s)), y \rangle_{Y'Y} = \sup_{\substack{y \in Y \\ \|y\|_Y=1}} b(u(s), y - \Pi_{Y_h} y) \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y=1}} b(u(s) - \Pi_{X_h} u(s), y) \leq \|L\|_{X;Y'} \|u(s) - \Pi_{X_h} u(s)\|_X. \end{aligned}$$

For  $x \in X$  and  $x_h \in X_h$ , the inf-sup condition (16.8) implies

$$\begin{aligned} \|x_h - \Pi_{X_h} x\|_X &\leq \frac{1}{\beta} \sup_{\substack{y_h \in Y_h \\ \|y_h\|_Y=1}} b(x_h - \Pi_{X_h} x, y_h) \\ &= \frac{1}{\beta} \sup_{\substack{y_h \in Y_h \\ \|y_h\|_Y=1}} b(x_h - x, y_h) \leq \|L\|_{X;Y'} \frac{1}{\beta} \|x_h - x\|_X. \end{aligned}$$

This last inequality and the triangular inequality give

$$(16.16) \quad \|x - \Pi_{X_h} x\|_X \leq \left(1 + \frac{\|L\|_{X;Y'}}{\beta}\right) \inf_{x_h \in X_h} \|x - x_h\|_X.$$

Then we use the estimate (16.15) with the (16.16) one and the consistency assumption (16.10) to check the consistency (12.10).



Noticing that for  $\delta \in \mathbb{R}$  and  $v \in X$

$$\begin{aligned} & D_{\lambda x}^2 \Phi_h(s, \lambda(s), u(s))(\delta, v) \\ &= (B(v), D_\lambda F_h(\lambda(s), u(s))\delta + D_x F_h(\lambda(s), u(s))v) \\ &= D_{\lambda x}^2 \Phi(s, \lambda(s), u(s))(\delta, v) + \left(0, (D_\lambda F_h(\lambda(s), u(s)) - D_\lambda F(\lambda(s), u(s)))\delta\right. \\ &\quad \left.+ (D_x F_h(\lambda(s), u(s)) - D_x F(\lambda(s), u(s)))v\right) \end{aligned}$$

and since  $D_{\lambda x}^2 \Phi(s, \lambda(s), u(s)) \in \mathcal{L}(\mathbb{R} \times X, Y')$  is an isomorphism, the stability assumption (12.11) is checked, see (1.1), if we prove

$$(16.17) \quad \lim_{h \rightarrow 0} \max_{s \in [-\epsilon, \epsilon]} \|D_\lambda F_h(\lambda(s), u(s)) - D_\lambda F(\lambda(s), u(s))\|_{Y'} = 0$$

and

$$(16.18) \quad \lim_{h \rightarrow 0} \max_{s \in [-\epsilon, \epsilon]} \|D_x F_h(\lambda(s), u(s)) - D_x F(\lambda(s), u(s))\|_{X; Y'} = 0.$$

We have successively

$$\begin{aligned} & \|D_\lambda F_h(\lambda(s), u(s)) - D_\lambda F(\lambda(s), u(s))\|_{Y'} \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} \langle D_\lambda F_h(\lambda(s), u(s)) - D_\lambda F(\lambda(s), u(s)), y \rangle_{Y'Y} \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} [\langle D_\lambda F(\lambda(s), u(s)), \Pi_{Y_h} y - y \rangle_{Y'Y}] \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} b(L^{-1} D_\lambda F(\lambda(s), u(s)), \Pi_{Y_h} y - y) \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} b((I - \Pi_{X_h})L^{-1} D_\lambda F(\lambda(s), u(s)), y) \\ &\leq \|L\|_{X; Y'} \|(I - \Pi_{X_h})L^{-1} D_\lambda F(\lambda(s), u(s))\|_X. \end{aligned}$$

Then with the estimates (16.16) and (16.10), we conclude that

$$\lim_{h \rightarrow 0} \max_{s \in [-\epsilon, \epsilon]} \|D_\lambda F_h(\lambda(s), u(s)) - D_\lambda F(\lambda(s), u(s))\|_{Y'} = 0.$$

Let us prove now the relation (16.18). For  $w \in X$  and  $y \in Y$  given we have

$$\begin{aligned} & \langle (D_x F_h(\lambda(s), u(s)) - D_x F(\lambda(s), u(s)))w, y \rangle_{Y'Y} \\ &= \langle D_x F(\lambda(s), u(s))w, \Pi_{Y_h} y - y \rangle_{Y'Y} + b(w, y - \Pi_{Y_h} y) \\ &= \langle (L + D_x G(\lambda(s), u(s)))w, \Pi_{Y_h} y - y \rangle_{Y'Y} + \langle Lw, y - \Pi_{Y_h} y \rangle_{Y'Y} \\ &= b(L^{-1} D_x G(\lambda(s), u(s))w, \Pi_{Y_h} y - y) \\ &= -b((I - \Pi_{X_h})L^{-1} D_x G(\lambda(s), u(s))w, y). \end{aligned}$$

Let us prove now that

$$\lim_{h \rightarrow 0} \|(I - \Pi_{X_h})L^{-1}D_x G(\lambda(s), u(s))\|_{X;X} = 0.$$

From the assumptions (16.3) and (16.4), we get that the operator  $L^{-1}D_x G(\lambda(s), u(s)) \in \mathcal{L}(X; X)$  is compact. Then from the consistency (16.10), the estimate (16.16), and the result (1.3), we get that for  $s \in [-\epsilon, \epsilon]$  the above limit is zero. With the compactness of the interval  $[-\epsilon, \epsilon]$  and the continuity of  $s \in [-\epsilon, \epsilon] \rightarrow (\lambda(s), u(s)) \in \mathbb{R} \times X$  we conclude that the relation (16.18) is true.

Finally we verify the Lipschitzian condition (12.9) for  $D_{\lambda x}^2 \Phi_h$ . For  $(\lambda, x) \in \mathbb{R} \times X$ ,  $s \in [-\epsilon, \epsilon]$ , and  $(\delta, v) \in \mathbb{R} \times X$

$$(16.19) \quad \begin{aligned} [D_{\lambda x}^2 \Phi_h(s, \lambda(s), u(s)) - D_{\lambda x}^2 \Phi_h(s, \lambda, x)](\delta, v) \\ = (0, [DF_h(\lambda(s), u(s)) - DF_h(\lambda, x)](\delta, v)). \end{aligned}$$

Moreover

$$(16.20) \quad \begin{aligned} & \| (DF_h(\lambda(s), u(s)) - DF_h(\lambda, x))(\delta, v) \|_{Y'} \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} \langle (DF_h(\lambda(s), u(s)) - DF_h(\lambda, x))(\delta, v), y \rangle_{Y'Y} \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} \langle (DF(\lambda(s), u(s)) - DF(\lambda, x))(\delta, v), \Pi_{Y_h} y \rangle_{Y'Y} \\ &\leq \|\Pi_{Y_h}\|_{Y;Y} \|DF(\lambda(s), u(s)) - DF(\lambda, x)\|_{\mathbb{R} \times X; Y'} (|\delta| + \|v\|_X). \end{aligned}$$

With the relation (16.19), the estimate (16.20) and the fact that  $DF$  is Lipschitzian at  $(\lambda(s), u(s))$ , ( $F$  is assumed to be  $C^2$ ), we conclude that (12.9) is verified for  $D_{\lambda x}^2 \Phi_h$  at the points  $(s, \lambda(s), u(s))$ ,  $s \in [-\epsilon, \epsilon]$ .

We can apply now Theorem 12.3.

**Theorem 16.1.** *Let  $F : \mathbb{R} \times X \rightarrow Y'$  satisfy the assumptions (16.2), (16.3), and (16.4) with  $p \geq 2$ . Let  $(\lambda_0, x_0) \in \mathbb{R} \times X$  be a simple limit point of problem (16.1). Under the hypotheses (16.8), (16.9), and (16.10), there exist  $h_0 > 0$ ,  $\delta_0 > 0$ , and for all  $0 < h \leq h_0$  a  $C^p$  mapping  $(\lambda_h, u_h) : s \in [-\epsilon, \epsilon] \rightarrow (\lambda_h(s), u_h(s)) \in \mathbb{R} \times X$  such that*

$$(\Phi_h(s, \lambda, x_h) = 0 \quad \text{and} \quad (\lambda, x_h) \in \overline{B}((\lambda(s), u(s)), \delta_0)) \iff (\lambda = \lambda_h(s), x_h = u_h(s)))$$

Moreover there exists a constant  $C > 0$  such that for all  $s \in [-\epsilon, \epsilon]$

$$(16.21) \quad |\lambda(s) - \lambda_h(s)| + \|u(s) - u_h(s)\|_X \leq C \inf_{x_h \in X_h} \|u(s) - x_h\|_X.$$

*Proof.* We have checked before that we can apply Theorem 12.3. The estimate (16.21) is a consequence of the (12.12) one and of the relations (16.15) and (16.16).  $\square$

We assume now that the mapping  $F$  is of class  $C^p$ ,  $p \geq 3$ , and that  $(\lambda_0, x_0)$  is a nondegenerate simple limit point, which means, see (15.7),

$$D_{xx}^2 F(\lambda_0, x_0)(\varphi, \varphi) \notin \text{Range}(D_x F(\lambda_0, x_0)).$$

We have for  $j = 0, 1, 2$

$$\begin{aligned} \left\| \frac{d^j}{ds^j} \Phi_h(s, \lambda(s), u(s)) \right\|_{\mathbb{R} \times Y'} &= \left\| \frac{d^j}{ds^j} F_h(\lambda(s), u(s)) \right\|_{Y'} \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} \left\langle \frac{d^j}{ds^j} F_h(\lambda(s), u(s)), y \right\rangle_{Y'Y} \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} b\left(\frac{d^j}{ds^j} u(s), y - \Pi_{Y_h} y\right) \\ &= \sup_{\substack{y \in Y \\ \|y\|_Y = 1}} b\left(\frac{d^j}{ds^j} u(s) - \Pi_{X_h} \frac{d^j}{ds^j} u(s), y\right). \end{aligned}$$

So with this last relation and the (16.16) one we get for  $j = 0, 1, 2$

$$(16.22) \quad \lim_{h \rightarrow 0} \max_{s \in [-\epsilon, \epsilon]} \left\| \frac{d^j}{ds^j} \Phi_h(s, \lambda(s), u(s)) \right\|_{\mathbb{R} \times Y'} = 0.$$

Then with the relation (16.22), Theorem 12.4 implies for  $j = 0, 1, 2$

$$(16.23) \quad \lim_{h \rightarrow 0} \max_{s \in [-\epsilon, \epsilon]} \left\{ \left| \frac{d^j}{ds^j} (\lambda(s) - \lambda_h(s)) \right| + \left\| \frac{d^j}{ds^j} (u(s) - u_h(s)) \right\|_X \right\} = 0.$$

We have assumed that  $(\lambda_0, x_0)$  is a nondegenerate simple limit point, so

$$\lambda(0) = \lambda_0, \quad \frac{d}{ds} \lambda(0) = 0, \quad \frac{d^2}{ds^2} \lambda(0) \neq 0.$$

There exist  $\epsilon_1$  with  $0 < \epsilon_1 \leq \epsilon$ ,  $h_1$  with  $0 < h_1 \leq h_0$ ,  $C > 0$ , and for  $0 < h \leq h_1$  a unique  $s_h \in [-\epsilon_1, \epsilon_1]$  such that

$$(16.24) \quad \frac{d}{ds} \lambda_h(s_h) = 0,$$

$$(16.25) \quad \left| \frac{d^2}{ds^2} \lambda_h(s) \right| \geq C \quad \text{for all } s \in [-\epsilon_1, \epsilon_1].$$

So the point  $(\lambda_h(s_h), u_h(s_h))$  is a nondegenerate simple limit point of the problem (16.13). We know then that the solutions  $(\lambda_h(s), u_h(s))$  form two regular solution branches parametrized by  $\lambda$  which meet at the point  $(\lambda_h(s_h), u_h(s_h))$ , see Remark 15.1.

Notice that the estimate (16.21) does not give in general an error between  $(\lambda_0, x_0)$  and  $(\lambda_h(s_h), u_h(s_h))$  where  $(\lambda_h(s_h), u_h(s_h))$  is the nondegenerate simple limit point of the approximate problem. Let us try to compare  $\lambda_0 \equiv \lambda(0)$  and  $\lambda_{0h} \equiv \lambda_h(s_h)$ . We write

$$\lambda_0 - \lambda_{0h} = \lambda_0 - \lambda_h(0) - \int_0^{s_h} \frac{d}{ds} \lambda_h(s) ds.$$

The estimate (16.25) implies that  $\frac{d^2}{ds^2} \lambda_h(s)$  does not change sign in  $[0, s_h]$  and with (16.24)

$$|\lambda_0 - \lambda_{0h}| \leq |\lambda_0 - \lambda_h(0)| + |s_h| \left| \frac{d}{ds} \lambda_h(0) \right|.$$

With both relations (16.24) and (16.25), we deduce that

$$(16.26) \quad |s_h| \leq \frac{1}{C} \left| \frac{d}{ds} \lambda_h(0) \right|$$

and then from the above inequality we get

$$(16.27) \quad |\lambda_0 - \lambda_{0h}| \leq |\lambda_0 - \lambda_h(0)| + \frac{1}{C} \left| \frac{d}{ds} \lambda(0) - \frac{d}{ds} \lambda_h(0) \right|^2.$$

From the estimate (16.26) and the boundedness of  $\frac{d}{ds} u_h(s)$ ,  $s \in [-\epsilon, \epsilon]$ , we get

$$(16.28) \quad \|u_h(0) - u_h(s_h)\|_X \leq \tilde{C} \left| \frac{d}{ds} \lambda_h(0) \right|,$$

and then

$$(16.29) \quad \|u(0) - u_h(s_h)\|_X \leq \|u(0) - u_h(0)\|_X + \tilde{C} \left| \frac{d}{ds} \lambda_h(0) - \frac{d}{ds} \lambda(0) \right|.$$

The estimates (16.27) and (16.29) measure the error between the nondegenerate simple limit point  $(\lambda_0, x_0)$  and its approximation  $(\lambda_h(s_h), u_h(s_h))$ .

It is often possible to improve the estimate for  $|\lambda(0) - \lambda_h(0)|$  along the following lines. The Taylor expansion gives

$$(16.30) \quad \begin{aligned} F(\lambda_h(0), u_h(0)) &= F(\lambda_0, x_0) + D_x F(\lambda_0, x_0)(u_h(0) - x_0) \\ &\quad + D_\lambda F(\lambda_0, x_0)(\lambda_h(0) - \lambda_0) \\ &\quad + \mathcal{O}(|\lambda_h(0) - \lambda_0|^2 + \|u_h(0) - x_0\|_X^2). \end{aligned}$$

On the other hand since  $F_h(\lambda_h(0), u_h(0)) = 0$ ,  $u_h(0)$  is a solution of (16.6) and

$$(16.31) \quad \text{for all } y \in Y \quad \langle F(\lambda_h(0), u_h(0)), \Pi_{Y_h} y \rangle_{Y'Y} = 0.$$

We choose now  $y_0 \in Y$ ,  $\|y_0\| = 1$ , such that

$$(16.32) \quad \text{for all } \psi \in X \quad \langle D_x F(\lambda_0, x_0) \psi, y_0 \rangle_{Y'Y} = 0.$$

From the relations (16.30), (16.31), and (16.32) we deduce

$$(16.33) \quad \langle D_\lambda F(\lambda_0, x_0), \Pi_{Y_h} y_0 \rangle_{Y'Y} (\lambda_h(0) - \lambda_0) \\ + \langle D_x F(\lambda_0, x_0)(u_h(0) - x_0), \Pi_{Y_h} y_0 - y_0 \rangle_{Y'Y} = \mathcal{O}(|\lambda_h(0) - \lambda_0|^2 + \|u_h(0) - x_0\|_X^2).$$

The point  $(\lambda_0, x_0)$  is a simple limit point, so  $D_\lambda F(\lambda_0, x_0) \notin \text{Range}(D_x F(\lambda_0, x_0))$  and then

$$(16.34) \quad \langle D_\lambda F(\lambda_0, x_0), y_0 \rangle_{Y'Y} \neq 0.$$

The consistency assumption (16.10) and the estimate (16.16) imply

$$(16.35) \quad \lim_{h \rightarrow 0} \langle D_\lambda F(\lambda_0, x_0), y_0 - \Pi_{Y_h} y_0 \rangle_{Y'Y} = 0.$$

From the relations (16.33), (16.34), and (16.35) we deduce the existence of a constant  $C > 0$  independent of  $h$  such that

$$(16.36) \quad |\lambda_0 - \lambda_h(0)| \leq C[\|u_h(0) - x_0\|_X + \|y_0 - \Pi_{Y_h} y_0\|_Y] \|u_h(0) - x_0\|_X.$$

With (16.36) the estimate (16.27) can be improved in the following way

$$(16.37) \quad |\lambda_0 - \lambda_{0h}| \leq C(\|u_h(0) - x_0\|_X^2 + \\ |\frac{d}{ds} \lambda(0) - \frac{d}{ds} \lambda_h(0)|^2 + \|u_h(0) - x_0\|_X \|y_0 - \Pi_{Y_h} y_0\|_Y).$$

*Remark 16.1.* So far we do not have presented a posteriori error estimates in the case of simple limit points. In fact under the assumptions of Theorem 16.1, we can prove there exists a constant  $C > 0$  such that for all  $s \in [-\epsilon, \epsilon]$

$$(16.38) \quad |\lambda_h(s) - \lambda(s)| + \|u_h(s) - u(s)\|_X \leq C \|F(\lambda_h(s), u_h(s))\|_{Y'}.$$

Indeed the estimate (12.13) of Theorem 12.3 with the mapping  $\Phi$  gives

$$|\lambda_h(s) - \lambda(s)| + \|u_h(s) - u(s)\|_X \leq C \|\Phi(s, \lambda_h(s), u_h(s))\|_{\mathbb{R} \times Y'}.$$

Going back to the definition of  $\Phi_h$  and  $\Phi$  we deduce that

$$B(u_h(s) - x_0) - s = B(u_h(s) - u(s)) + B(u(s) - x_0) = 0$$

and then

$$\Phi(s, \lambda_h(s), u_h(s)) = (0, F(\lambda_h(s), u_h(s))). \quad \square$$

*Remark 16.2.* The goal of this section was to present a general but simple method to study approximations of simple limit points. We have not emphasized on what are the weakest assumptions on  $F$ .

The assumption (16.2) on the shape of  $F$  can be avoided, but then the study is different, see CALOZ [1994].

The function  $F$  does not need to be a  $C^2$  function but  $DF$  needs to be Lipschitzian at  $(\lambda(s), u(s))$  for  $s \in [-\epsilon, \epsilon]$ .

When we have applied the results of Section 12, we have decided to associate to each solution  $(s, \lambda(s), u(s))$ ,  $s \in [-\epsilon, \epsilon]$ , the solution  $(\lambda(s), u(s))$  itself instead of some approximations of it, say  $(\tilde{\lambda}_h(s), \tilde{u}_h(s))$ . Adding the appropriate assumptions on these approximations, see Sections 12 and 14, we can reproduce the approach of the present section to get estimates better suited to study the error in different norms. An other method would be to generalize Theorem 6.4.  $\square$

### 17. Approximation of simple limit points. An example

The goal of the section is to apply to a simple and concrete example the general results on the approximation of simple limit points in Section 16.

For simplicity we consider again problem (14.3). Let  $\Omega \subset \mathbb{R}^2$  be a regular convex open bounded set and let  $f : H_0^1(\Omega) \rightarrow L^2(\Omega)$  be given by: for all  $v \in H_0^1(\Omega)$

$$(17.1) \quad \text{for all } x \in \Omega \quad f(v)(x) = v^2(x) + 1.$$

We define the mapping  $F : \mathbb{R} \times H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  by: for all  $(\lambda, v) \in \mathbb{R} \times H_0^1(\Omega)$ ,  $w \in H_0^1(\Omega)$

$$(17.2) \quad \langle F(\lambda, v), w \rangle_{H^{-1}(\Omega); H_0^1(\Omega)} = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w \, dx - \lambda \int_{\Omega} f(v) w \, dx.$$

There exists a solution  $(\lambda^*, u^*) \in \mathbb{R} \times H_0^1(\Omega)$ ,  $u^*$  positive in  $\Omega$ , to the problem

$$(17.3) \quad F(\lambda, v) = 0,$$

such that  $D_v F(\lambda^*, u^*) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is not an isomorphism, see Remark 3.2. We can study more precisely the mapping  $D_v F(\lambda^*, u^*)$ . If  $T \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$  is the inverse of the  $(-\Delta)$  with homogeneous Dirichlet boundary conditions, then for  $w \in H_0^1(\Omega)$

$$TD_v F(\lambda^*, u^*)w = w - \lambda^* TDf(u^*)w$$

and the mapping  $TD_v F(\lambda^*, u^*) \in \mathcal{L}(H_0^1(\Omega); H_0^1(\Omega))$  has a kernel of dimension 1, namely  $\text{span}\{\varphi\}$ , where the function  $\varphi \in H_0^1(\Omega)$ ,  $\varphi > 0$  in  $\Omega$ , is an eigenvector corresponding to the least eigenvalue  $\lambda^*$  of the problem: find  $\delta \in \mathbb{R}$  and  $w \in H_0^1(\Omega)$ ,  $w \neq 0$ , such that

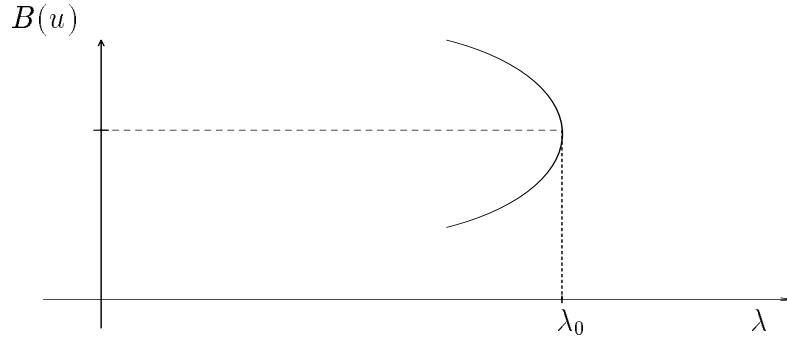
$$-\Delta w = \delta Df(u^*)w \quad \text{in } \Omega,$$

(see Remark 3.4, the proof of Theorem 3.4, and Theorem 3.2). Therefore the mapping  $TD_v F(\lambda^*, u^*)$  is a Fredholm operator of index 0. The range of  $TD_v F(\lambda^*, u^*)$  is orthogonal to  $\varphi \in H_0^1(\Omega)$ , with respect to the scalar product  $(u, v)_{1, \Omega} = \int_{\Omega} \mathbf{grad} u \mathbf{grad} v \, dx$ . Since both  $u^*$  and  $\varphi$  are positive in  $\Omega$ , we easily check that  $TDF(\lambda^*, u^*) \in \mathcal{L}(\mathbb{R} \times H_0^1(\Omega); H_0^1(\Omega))$  is a Fredholm operator of index 1 and  $H_0^1(\Omega) = \text{Range}(TDF(\lambda^*, u^*))$ . The operator  $T \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$  is an isomorphism, so  $DF(\lambda^*, u^*) \in \mathcal{L}(\mathbb{R} \times H_0^1(\Omega); H^{-1}(\Omega))$  is a Fredholm operator of index 1.

The pair  $(\lambda^*, u^*) \in \mathbb{R} \times H_0^1(\Omega)$  satisfies to the condition (15.2). It is a simple limit point of the problem (17.3). Notice that this result has been already checked in Remark 3.4. It is a simple matter to prove that  $D_{vv}^2 F(\lambda^*, u^*)(\varphi, \varphi)$  does not belong to the range of  $D_v F(\lambda^*, u^*)$ , so the relation (15.7) is satisfied and  $(\lambda^*, u^*)$  is a nondegenerate simple limit point.

Let  $B \in H^{-1}(\Omega)$  with  $B(\varphi) = 1$  denote a functional appearing in the definition of the mapping  $\Phi$  in (15.3). For instance we could consider

$$B(v) = \frac{1}{\|\varphi\|_{0, \Omega}^2} \int_{\Omega} \varphi v \, dx.$$

FIGURE 17.1:  $(\lambda^*, u^*)$  NONDEGENERATE SIMPLE LIMIT POINT.

Then graphically we have the situation of Figure 17.1.

Here  $(\lambda, u) : s \in [-\epsilon, \epsilon] \rightarrow (\lambda(s), u(s)) \in \mathbb{R} \times H_0^1(\Omega)$  denotes the  $C^\infty$  mapping such that for all  $s \in [-\epsilon, \epsilon]$

$$\begin{aligned} F(\lambda(s), u(s)) &= 0 \quad \text{and} \quad B(u(s) - u^*) = s, \\ \lambda(0) &= \lambda^* \quad \text{and} \quad u(0) = u^*. \end{aligned}$$

Given a family  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  of isoparametric triangulations of  $\Omega$ , we consider the corresponding family  $\{V_h\}_{0 < h \leq 1}$  of finite element subspaces of degree  $k = 1$  or  $2$

$$(17.4) \quad V_h = \{v \in C^0(\mathbb{R}^2); v(x) = 0 \text{ for all } x \notin \Omega_h, v|_T \in P_T \text{ for all } T \in \mathcal{T}_h\},$$

see (3.27). A Galerkin approximation to the problem (17.3) consists in finding  $\lambda \in \mathbb{R}$  and  $u_h \in V_h$  such that

$$(17.5) \quad \text{for all } v_h \in V_h \quad \langle F(\lambda, u_h), v_h \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} = 0.$$

We will study problem (17.5) in a neighborhood of the nondegenerate simple limit point  $(\lambda^*, u^*)$  by applying the results of Section 16 with  $X = Y = H_0^1(\Omega)$ .

It is not difficult to check that the mapping  $F$  in (17.2) has the structure (16.2) with the isomorphism  $L = (-\Delta) : H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  and the  $C^\infty$  mapping  $G(\lambda, \cdot) = -\lambda f(\cdot)$  such that  $\lambda Df(v) \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  is compact for all  $v \in H_0^1(\Omega)$ .

The bilinear form  $b : H_0^1(\Omega) \times H_0^1(\Omega) \rightarrow \mathbb{R}$  in our example is

$$\text{for all } v, w \in H_0^1(\Omega) \quad b(v, w) = \int_{\Omega} \mathbf{grad} v \mathbf{grad} w \, dx.$$

Then problem (17.5) can be written equivalently in the form: find  $\lambda \in \mathbb{R}$  and  $u_h \in H_0^1(\Omega)$  such that for all  $v \in H_0^1(\Omega)$

$$(17.6) \quad \langle F_h(\lambda, u_h), v \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} \equiv \langle F(\lambda, u_h), \Pi_{V_h} v \rangle_{H^{-1}(\Omega) H_0^1(\Omega)} + b(u_h, v - \Pi_{V_h} v) = 0$$

or equivalently

$$\int_{\Omega} \mathbf{grad} u_h \mathbf{grad} v \, dx - \lambda \int_{\Omega} f(u_h) \Pi_{V_h} v \, dx = 0,$$

with the elliptic projector  $\Pi_{V_h} \in \mathcal{L}(H_0^1(\Omega); V_h)$ . Then the stability assumptions (16.8) and (16.9) are immediate. Furthermore the approximability hypothesis (16.10) is a consequence of the classical polynomial interpolation results.

We can apply Theorem 16.1.

**Theorem 17.1.** *There exist  $h_0 > 0$ ,  $\delta_0 > 0$ , and for all  $0 < h \leq h_0$  a  $C^\infty$  mapping  $(\lambda_h, u_h) : s \in [-\epsilon, \epsilon] \rightarrow (\lambda_h(s), u_h(s)) \in \mathbb{R} \times H_0^1(\Omega)$  such that*

$$(17.7) \quad \left( (\lambda, u_h) \in \mathbb{R} \times V_h \text{ is a solution to (17.5), } B(u_h - u(s)) = 0, \text{ and} \right. \\ \left. (\lambda, u_h) \in \overline{B}((\lambda(s), u(s)), \delta_0) \right) \iff (\lambda = \lambda_h(s), u_h = u_h(s)).$$

Moreover there exists a constant  $C > 0$  such that for all  $s \in [-\epsilon, \epsilon]$  and  $0 < h \leq h_0$

$$(17.8) \quad |\lambda(s) - \lambda_h(s)| + |u(s) - u_h(s)|_{1,\Omega} \leq Ch^k,$$

$k = 1$  or  $2$  depending on  $V_h$ .

There exists a  $s_h \in [-\epsilon, \epsilon]$  such that the solution  $(\lambda_h(s_h), u_h(s_h))$  is a nondegenerate simple limit point of the problem (17.6) and the following estimates hold

$$(17.9) \quad |\lambda^* - \lambda_h(0)| \leq Ch^{2k}$$

and

$$(17.10) \quad |\lambda^* - \lambda_h(s_h)| \leq Ch^{2k}.$$

*Proof.* We have checked before that we can apply the results of the previous section. The estimate (17.8) is a consequence of the (16.21) one. The estimates (17.9), (17.10) are got from the (16.36), (16.37) ones. Notice that an estimate for the term

$$\left| \frac{d}{ds} \lambda(0) - \frac{d}{ds} \lambda_h(0) \right|$$

is a consequence of Theorem 12.4 via the estimate (12.14).  $\square$



**APPROXIMATION OF SIMPLE BIFURCATION POINTS**

Let  $X, Z$  be two real Banach spaces and  $F : X \rightarrow Z$  be a  $C^p$  mapping,  $p \geq 1$ . In this chapter we shall analyze approximations of the problem

$$F(x) = 0$$

in a neighborhood of a solution  $x_0 \in X$  when  $DF(x_0)$  is a not surjective Fredholm operator of index greater than 1, that is to say

- (i)  $DF(x_0) \in \mathcal{L}(X; Z)$  has a finite-dimensional kernel,
- (ii)  $\text{Range}(DF(x_0))$  is closed, different from  $Z$ ,
- (iii)  $\dim \text{Ker}(DF(x_0)) - \dim(Z/\text{Range}(DF(x_0))) \geq 1$ ,

where  $Z/\text{Range}(DF(x_0))$  denotes the quotient space of  $Z$  by  $\text{Range}(DF(x_0))$ . Such a solution  $x_0$  is called singular. Compared to the case treated in the previous chapter, where the solution  $x_0$  was supposed to be regular, that is satisfying the assumption (11.1), the study is much more involved. Our goal is to present the method in a simple way and to complete it by bibliographical comments.

In Section 18 we develop a procedure to study the exact problem. This method is the standard Lyapunov-Schmidt decomposition, leading to the bifurcation equation. The common case when  $X = \mathbb{R} \times \mathcal{X}$ ,  $\mathbb{R}$  representing the parameter space and the Banach space  $\mathcal{X}$  being the state space, is detailed in Section 19. Then approximations are analyzed in Section 20 where we get estimates between the bifurcation equation and the approximation bifurcation equations. In Section 21, we compare thoroughly the solutions of the bifurcation equation and the ones of the approximation bifurcation equations when these equations are defined in  $\mathbb{R}^2$ . A simple model example is presented in Section 22 to illustrate the general theory. Finally in Section 23 we address some further comments and bibliographical complements.

**18. Lyapunov-Schmidt procedure. Bifurcation equations**

Let  $X, Z$  be two real Banach spaces and  $F : X \rightarrow Z$  be a  $C^p$  mapping,  $p \geq 1$ . A point  $x \in X$  is called singular if

$$(18.1) \quad \begin{aligned} DF(x) \in \mathcal{L}(X; Z) \text{ is a Fredholm operator of index } m, m \in \mathbb{N}, m \geq 1, \\ \text{and } \text{Range}(DF(x_0)) \neq Z; \end{aligned}$$

the index of the operator  $DF(x)$  is

$$\text{index}(DF(x)) = \dim \text{Ker}(DF(x)) - \dim \text{Coker}(DF(x)) = m,$$

where  $\text{Coker}(DF(x))$  is the quotient space  $Z/\text{Range}(DF(x))$ , and the codimension of  $\text{Range}(DF(x))$  (the dimension of the cokernel) is equal to  $k \geq 1$ .

In this section, we will present a method to analyze the solution set  $S$  of the problem

$$(18.2) \quad F(x) = 0,$$

in a neighborhood of the singular point  $x_0 \in X$  satisfying  $F(x_0) = 0$ . Such a point  $x_0$  is called a singular solution.

Let  $m$  be the index of  $DF(x_0)$  and  $k$  be the dimension of  $\text{Coker}(DF(x_0))$ .

Since the dimension of  $\text{Ker}(DF(x_0))$  is finite this subspace is direct, that is there exists a closed subspace  $X_1 \subset X$  such that

$$X = X_1 \oplus \text{Ker}(DF(x_0)).$$

The range of  $DF(x_0) \subset Z$  is of finite codimension  $k$ , so this subspace is direct. There exists a subspace  $Z_1 \subset Z$ , of dimension  $k$ , such that

$$Z = Z_1 \oplus \text{Range}(DF(x_0)).$$

Subsequently in this chapter we shall use the following notations. A point  $x \in X$  has a unique representation of the form

$$x = x_1 + x_2 \quad \text{with } x_1 \in X_1 \text{ and } x_2 \in \text{Ker}(DF(x_0))$$

and  $P_1 \in \mathcal{L}(X; X_1)$ ,  $P_2 \in \mathcal{L}(X; \text{Ker}(DF(x_0)))$  denote the projectors associated to the decomposition of  $X$ , given by

$$P_1(x) = x_1, \quad P_2(x) = x_2.$$

In a similar way a point  $z \in Z$  has a unique representation of the form

$$z = z_1 + z_2 \quad \text{with } z_1 \in Z_1 \text{ and } z_2 \in \text{Range}(DF(x_0))$$

and  $Q_1 \in \mathcal{L}(Z; Z_1)$ ,  $Q_2 \in \mathcal{L}(Z; \text{Range}(DF(x_0)))$  denote the projectors associated to the decomposition of  $Z$ , given by

$$Q_1(z) = z_1, \quad Q_2(z) = z_2.$$

With these notations problem (18.2) can be written in the equivalent form

$$Q_1 F(x) = 0 \quad \text{and} \quad Q_2 F(x) = 0.$$

Since we study problem (18.2) in a neighborhood of the singular solution  $x_0$ , we write it in the following equivalent way: find  $x_1 \in X_1$ ,  $x_2 \in \text{Ker}(DF(x_0))$  such that

$$(18.3) \quad Q_1 F(x_0 + x_1 + x_2) = 0$$

and

$$(18.4) \quad Q_2 F(x_0 + x_1 + x_2) = 0.$$

Let us look first at the equation (18.4). The mapping  $\mathcal{F} : X_1 \times \text{Ker}(DF(x_0)) \rightarrow \text{Range}(DF(x_0))$  is defined for  $(x_1, x_2) \in X_1 \times \text{Ker}(DF(x_0))$  by

$$(18.5) \quad \mathcal{F}(x_1, x_2) = Q_2 F(x_0 + x_1 + x_2).$$

We have  $\mathcal{F}(0, 0) = 0$  and the derivative  $D_{x_1} \mathcal{F}(0, 0) \in \mathcal{L}(X_1; \text{Range}(DF(x_0)))$  is equal to  $Q_2 DF(x_0)$ , which is an isomorphism from  $X_1$  onto  $\text{Range}(DF(x_0))$ .

We apply the implicit function theorem 2.3. There exist positive constants  $\epsilon$ ,  $\eta$ , and a mapping  $g : B(0, \epsilon) \subset \text{Ker}(DF(x_0)) \rightarrow X_1$  which is of class  $C^p$ , such that for all  $x_2 \in B(0, \epsilon)$

$$(\mathcal{F}(x_1, x_2) = 0 \quad \text{and} \quad x_1 \in B(0, \eta) \subset X_1) \iff x_1 = g(x_2).$$

In a neighborhood of  $x_0$ , the solution set to  $Q_2 F(x) = 0$  is the set

$$S_{\mathcal{F}} = \{x \in X; x = x_0 + g(x_2) + x_2, x_2 \in B(0, \epsilon)\},$$

so

$$(18.6) \quad Q_2 F(x_0 + g(x_2) + x_2) = 0 \quad \text{for all } x_2 \in B(0, \epsilon).$$

To complete our study, we need to analyze the equation (18.3) on the solution set  $S_{\mathcal{F}}$ . The mapping  $\mathcal{G} : B(0, \epsilon) \subset \text{Ker}(DF(x_0)) \rightarrow Z_1$  is defined for  $x_2 \in B(0, \epsilon)$  by

$$(18.7) \quad \mathcal{G}(x_2) = Q_1 F(x_0 + g(x_2) + x_2).$$

Notice that the mapping  $\mathcal{G}$  can be identified to a mapping defined in  $\mathbb{R}^{k+m}$  with values in  $\mathbb{R}^k$ . Indeed let  $\zeta_1, \dots, \zeta_k$  be a basis of  $Z_1$  and  $\varphi_1, \dots, \varphi_{k+m}$  be a basis of  $\text{Ker}(DF(x_0))$ . Then we can consider for a convenient  $\delta > 0$  the function  $\tilde{\mathcal{G}} : B(0, \delta) \subset \mathbb{R}^{k+m} \rightarrow \mathbb{R}^k$  given by: for  $\alpha = (\alpha_1, \dots, \alpha_{k+m}) \in B(0, \delta)$

$$\tilde{\mathcal{G}}(\alpha) = \left( \mathcal{G}_1 \left( \sum_{j=1}^{k+m} \alpha_j \varphi_j \right), \dots, \mathcal{G}_k \left( \sum_{j=1}^{k+m} \alpha_j \varphi_j \right) \right)$$

where

$$\mathcal{G}(x_2) = \sum_{\ell=1}^k \mathcal{G}_\ell(x_2)\zeta_\ell.$$

Therefore solving problem (18.2) in a neighborhood of the singular solution  $x_0$  has been reduced to solving in a neighborhood of the origin the equation

$$(18.8) \quad \mathcal{G}(x_2) = 0$$

or the finite nonlinear system of  $k$  equations with  $k + m$  unknowns

$$(18.8') \quad \tilde{\mathcal{G}}(\boldsymbol{\alpha}) = 0.$$

The equation (18.8) (or (18.8')) is called the bifurcation equation of problem (18.2) at  $x_0$ . The method we have used to get the bifurcation equation is called the Lyapunov-Schmidt procedure.

*Remark 18.1.* There is a different way to derive the bifurcation equation, as proposed in CROUZEIX and RAPPAZ [1989]. The idea is similar to the one we have used to study a solution set containing a simple limit point. It consists in transforming the problem to reduce it to a simpler one.

Let the mapping  $\Phi : \mathbb{R}^k \times X \rightarrow Z$  be given by: for  $(\boldsymbol{\mu}, x) \in \mathbb{R}^k \times X$

$$\Phi(\boldsymbol{\mu}, x) = F(x) + \sum_{j=1}^k \mu_j \zeta_j.$$

Recall that the elements  $\zeta_1, \dots, \zeta_k$  form a basis of  $Z_1$ . Then finding the solutions of problem (18.2) in a neighborhood of  $x_0$  is equivalent to find the solutions  $(\boldsymbol{\mu}, x) \in \mathbb{R}^k \times X$  of

$$(18.9) \quad \Phi(\boldsymbol{\mu}, x) = 0$$

in a neighborhood of  $(0, x_0)$  satisfying  $\boldsymbol{\mu} = 0$ . Then problem (18.9) is solved in two steps. First we look for the solutions of (18.9) in a neighborhood of  $(0, x_0)$ . Secondly we determine all the solutions  $(\boldsymbol{\mu}, x)$  with  $\boldsymbol{\mu} = 0$ .

Remark that  $D\Phi(0, x_0) \in \mathcal{L}(\mathbb{R}^k \times X; Z)$  is a Fredholm operator with index  $k + m$ . Indeed  $\text{Ker}(D\Phi(0, x_0)) = \{(0, x) \in \mathbb{R}^k \times X; x \in \text{Ker}(DF(x_0))\}$ . Moreover the range of  $D\Phi(0, x_0)$  is equal to  $Z$ . We have a problem studied in Section 15 when  $k = 1$ . The case  $k \geq 2$  can be treated in a similar way.

In the second step we look at the solutions of (18.9) with  $\boldsymbol{\mu} = 0$ , which is precisely the bifurcation equation.  $\square$

*Remark 18.2.* Since  $F(x_0) = 0$ , we easily check that  $g(0) = 0$  and then  $\mathcal{G}(0) = 0$ . Moreover when differentiating (18.7) with respect to  $x_2$ , we get  $D\mathcal{G}(0) = 0$ . Hence problem (18.8) is singular. In the next section we will study the bifurcation equation more carefully.  $\square$

### 19. The case with parameter and state spaces

In this section, we analyze more carefully the bifurcation equation deduced from the Lyapunov-Schmidt procedure when the  $C^p$  mapping  $F$ ,  $p \geq 2$ , is

$$F : (\lambda, \kappa) \in \mathbb{R} \times \mathcal{X} \rightarrow F(\lambda, \kappa) \in Z$$

where  $\mathbb{R}$  represents the parameter space while  $\mathcal{X}$  is the state space.

Let  $(\lambda_0, \kappa_0) \in \mathbb{R} \times \mathcal{X}$  be a singular solution of

$$(19.1) \quad F(\lambda, \kappa) = 0$$

such that

$$(19.2) \quad \begin{cases} \text{(i)} & \dim \text{Ker}(D_\kappa F(\lambda_0, \kappa_0)) = 1, \\ \text{(ii)} & \text{codim}(D_\kappa F(\lambda_0, \kappa_0)) = 1, \\ \text{(iii)} & D_\lambda F(\lambda_0, \kappa_0) = 0. \end{cases}$$

The last assumption in (19.2) is satisfied when we study bifurcation from the trivial branch of solutions  $\{(\lambda, 0); \lambda \in \mathbb{R}\}$ , when  $F$  satisfies  $F(\lambda, 0) = 0$  for all  $\lambda \in \mathbb{R}$ . At the end of the section, we shall extend our results with weaker assumptions than (19.2).

Under the assumptions (19.2), we can write the decomposition of both spaces  $X$  and  $Z$  in a more detailed way than in the previous section. Notice that the derivative  $DF(\lambda_0, \kappa_0) \in \mathcal{L}(\mathbb{R} \times \mathcal{X}; Z)$  satisfies: for all  $(\delta, w) \in \mathbb{R} \times \mathcal{X}$

$$DF(\lambda_0, \kappa_0)(\delta, w) = D_\kappa F(\lambda_0, \kappa_0)w.$$

Since the kernel of  $D_\kappa F(\lambda_0, \kappa_0)$  is of dimension 1, there is a closed subspace  $\mathcal{X}_1 \subset \mathcal{X}$  and  $\xi \in \text{Ker}(D_\kappa F(\lambda_0, \kappa_0))$ ,  $\xi \neq 0$ , such that

$$(19.3) \quad \begin{aligned} \mathcal{X} &= \mathcal{X}_1 \oplus \text{Ker}(D_\kappa F(\lambda_0, \kappa_0)) \\ &\equiv \mathcal{X}_1 \oplus \text{span}\{\xi\}. \end{aligned}$$

From the assumption (19.2), we deduce that the kernel of  $DF(\lambda_0, \kappa_0)$  is of dimension 2 and that

$$\text{Ker}(DF(\lambda_0, \kappa_0)) = \mathbb{R} \times \text{Ker}(D_\kappa F(\lambda_0, \kappa_0)).$$

Hence we can write a decomposition of  $X$

$$(19.4) \quad \begin{aligned} X &= X_1 \oplus \text{Ker}(DF(\lambda_0, \kappa_0)) \\ &\equiv X_1 \oplus \text{span}\{\varphi_1, \varphi_2\}, \end{aligned}$$

where  $X_1 = \{0\} \times \mathcal{X}_1$ ,  $\varphi_1 = (1, 0) \in \mathbb{R} \times \mathcal{X}$ , and  $\varphi_2 = (0, \xi) \in \mathbb{R} \times \mathcal{X}$ .

Since the range of  $DF(\lambda_0, \kappa_0)$  is equal to the range of  $D_\kappa F(\lambda_0, \kappa_0)$  and the codimension of  $D_\kappa F(\lambda_0, \kappa_0)$  is 1, the following decomposition of  $Z$  holds

$$(19.5) \quad \begin{aligned} Z &= Z_1 \oplus \text{Range}(D_\kappa F(\lambda_0, \kappa_0)) \\ &\equiv \text{span}\{\zeta_1\} \oplus \text{Range}(D_\kappa F(\lambda_0, \kappa_0)), \end{aligned}$$

where  $\zeta_1$  does not belong to the range of  $D_\kappa F(\lambda_0, \kappa_0)$ . Let  $P_1, P_2$  denote the two projectors associated to the decomposition (19.4) of  $X$  and  $Q_1, Q_2$  associated to the decomposition (19.5) of  $Z$ .

We can apply the Lyapunov-Schmidt procedure given in Section 18. Here we have  $k = m = 1$ . There exist positive constants  $\epsilon, \eta$ , and a mapping

$$g : (\alpha_1, w) \in B(0, \epsilon) \subset \mathbb{R} \times \text{Ker}(D_\kappa F(\lambda_0, \kappa_0)) \rightarrow g(\alpha_1, w) \in \mathcal{X}_1$$

of class  $C^p$  such that for all  $(\alpha_1, w) \in B(0, \epsilon)$

$$(Q_2 F(\lambda_0 + \alpha_1, \kappa_0 + \kappa_1 + w) = 0 \text{ and } \kappa_1 \in B(0, \eta) \subset \mathcal{X}_1) \iff \kappa_1 = g(\alpha_1, w).$$

In fact  $w = \alpha_2 \xi$  and  $g$  can be identified to a function defined on  $B(0, \delta) \subset \mathbb{R}^2$ ,  $\delta$  small enough, by setting

$$g(\alpha_1, \alpha_2) = g(\alpha_1, \alpha_2 \xi).$$

Then  $g$  satisfies for all  $(\alpha_1, \alpha_2) \in B(0, \delta)$

$$(19.6) \quad Q_2 F(\lambda_0 + \alpha_1, \kappa_0 + g(\alpha_1, \alpha_2) + \alpha_2 \xi) = 0.$$

In particular  $g(0, 0) = 0$ . Let the functional  $B$  in  $Z'$  satisfy

$$B(\zeta_1) = 1, \quad B(z) = 0 \quad \text{for all } z \in \text{Range}(DF(\lambda_0, \kappa_0)).$$

We can write then the bifurcation equation in the following way

$$(19.7) \quad f(\alpha_1, \alpha_2) \equiv B(F(\lambda_0 + \alpha_1, \kappa_0 + g(\alpha_1, \alpha_2) + \alpha_2 \xi)) = 0,$$

where  $f : B(0, \delta) \subset \mathbb{R}^2 \rightarrow \mathbb{R}$  is of class  $C^p$ . We easily check that

$$f(0, 0) = 0, \quad \frac{\partial f}{\partial \alpha_1}(0, 0) = 0, \quad \text{and} \quad \frac{\partial f}{\partial \alpha_2}(0, 0) = 0,$$

which are relations we expect for a bifurcation equation, see Remark 18.2. To study the solutions of the bifurcation equation in a neighborhood of the origin, we need to compute the second derivative of  $f$  with respect to  $\alpha_1$  and  $\alpha_2$ , the Hessian matrix  $D^2 f(0, 0)$ . When deriving the equation (19.6) with respect to  $\alpha_1$  and  $\alpha_2$ , at  $(0, 0)$ , we get with (19.2) (iii)

$$D_\kappa F(\lambda_0, \kappa_0) \frac{\partial g}{\partial \alpha_1}(0, 0) = 0, \quad D_\kappa F(\lambda_0, \kappa_0) \left( \frac{\partial g}{\partial \alpha_2}(0, 0) + \xi \right) = 0,$$

which implies

$$(19.8) \quad \frac{\partial g}{\partial \alpha_1}(0, 0) = 0 = \frac{\partial g}{\partial \alpha_2}(0, 0).$$

Then a simple calculation gives

$$\begin{aligned} \frac{\partial^2 f}{\partial \alpha_1^2}(0, 0) &= B(D_{\lambda\lambda}^2 F(\lambda_0, \kappa_0)), \\ \frac{\partial^2 f}{\partial \alpha_1 \partial \alpha_2}(0, 0) &= B(D_{\lambda\kappa}^2 F(\lambda_0, \kappa_0)\xi), \\ \frac{\partial^2 f}{\partial \alpha_2^2}(0, 0) &= B(D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\xi, \xi)). \end{aligned}$$

The determinant of the Hessian matrix  $D^2 f(0, 0)$  is

$$(19.9) \quad \det(D^2 f(0, 0)) = B(D_{\lambda\lambda}^2 F(\lambda_0, \kappa_0))B(D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\xi, \xi)) \\ - (B(D_{\lambda\kappa}^2 F(\lambda_0, \kappa_0)\xi))^2.$$

The bifurcation equation is commonly studied with the Morse lemma. In fact the results are strongly depending on the sign of the determinant of  $D^2 f(0, 0)$ .

**Theorem 19.1.** (*Morse Lemma*). *Let  $H$  be a Hilbert space with the scalar product  $(\cdot, \cdot)_H$  and the norm  $\|\cdot\|_H$ . Let  $f : V \subset H \rightarrow \mathbb{R}$  be a mapping of class  $C^1$  defined on a neighborhood  $V$  of the origin. We assume that  $D^2 f(0)$  exists. Let the continuous symmetric linear operator  $S : H \rightarrow H$  be defined by: for all  $u, v$  in  $H$*

$$(Su, v)_H = D^2 f(0)(u, v).$$

We assume that  $Df(0) = 0$  and that  $S$  is an isomorphism onto  $H$ .

Then there exist two neighborhoods  $\mathcal{V}, \mathcal{W}$  of the origin in  $H$ , and a  $C^1$  diffeomorphism  $\Phi : \mathcal{V} \subset V \rightarrow \mathcal{W}$  satisfying

$$\Phi(0) = 0, \quad D\Phi(0) = I,$$

and for all  $v \in \mathcal{V}$

$$f(v) = f(0) + \frac{1}{2}D^2 f(0)(\Phi(v), \Phi(v)).$$

Moreover if  $f$  is in  $C^m(V; \mathbb{R})$  with  $m \geq 1$  and if  $D^{m+1} f(0)$  exists, then  $\Phi$  is of class  $C^m$ .  $\square$

We can find proofs of the Morse lemma for instance in NIRENBERG [1974] or in CROUZEIX and RAPPAZ [1989]. We do not repeat the proof here.

We assume from now on that the determinant of the Hessian matrix  $D^2 f(0,0)$  is different from 0 and apply the Morse lemma to the bifurcation equation (19.7). There exist two neighborhoods  $\mathcal{V}$ ,  $\mathcal{W}$  of the origin in  $\mathbb{R}^2$  and a  $C^1$  diffeomorphism  $\Phi : \mathcal{V} \subset B(0, \delta) \subset \mathbb{R}^2 \rightarrow \mathcal{W} \subset \mathbb{R}^2$  satisfying

$$\Phi(0,0) = 0, \quad D\Phi(0,0) = I,$$

and for all  $(\alpha_1, \alpha_2) \in \mathcal{V}$

$$f(\alpha_1, \alpha_2) = \frac{1}{2} (D^2 f(0,0)\Phi(\alpha_1, \alpha_2), \Phi(\alpha_1, \alpha_2)),$$

here  $(\cdot, \cdot)$  is the scalar product in  $\mathbb{R}^2$ .

We analyze now the two different cases which can occur depending on the sign of the determinant of  $D^2 f(0,0)$ . First let us assume that

$$(19.10) \quad \det(D^2 f(0,0)) > 0,$$

which is the elliptic case. Since the matrix  $D^2 f(0,0)$  is symmetric, there exists a orthogonal matrix  $Q$  such that

$$Q^t DQ = D^2 f(0,0),$$

where the diagonal matrix  $D$  has the diagonal elements  $d_1$  and  $d_2$ ,  $d_i$  eigenvalue of  $D^2 f(0,0)$ . Let the vector  $\Psi$  of components  $\Psi_1, \Psi_2$ , be

$$\Psi(\alpha_1, \alpha_2) = Q\Phi(\alpha_1, \alpha_2),$$

then for  $(\alpha_1, \alpha_2) \in \mathcal{V}$

$$f(\alpha_1, \alpha_2) = \frac{d_1}{2} \Psi_1^2(\alpha_1, \alpha_2) + \frac{d_2}{2} \Psi_2^2(\alpha_1, \alpha_2).$$

Then the bifurcation equation reads

$$d_1 \Psi_1^2(\alpha_1, \alpha_2) + d_2 \Psi_2^2(\alpha_1, \alpha_2) = 0,$$

where  $d_1$  and  $d_2$  have the same sign. Since we have

$$\Psi(0,0) = 0 \quad \text{and} \quad D\Psi(0,0) \text{ regular,}$$

so in a neighborhood of  $(0,0)$ , the bifurcation equation has the only solution  $(0,0)$  and the solution set

$$S = \{(\lambda, \kappa) \in \mathbb{R} \times \mathcal{X}; F(\lambda, \kappa) = 0\}$$



is reduced to  $\{(\lambda_0, \kappa_0)\}$  in a neighborhood of  $(\lambda_0, \kappa_0)$ . Such a solution  $(\lambda_0, \kappa_0)$  to problem (19.1) is called an isolated solution.

We assume now that

$$(19.11) \quad \det(D^2 f(0, 0)) < 0,$$

which is the hyperbolic case. Since the matrix  $Df(0, 0)$  is regular and symmetric, there exists a regular matrix  $R$  such that

$$R^t \begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix} R = D^2 f(0, 0).$$

Let the vector  $\Psi$  of components  $\Psi_1, \Psi_2$ , be

$$\Psi(\alpha_1, \alpha_2) = R\Phi(\alpha_1, \alpha_2),$$

then the bifurcation equation reads: for all  $(\alpha_1, \alpha_2) \in \mathcal{V}$

$$(19.12) \quad f(\alpha_1, \alpha_2) \equiv \Psi_1(\alpha_1, \alpha_2)\Psi_2(\alpha_1, \alpha_2) = 0.$$

Notice that

$$(19.13) \quad \Psi(0, 0) = 0 \quad \text{and} \quad D\Psi(0, 0) \text{ is regular.}$$

From (19.12), (19.13), and the Morse lemma, we deduce that the solution set of  $f$

$$\mathcal{S}_f (\equiv \{(\alpha_1, \alpha_2) \in B(0, \delta); f(\alpha_1, \alpha_2) = 0\}) \cap \mathcal{V}$$

is  $C^{p-1}$  diffeomorphic to two parts of straight lines (in the coordinates  $\Psi_1, \Psi_2$ ) which intersect at  $(0, 0) \in \mathbb{R}^2$ . Consequently in a neighborhood of  $(\lambda_0, \kappa_0)$  the solution set  $S$  of (19.1) consists in two curves which intersect at the point  $(\lambda_0, \kappa_0)$ .

More precisely let  $\alpha(\Psi)$  denote the inverse of the  $C^{p-1}$  diffeomorphism  $\Psi(\alpha)$  defined for  $\Psi$  in  $R\mathcal{W}$ , say for  $|\Psi_1| + |\Psi_2| \leq \epsilon$  for simplicity. Then

$$\mathcal{S}_f = \{(\alpha_1(s, 0), \alpha_2(s, 0)); |s| \leq \epsilon\} \cup \{(\alpha_1(0, s), \alpha_2(0, s)); |s| \leq \epsilon\}.$$

For the parameter value  $s$ , with  $|s| \leq \epsilon$ , we set

$$\begin{aligned} \lambda_1(s) &= \lambda_0 + \alpha_1(s, 0), & \kappa_1(s) &= \kappa_0 + g(\alpha_1(s, 0), \alpha_2(s, 0)) + \alpha_2(s, 0)\xi, \\ \lambda_2(s) &= \lambda_0 + \alpha_1(0, s), & \kappa_2(s) &= \kappa_0 + g(\alpha_1(0, s), \alpha_2(0, s)) + \alpha_2(0, s)\xi. \end{aligned}$$

Clearly for  $i = 1$  or  $2$  and  $|s| \leq \epsilon$

$$\begin{aligned} F(\lambda_i(s), \kappa_i(s)) &= 0, \\ \lambda_i(0) = \lambda_0, \kappa_i(0) = \kappa_0 & \quad \text{and if } s \neq 0 \quad (\lambda_i(s), \kappa_i(s)) \neq (\lambda_0, \kappa_0). \end{aligned}$$

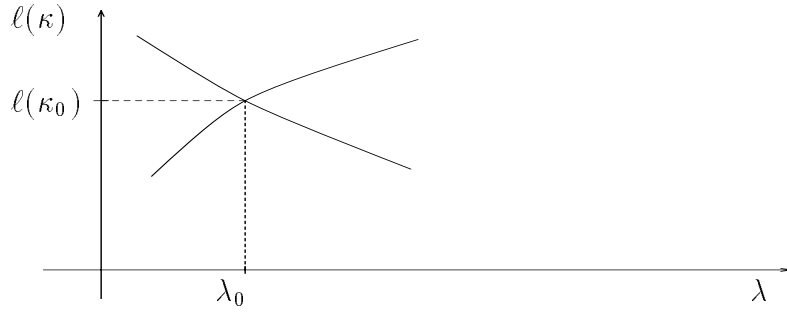


FIGURE 19.1: SIMPLE BIFURCATION POINT.

Moreover the two branches of solutions intersect at  $(\lambda_0, \kappa_0)$ . If  $\ell$  is a linear form on  $\mathcal{X}$ , then in a neighborhood of  $(\lambda_0, \kappa_0)$  the solution set of (19.1) can be represented as in the figure 19.1.

Such a solution  $(\lambda_0, \kappa_0)$  of (19.1) is called a simple bifurcation point.

Before studying bifurcation approximations of simple bifurcation points, we end up the section with some remarks.

*Remark 19.1.* If the mapping  $F : \mathbb{R} \times \mathcal{X} \rightarrow Z$  satisfies

$$F(\lambda, 0) = 0 \quad \text{for all } \lambda \in \mathbb{R}$$

and if  $(\lambda_0, 0)$  satisfies (19.2) (i) and (ii), then we have

$$D_\lambda F(\lambda_0, 0) = 0 \quad \text{and} \quad D_{\lambda\lambda}^2 F(\lambda_0, 0) = 0.$$

Clearly the assumption (19.2) (iii) is satisfied. The determinant of the Hessian matrix  $D^2 f(0, 0)$  given in (19.9) reads

$$\det(D^2 f(0, 0)) = -(B(D_{\lambda\kappa}^2 F(\lambda_0, 0)\xi))^2.$$

So  $(\lambda_0, 0)$  is a simple bifurcation point if and only if  $B(D_{\lambda\kappa}^2 F(\lambda_0, 0)\xi) \neq 0$  or in other words

$$D_{\lambda\kappa}^2 F(\lambda_0, 0)\xi \notin \text{Range}(D_\kappa F(\lambda_0, 0)),$$

which is the Crandall-Rabinowitz condition, see CRANDALL and RABINOWITZ [1973]. Notice that approximations of simple bifurcation points on the trivial branch  $\{(\lambda, 0); \lambda \in \mathbb{R}\}$  is studied in CROUZEIX and RAPPAZ [1989], Chapter 5, where they use a simpler approach than the Lyapunov-Schmidt procedure.  $\square$

*Remark 19.2.* If we replace the assumption (19.2) (iii), namely  $D_\lambda F(\lambda_0, \kappa_0) = 0$ , by the weaker one

$$(19.14) \quad D_\lambda F(\lambda_0, \kappa_0) \in \text{Range}(D_\kappa F(\lambda_0, \kappa_0)),$$

then we can use the above development with minor changes and get a similar conclusion.

Indeed the kernel of  $DF(\lambda_0, \kappa_0)$  is of dimension 2 and

$$\text{Ker}(DF(\lambda_0, \kappa_0)) = \text{span}\{\varphi_1, \varphi_2\}$$

with  $\varphi_1 = (1, \eta) \in \mathbb{R} \times \mathcal{X}_1$ ,  $\eta$  the unique element satisfying

$$D_\lambda F(\lambda_0, \kappa_0) + D_\kappa F(\lambda_0, \kappa_0)\eta = 0,$$

and with  $\varphi_2 = (0, \xi) \in \mathbb{R} \times \mathcal{X}$ ,  $\xi \neq 0$  and

$$D_\kappa F(\lambda_0, \kappa_0)\xi = 0.$$

With that choice of  $\varphi_1$  and  $\varphi_2$ , the decomposition (19.4) for  $X$  still remains valid, that is

$$(19.15) \quad X = X_1 \oplus \text{span}\{\varphi_1, \varphi_2\}$$

where  $X_1 = \{0\} \times \mathcal{X}_1$ . Then the equation (19.6) is modified in

$$(19.16) \quad Q_2 F(\lambda_0 + \alpha_1, \kappa_0 + g(\alpha_1, \alpha_2) + \alpha_1 \eta + \alpha_2 \xi) = 0$$

and the bifurcation equation (19.7) in

$$(19.17) \quad f(\alpha_1, \alpha_2) \equiv B(F(\lambda_0 + \alpha_1, \kappa_0 + g(\alpha_1, \alpha_2) + \alpha_1 \eta + \alpha_2 \xi)) = 0.$$

When deriving the equation (19.16) with respect to  $\alpha_1$  and  $\alpha_2$ , at  $(0, 0)$ , we get

$$\begin{aligned} D_\lambda F(\lambda_0, \kappa_0) + D_\kappa F(\lambda_0, \kappa_0) \left( \frac{\partial g}{\partial \alpha_1}(0, 0) + \eta \right) &= 0, \\ D_\kappa F(\lambda_0, \kappa_0) \left( \frac{\partial g}{\partial \alpha_2}(0, 0) + \xi \right) &= 0, \end{aligned}$$

which implies

$$\frac{\partial g}{\partial \alpha_1}(0, 0) = 0 = \frac{\partial g}{\partial \alpha_2}(0, 0).$$

Then a simple calculation gives the following expressions for the elements of the Hessian matrix  $D^2 f(0, 0)$

$$\begin{aligned} \frac{\partial^2 f}{\partial \alpha_1^2}(0, 0) &= B(D_{\lambda\lambda}^2 F(\lambda_0, \kappa_0) + 2D_{\lambda\kappa}^2 F(\lambda_0, \kappa_0)\eta + D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\eta, \eta)), \\ \frac{\partial^2 f}{\partial \alpha_1 \partial \alpha_2}(0, 0) &= B(D_{\lambda\kappa}^2 F(\lambda_0, \kappa_0)\xi + D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\eta, \xi)), \\ \frac{\partial^2 f}{\partial \alpha_2^2}(0, 0) &= B(D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\xi, \xi)). \end{aligned}$$

Then the study of the bifurcation equation goes along the same line as before, depending on the sign of the determinant of  $D^2 f(0, 0)$ .  $\square$

*Remark 19.3.* A similar development can be carried out if we replace the assumption (19.2) by

$$(19.18) \quad \begin{cases} \text{(i)} & \dim \text{Ker}(D_\kappa F(\lambda_0, \kappa_0)) = 2, \\ \text{(ii)} & \text{codim}(D_\kappa F(\lambda_0, \kappa_0)) = 2, \\ \text{(iii)} & D_\lambda F(\lambda_0, \kappa_0) \notin \text{Range}(D_\kappa F(\lambda_0, \kappa_0)). \end{cases}$$

Indeed the kernel of  $DF(\lambda_0, \kappa_0)$  is of dimension 2 and

$$\text{Ker}(DF(\lambda_0, \kappa_0)) = \text{span}\{\varphi_1, \varphi_2\}$$

with  $\varphi_1 = (0, \xi_1) \in \mathbb{R} \times \mathcal{X}$  and  $\varphi_2 = (0, \xi_2) \in \mathbb{R} \times \mathcal{X}$  such that

$$\text{Ker}(D_\kappa F(\lambda_0, \kappa_0)) = \text{span}\{\xi_1, \xi_2\}.$$

There exists a closed subspace  $\mathcal{X}_1 \subset \mathcal{X}$  such that

$$\mathcal{X} = \mathcal{X}_1 \oplus \text{span}\{\xi_1, \xi_2\}$$

and we can write a decomposition of  $X$

$$(19.19) \quad X = X_1 \oplus \text{span}\{\varphi_1, \varphi_2\}$$

with  $X_1 = \mathbb{R} \times \mathcal{X}_1$ . The codimension of  $DF(\lambda_0, \kappa_0)$  is 1, so the following decomposition of  $Z$  holds

$$(19.20) \quad Z = \text{span}\{\zeta_1\} \oplus \text{Range}(DF(\lambda_0, \kappa_0)).$$

The mapping  $g$  can be identified to a function defined on  $B(0, \delta) \subset \mathbb{R}^2$  with values in  $\mathbb{R} \times \mathcal{X}_1$ , that is

$$g : (\alpha_1, \alpha_2) \in B(0, \delta) \subset \mathbb{R}^2 \rightarrow (g_1(\alpha_1, \alpha_2), g_2(\alpha_1, \alpha_2)) \in \mathbb{R} \times \mathcal{X}_1,$$

and satisfies for all  $(\alpha_1, \alpha_2) \in B(0, \delta)$

$$(19.21) \quad Q_2 F(\lambda_0 + g_1(\alpha_1, \alpha_2), \kappa_0 + g_2(\alpha_1, \alpha_2) + \alpha_1 \xi_1 + \alpha_2 \xi_2) = 0.$$

The bifurcation equation reads

$$(19.22) \quad f(\alpha_1, \alpha_2) \equiv B(F(\lambda_0 + g_1(\alpha_1, \alpha_2), \kappa_0 + g_2(\alpha_1, \alpha_2) + \alpha_1 \xi_1 + \alpha_2 \xi_2)) = 0$$

where  $f : B(0, \delta) \subset \mathbb{R}^2 \rightarrow \mathbb{R}$  is of class  $C^p$ .

When deriving the equation (19.21) with respect to  $\alpha_1$  and  $\alpha_2$ , at  $(0,0)$ , we get

$$\begin{aligned} D_\lambda F(\lambda_0, \kappa_0) \frac{\partial g_1}{\partial \alpha_1}(0,0) + D_\kappa F(\lambda_0, \kappa_0) \left( \frac{\partial g_2}{\partial \alpha_1}(0,0) + \xi_1 \right) &= 0, \\ D_\lambda F(\lambda_0, \kappa_0) \frac{\partial g_1}{\partial \alpha_2}(0,0) + D_\kappa F(\lambda_0, \kappa_0) \left( \frac{\partial g_2}{\partial \alpha_2}(0,0) + \xi_2 \right) &= 0. \end{aligned}$$

Since

$$\text{Range}(DF(\lambda_0, \kappa_0)) = \text{span}\{D_\lambda F(\lambda_0, \kappa_0)\} \oplus \text{Range}(D_\kappa F(\lambda_0, \kappa_0))$$

and since  $\xi_1$  and  $\xi_2$  are in the kernel of  $D_\kappa F(\lambda_0, \kappa_0)$  we deduce that

$$\frac{\partial g_i}{\partial \alpha_j}(0,0) = 0, \quad 1 \leq i, j \leq 2.$$

Then a simple calculation gives the following expressions

$$\begin{aligned} \frac{\partial^2 f}{\partial \alpha_1^2}(0,0) &= B(D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\xi_1, \xi_1)), \\ \frac{\partial^2 f}{\partial \alpha_1 \partial \alpha_2}(0,0) &= B(D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\xi_1, \xi_2)), \\ \frac{\partial^2 f}{\partial \alpha_2^2}(0,0) &= B(D_{\kappa\kappa}^2 F(\lambda_0, \kappa_0)(\xi_2, \xi_2)). \end{aligned}$$

The study of the bifurcation equation goes along the same line as before, depending on the sign of the determinant of  $D^2 f(0,0)$ .

In the hyperbolic case, the solution set of  $F(\lambda, x) = 0$  in a neighborhood of  $(\lambda_0, \kappa_0)$  is represented in the figure 19.2, where  $\ell$  is a linear form on  $\mathcal{X}$ .

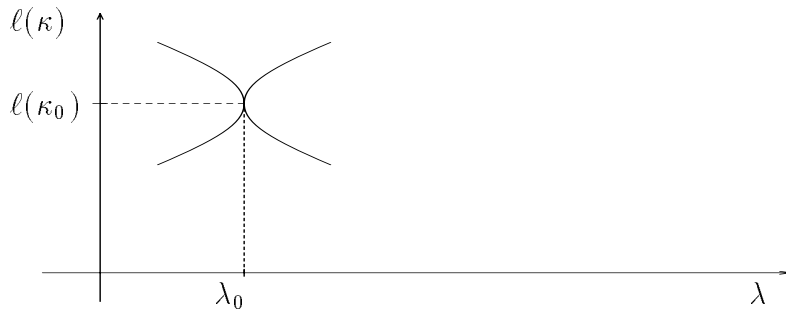


FIGURE 19.2: DOUBLE NONDEGENERATE LIMIT POINT.

Such a solution  $(\lambda_0, \kappa_0)$  of  $F(\lambda, \kappa) = 0$  is called a double nondegenerate limit point.  $\square$

## 20. Approximate bifurcation equations

Let  $X, Z$  be two real Banach spaces,  $F : X \rightarrow Z$  be a  $C^p$  mapping,  $p \geq 1$ , and  $x_0 \in X$  be a singular solution of the problem: find  $x \in X$  such that

$$(20.1) \quad F(x) = 0.$$

We assume that the solution  $x_0$  satisfies

$$(20.2) \quad DF(x_0) \in \mathcal{L}(X; Z) \text{ is a not surjective Fredholm operator of index } m, m \geq 1.$$

The study of the solution set  $S$  of (20.1) in a neighborhood of  $x_0$  has been done in Sections 18 and 19. We will keep the notations introduced there. Let  $P_1, P_2$  and  $Q_1, Q_2$  be the projectors associated to the decomposition of  $X$  and  $Z$ ,

$$X = X_1 \oplus \text{Ker}(DF(x_0)), \quad Z = Z_1 \oplus \text{Range}(DF(x_0)).$$

We introduce now a family  $\{F_h\}_{0 < h \leq 1}$  of  $C^p$  mappings,  $p \geq 1$ ,

$$F_h : x \in X \rightarrow F_h(x) \in Z$$

which are approximations of  $F$ . Our goal is to study the existence and the convergence of solutions of the problem: find  $x \in X$  such that

$$(20.3) \quad F_h(x) = 0$$

in a neighborhood of the singular solution  $x_0$ . We assume that the family  $\{F_h\}_{0 < h \leq 1}$  satisfies: there exists a real number  $r > 0$  such that

$$(20.4) \quad \lim_{h \rightarrow 0} \sup_{x \in B(x_0, r)} \|D^\ell F(x) - D^\ell F_h(x)\| = 0,$$

for  $\ell = 0, 1$ .

To study problem (20.3) under the assumption (20.4), we will apply the Lyapunov-Schmidt procedure to the mapping  $F_h$  with respect to the decomposition of  $X$  and  $Z$  given in the analysis of problem (20.1). Since we study problem (20.3) in a neighborhood of  $x_0$ , we write it in the following way: find  $x_1 \in X_1, x_2 \in \text{Ker}(DF(x_0))$  such that

$$(20.5) \quad Q_1 F_h(x_0 + x_1 + x_2) = 0$$

and

$$(20.6) \quad Q_2 F_h(x_0 + x_1 + x_2) = 0.$$

Corresponding to the mapping  $\mathcal{F}$  in (18.5) for the exact problem, we define the mapping  $\mathcal{F}_h : X_1 \times \text{Ker}(DF(x_0)) \rightarrow \text{Range}(DF(x_0))$  by: for  $x_1 \in X_1, x_2 \in \text{Ker}(DF(x_0))$ ,

$$\mathcal{F}_h(x_1, x_2) = Q_2 F_h(x_0 + x_1 + x_2).$$

To study the solutions of (20.6) or of  $\mathcal{F}_h(x_1, x_2) = 0$ , we use Theorem 12.3, Remark 12.1, and Theorem 12.4, whose assumptions are not difficult to check with the hypothesis (20.4), since  $D_{x_1} \mathcal{F}(0, 0) = Q_2 DF(x_0) \in \mathcal{L}(X_1; \text{Range}(DF(x_0)))$  is an isomorphism.

Let  $g : B(0, \epsilon) \subset \text{Ker}(DF(x_0)) \rightarrow X_1$  be the  $C^p$  mapping corresponding to the Lyapunov-Schmidt decomposition for  $F$ , satisfying to (18.6). By Theorem 12.3, there exist positive constants  $h_0 \leq 1$ ,  $\eta$ , and a  $C^p$  mapping  $g_h : B(0, \epsilon) \subset \text{Ker}(DF(x_0)) \rightarrow X_1$  such that for all  $x_2 \in B(0, \epsilon)$ ,  $0 < h \leq h_0$ ,

$$(\mathcal{F}_h(x_1, x_2) = 0 \quad \text{and} \quad x_1 \in B(0, \eta) \subset X_1) \iff x_1 = g_h(x_2).$$

Moreover the following estimates hold: for all  $x_2 \in B(0, \epsilon)$

$$(20.7) \quad \|g(x_2) - g_h(x_2)\|_X \leq C \|F_h(x_0 + g(x_2) + x_2)\|_Z,$$

$$(20.8) \quad \|Dg_h(x_2)\|_{X;X} \leq C.$$

Notice that we also have estimates for higher order derivatives, see (12.14) and (12.15), which are not needed here for our purpose.

To complete our analysis we need to study the equation (20.5) on the solution set of  $\mathcal{F}_h$  in  $(B(0, \epsilon) \times B(0, \eta))$ . The approximate analogue of the mapping  $\mathcal{G}$ , is  $\mathcal{G}_h : B(0, \epsilon) \subset \text{Ker}(DF(x_0)) \rightarrow Z_1$  defined by: for  $x_2 \in B(0, \epsilon)$

$$(20.9) \quad \mathcal{G}_h(x_2) = Q_1 F_h(x_0 + g_h(x_2) + x_2).$$

Then the approximate bifurcation equation reads: find  $x_2 \in B(0, \epsilon)$  such that

$$(20.10) \quad \mathcal{G}_h(x_2) = 0.$$

Before undertaking with more details the study of (20.10), which is the object of Section 21 in the case developed in Section 19, we still need a comparison result between the solutions of the exact problem (20.1) and the ones of the approximate problems (20.3). For that purpose we define the sets

$$\mathcal{S}_g = \{x_2 \in \mathcal{V}; \mathcal{G}(x_2) = 0\},$$

$$\mathcal{S}_{g_h} = \{x_2 \in \mathcal{V}; \mathcal{G}_h(x_2) = 0\},$$

where  $\mathcal{V} \subset B(0, \epsilon) \subset \text{Ker}(DF(x_0))$  is a compact convex neighborhood in  $\text{Ker}(DF(x_0))$  of 0; we also introduce

$$\mathcal{S} = \{x = x_0 + g(x_2) + x_2 \in X; x_2 \in \mathcal{S}_g\},$$

$$\mathcal{S}_h = \{x = x_0 + g_h(x_2) + x_2 \in X; x_2 \in \mathcal{S}_{g_h}\}.$$

Notice that for  $x \in \mathcal{S}$ , we have  $F(x) = 0$  and for  $x_h \in \mathcal{S}_h$ , we have  $F_h(x_h) = 0$ . The distance we will use to compare the sets  $\mathcal{S}_g$ ,  $\mathcal{S}_{g_h}$  and  $\mathcal{S}$ ,  $\mathcal{S}_h$  is the Hausdorff metric.

Let  $Y$  be a vector space endowed with the metric  $d$  and  $\mathcal{Y}$  denote the set of all nonempty bounded closed sets in  $Y$ . For  $(A, B) \in \mathcal{Y} \times \mathcal{Y}$ , we define

$$\rho(A, B) = \sup_{y \in A} d(y, B) \equiv \sup_{y \in A} (\inf_{z \in B} d(y, z)),$$

$$\delta(A, B) = \max(\rho(A, B), \rho(B, A)).$$

$\delta$  is a metric in  $\mathcal{Y}$  called the Hausdorff metric.

**Theorem 20.1.** *In the framework developed above, there exist positive constants  $C$  and  $h_0$  such that for all  $0 < h \leq h_0$*

$$(20.11) \quad \delta(\mathcal{S}, \mathcal{S}_h) \leq C \left[ \delta(\mathcal{S}_{\mathcal{G}}, \mathcal{S}_{\mathcal{G}_h}) + \sup_{x \in \mathcal{S}} \|F_h(x)\|_Z \right].$$

*Proof.* Let  $x_2 \in \mathcal{S}_{\mathcal{G}}$ ,  $x_{2h} \in \mathcal{S}_{\mathcal{G}_h}$ , and

$$x = x_0 + g(x_2) + x_2, \quad x_h = x_0 + g_h(x_{2h}) + x_{2h}.$$

By definition we have  $x \in \mathcal{S}$  and  $x_h \in \mathcal{S}_h$ . Furthermore

$$(20.12) \quad \|x - x_h\|_X \leq \|x_2 - x_{2h}\|_X + \|g(x_2) - g_h(x_2)\|_X + \|g_h(x_2) - g_h(x_{2h})\|_X.$$

We introduce in (20.12) both estimates (20.7) and (20.8) and deduce

$$\|x - x_h\|_X \leq C(\|x_2 - x_{2h}\|_X + \|F_h(x_0 + g(x_2) + x_2)\|_Z)$$

and

$$(20.13) \quad \|x - x_h\|_X \leq C(\|x_2 - x_{2h}\|_X + \sup_{x \in \mathcal{S}} \|F_h(x)\|_Z).$$

The infimum over  $x_h \in \mathcal{S}_h$  in (20.13) leads to

$$\inf_{x_h \in \mathcal{S}_h} \|x - x_h\|_X \leq C(\rho(\mathcal{S}_{\mathcal{G}}, \mathcal{S}_{\mathcal{G}_h}) + \sup_{x \in \mathcal{S}} \|F_h(x)\|_Z)$$

and

$$(20.14) \quad \rho(\mathcal{S}, \mathcal{S}_h) \leq C(\rho(\mathcal{S}_{\mathcal{G}}, \mathcal{S}_{\mathcal{G}_h}) + \sup_{x \in \mathcal{S}} \|F_h(x)\|_Z).$$

In the same way we can prove

$$(20.15) \quad \rho(\mathcal{S}_h, \mathcal{S}) \leq C(\rho(\mathcal{S}_{\mathcal{G}_h}, \mathcal{S}_{\mathcal{G}}) + \sup_{x \in \mathcal{S}} \|F_h(x)\|_Z).$$

The two inequalities (20.14) and (20.15), and the definition of the Hausdorff metric imply the estimate (20.11).  $\square$

To obtain some estimates with respect to  $h$  for  $\delta(\mathcal{S}, \mathcal{S}_h)$ , it remains with (20.11) to bound the quantity  $\delta(\mathcal{S}_{\mathcal{G}}, \mathcal{S}_{\mathcal{G}_h})$  to some power of  $h$ . In the next section we will study this quantity when the functions  $\mathcal{G}$  and  $\mathcal{G}_h$  are defined on a neighborhood of 0 in  $\mathbb{R}^2$  with values in  $\mathbb{R}$ .



### 21. Approximation of simple bifurcation points

To complete our study we still need to present a precise comparison between the solution sets of the bifurcation equation  $\mathcal{G}(x_2) = 0$  and of the approximate bifurcation equation  $\mathcal{G}_h(x_2) = 0$ . Then from the estimate (20.11), we could derive error estimates.

For simplicity we analyze the case of a simple bifurcation point which is developed in Section 19. Let  $F : \mathbb{R} \times \mathcal{X} \rightarrow Z$  be a  $C^p$  mapping,  $p \geq 2$ , and let  $(\lambda_0, \kappa_0) \in \mathbb{R} \times \mathcal{X}$  be a singular solution to (19.1) satisfying to (19.2). Then the bifurcation equation reads

$$(21.1) \quad f(\alpha_1, \alpha_2) = 0,$$

where  $f : B(0, \delta) \subset \mathbb{R}^2 \rightarrow \mathbb{R}$  is of class  $C^p$ ,  $p \geq 2$ , and satisfies

$$f(0, 0) = 0, \quad \frac{\partial f}{\partial \alpha_1}(0, 0) = 0, \quad \text{and} \quad \frac{\partial f}{\partial \alpha_2}(0, 0) = 0.$$

The simple bifurcation case occurs under the assumption

$$(21.2) \quad \det(D^2 f(0, 0)) < 0.$$

In the following we will compare the solutions of approximate bifurcation equations with respect to the solution of the given bifurcation equation (21.1).

We can apply the theory of Section 20 to a family  $\{F_h\}_{0 < h \leq 1}$  of  $C^p$  mappings,  $p \geq 2$ ,

$$F_h : (\lambda, \kappa) \in \mathbb{R} \times \mathcal{X} \rightarrow F_h(\lambda, \kappa) \in Z,$$

satisfying

$$(21.3) \quad \lim_{h \rightarrow 0} \sup_{(\lambda, \kappa) \in B((\lambda_0, \kappa_0), r)} \|D^\ell F(\lambda, \kappa) - D^\ell F_h(\lambda, \kappa)\| = 0,$$

for some  $r > 0$ ,  $\ell = 0, 1$ . Then we get a family of bifurcation equations, say for  $h \in (0, h_0]$

$$(21.4) \quad f_h(\alpha_1, \alpha_2) = 0,$$

where  $f_h : B(0, \delta) \subset \mathbb{R}^2 \rightarrow \mathbb{R}$  is of class  $C^p$ , only to restrict  $\delta$  we keep the same as the one in the definition of  $f$ . With the assumption (21.3) we deduce that

$$(21.5) \quad \lim_{h \rightarrow 0} \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} |f(\alpha_1, \alpha_2) - f_h(\alpha_1, \alpha_2)| = 0.$$

We define the two solution sets

$$\begin{aligned} \mathcal{S}_f &= \{(\alpha_1, \alpha_2) \in B(0, \delta); f(\alpha_1, \alpha_2) = 0\}, \\ \mathcal{S}_{f_h} &= \{(\alpha_1, \alpha_2) \in B(0, \delta); f_h(\alpha_1, \alpha_2) = 0\}. \end{aligned}$$

We compare now the behavior of  $\mathcal{S}_{f_h}$  and of  $\mathcal{S}_f$ . Our goal is to get error estimates for the Hausdorff distance  $\delta(\mathcal{S}_f, \mathcal{S}_{f_h})$  and to give qualitative results for  $\mathcal{S}_{f_h}$ .

**Theorem 21.1.** *Under the assumptions (21.2) and (21.5), there exists a constant  $C > 0$  independent of  $h$  such that*

$$(21.6) \quad \delta(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq C \left( \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} |f(\alpha_1, \alpha_2) - f_h(\alpha_1, \alpha_2)| \right)^{1/2}.$$

*Proof.* For notation simplicity we set

$$\epsilon_h = \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} |f(\alpha_1, \alpha_2) - f_h(\alpha_1, \alpha_2)|.$$

Applying the Morse lemma, we deduce that there exists a  $C^{p-1}$  diffeomorphism  $\Psi$  from  $B(0, \delta)$ , only to restrict  $\delta$ , into  $\mathbb{R}^2$  such that the bifurcation equation reads

$$f(\alpha_1, \alpha_2) \equiv \Psi_1(\alpha_1, \alpha_2)\Psi_2(\alpha_1, \alpha_2) = 0,$$

see (19.12). Consequently without loss of generality, we can assume that the function  $f$  is, for  $(\alpha_1, \alpha_2) \in B(0, \delta)$

$$f(\alpha_1, \alpha_2) = \alpha_1\alpha_2.$$

Then  $\rho(\mathcal{S}_{f_h}, \mathcal{S}_f)$  given by

$$\rho(\mathcal{S}_{f_h}, \mathcal{S}_f) = \sup_{(\alpha_1, \alpha_2) \in \mathcal{S}_{f_h}} \left( \inf_{(\beta_1, \beta_2) \in \mathcal{S}_f} \|(\alpha_1, \alpha_2) - (\beta_1, \beta_2)\| \right),$$

where  $\|\cdot\|$  denotes the Euclidian norm in  $\mathbb{R}^2$ , can be bounded by

$$(21.7) \quad \rho(\mathcal{S}_{f_h}, \mathcal{S}_f) \leq \epsilon_h.$$

Let us estimate now  $\rho(\mathcal{S}_f, \mathcal{S}_{f_h})$ . We define the sets

$$\begin{aligned} C_h^+ &= \{(\alpha_1, \alpha_2) \in B(0, \delta); f(\alpha_1, \alpha_2) \geq \epsilon_h\}, \\ C_h^- &= \{(\alpha_1, \alpha_2) \in B(0, \delta); f(\alpha_1, \alpha_2) \leq -\epsilon_h\}. \end{aligned}$$

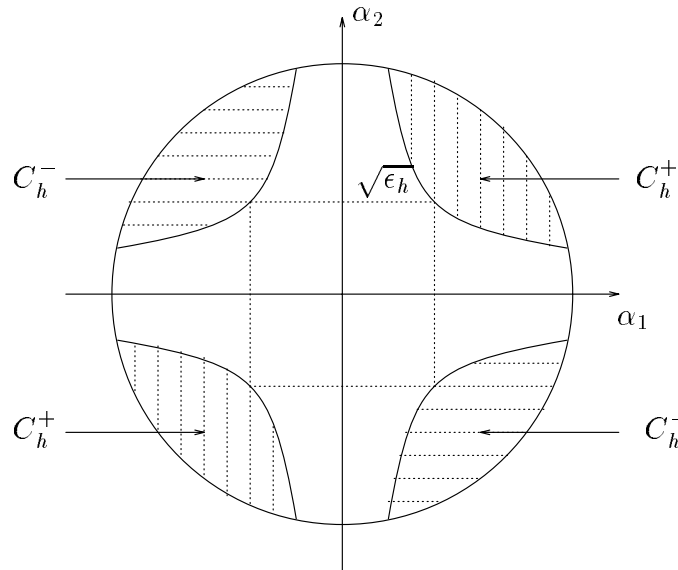


FIGURE 21.1: SETS  $C_h^+$  AND  $C_h^-$ .

From the assumption (21.5) we know that  $\epsilon_h$  tends to 0 when  $h$  tends to 0, so for  $h$  small enough both sets  $C_h^+$  and  $C_h^-$  are not empty. Let  $(\alpha_1, \alpha_2)$  be in  $\mathcal{S}_f$ . There exist  $(\alpha_{1h}^+, \alpha_{2h}^+)$  in  $C_h^+$  and  $(\alpha_{1h}^-, \alpha_{2h}^-)$  in  $C_h^-$  satisfying

$$\|(\alpha_1, \alpha_2) - (\alpha_{1h}^+, \alpha_{2h}^+)\| \leq \sqrt{2\epsilon_h}, \quad \|(\alpha_1, \alpha_2) - (\alpha_{1h}^-, \alpha_{2h}^-)\| \leq \sqrt{2\epsilon_h}.$$

By the definition of  $\epsilon_h$  we have

$$\begin{aligned} f_h(\alpha_{1h}^+, \alpha_{2h}^+) &\geq f(\alpha_{1h}^+, \alpha_{2h}^+) - \epsilon_h \geq 0, \\ f_h(\alpha_{1h}^-, \alpha_{2h}^-) &\leq f(\alpha_{1h}^-, \alpha_{2h}^-) + \epsilon_h \leq 0. \end{aligned}$$

So there exists  $\theta \in [0, 1]$  such that

$$f_h(\theta(\alpha_{1h}^+, \alpha_{2h}^+) + (1 - \theta)(\alpha_{1h}^-, \alpha_{2h}^-)) = 0$$

and if we denote by  $(\alpha_{1h}^0, \alpha_{2h}^0)$  that solution of  $f_h$ , we have

$$\|(\alpha_1, \alpha_2) - (\alpha_{1h}^0, \alpha_{2h}^0)\| \leq \sqrt{2\epsilon_h}.$$

Then we deduce

$$(21.8) \quad \rho(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq \sqrt{2\epsilon_h}.$$

Finally from both estimates (21.7) and (21.8) we can conclude.  $\square$

*Remark 21.1.* If we do not add any assumptions to (21.2) and (21.5), then the estimate (21.6) is optimal. Consider  $f(\alpha_1, \alpha_2) = \alpha_1\alpha_2$  and  $f_h(\alpha_1, \alpha_2) = \alpha_1\alpha_2 + \epsilon_h$ ,  $\epsilon_h > 0$  for  $h \in (0, 1]$ , then we have precisely

$$\delta(\mathcal{S}_f, \mathcal{S}_{f_h}) = \sqrt{2\epsilon_h}. \quad \square$$

*Remark 21.2.* In the proof of Theorem 21.1, we have only used the continuity of  $f_h$  and both relations (21.2) and (21.5). In fact with the assumption (21.3), we also have

$$(21.9) \quad \lim_{h \rightarrow 0} \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\| = 0.$$

Then we will derive in the following the estimate

$$(21.10) \quad \delta(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq C(\sqrt{f_h(0, 0)} + \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\|).$$

It is sometimes possible to improve the error analysis by using the estimate (21.10) rather than the (21.6) one.

To check the estimate (21.10) we proceed like in the proof of Theorem 21.1. Notice first that with the Taylor expansion we get for  $(\alpha_1, \alpha_2) \in B(0, \delta)$

$$\begin{aligned} f(\alpha_1, \alpha_2) - f_h(\alpha_1, \alpha_2) &= f(0, 0) - f_h(0, 0) \\ &\quad + \int_0^1 (Df(t(\alpha_1, \alpha_2)) - Df_h(t(\alpha_1, \alpha_2))) (\alpha_1, \alpha_2) dt. \end{aligned}$$

Since  $f(0, 0) = 0$ , we deduce that

$$|f(\alpha_1, \alpha_2) - f_h(\alpha_1, \alpha_2)| \leq |f_h(0, 0)| + \epsilon_h^1 \|(\alpha_1, \alpha_2)\|,$$

where  $\epsilon_h^1$  stands for

$$\epsilon_h^1 = \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\|.$$

As in the proof of Theorem 21.1, we can assume that  $f$  is given by  $f(\alpha_1, \alpha_2) = \alpha_1 \alpha_2$  and satisfies

$$|f(\alpha_1, \alpha_2) - f_h(\alpha_1, \alpha_2)| \leq |f_h(0, 0)| + C\epsilon_h^1(|\alpha_1| + |\alpha_2|).$$

Then we have

$$(21.11) \quad \rho(\mathcal{S}_{f_h}, \mathcal{S}_f) \leq C(|f_h(0, 0)| + \epsilon_h^1).$$

To get an estimate for the distance  $\rho(\mathcal{S}_f, \mathcal{S}_{f_h})$ , we define the two sets

$$\begin{aligned} C_h^+ &= \{(\alpha_1, \alpha_2) \in B(0, \delta); f(\alpha_1, \alpha_2) \geq |f_h(0, 0)| + C\epsilon_h^1(|\alpha_1| + |\alpha_2|)\}, \\ C_h^- &= \{(\alpha_1, \alpha_2) \in B(0, \delta); f(\alpha_1, \alpha_2) \leq -[|f_h(0, 0)| + C\epsilon_h^1(|\alpha_1| + |\alpha_2|)]\}. \end{aligned}$$

From both assumptions (21.5) and (21.9), we know that the sets  $C_h^+$  and  $C_h^-$  are not empty. Using the same method as in the proof of Theorem 21.1, we get

$$(21.12) \quad \rho(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq C(\sqrt{|f_h(0, 0)|} + \epsilon_h^1).$$

The two estimates (21.11) and (21.12) imply the (21.10) one.  $\square$

*Remark 21.3.* We still can modify the estimate (21.10) in the following way. We assume (21.2), (21.5), (21.9), and

$$(21.13) \quad \lim_{h \rightarrow 0} \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} \|D^2 f(\alpha_1, \alpha_2) - D^2 f_h(\alpha_1, \alpha_2)\| = 0.$$

Notice that the assumption (21.13) is true if we assume (21.3) with  $\ell = 2$ . Then the following estimate holds

$$(21.14) \quad \delta(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq C(\sqrt{|f_h(0, 0)|} + \sup_{(\alpha_1, \alpha_2) \in \mathcal{S}_f} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\|).$$

Indeed let  $(\alpha_1, \alpha_2)$  be in  $B(0, \delta)$  and  $(\alpha_1^0, \alpha_2^0)$  be in  $\mathcal{S}_f$ . From the Taylor expansion we deduce that

$$\begin{aligned} f_h(\alpha_1, \alpha_2) - f(\alpha_1, \alpha_2) &= \\ f_h(\alpha_1^0, \alpha_2^0) - f(\alpha_1^0, \alpha_2^0) &+ (Df_h(\alpha_1^0, \alpha_2^0) - Df(\alpha_1^0, \alpha_2^0))(\alpha_1 - \alpha_1^0, \alpha_2 - \alpha_2^0) \\ &+ \frac{1}{2} (D^2 f_h(\tilde{\alpha}_1, \tilde{\alpha}_2) - D^2 f(\tilde{\alpha}_1, \tilde{\alpha}_2))((\alpha_1 - \alpha_1^0, \alpha_2 - \alpha_2^0), (\alpha_1 - \alpha_1^0, \alpha_2 - \alpha_2^0)) \end{aligned}$$

where  $(\tilde{\alpha}_1, \tilde{\alpha}_2)$  belongs to the segment passing by  $(\alpha_1, \alpha_2)$  and  $(\alpha_1^0, \alpha_2^0)$ . So

$$\begin{aligned} |f_h(\alpha_1, \alpha_2) - f(\alpha_1, \alpha_2)| &\leq |f_h(\alpha_1^0, \alpha_2^0) - f(\alpha_1^0, \alpha_2^0)| \\ &+ \epsilon_h^1(\mathcal{S}_f) \|(\alpha_1 - \alpha_1^0, \alpha_2 - \alpha_2^0)\| + \frac{1}{2} \epsilon_h^2 \|(\alpha_1 - \alpha_1^0, \alpha_2 - \alpha_2^0)\|^2, \end{aligned}$$

where

$$\begin{aligned} \epsilon_h^1(\mathcal{S}_f) &= \sup_{(\alpha_1, \alpha_2) \in \mathcal{S}_f} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\| \\ \epsilon_h^2 &= \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} \|D^2 f(\alpha_1, \alpha_2) - D^2 f_h(\alpha_1, \alpha_2)\|. \end{aligned}$$

Noticing that

$$\begin{aligned} |f_h(\alpha_1^0, \alpha_2^0) - f(\alpha_1^0, \alpha_2^0)| &\leq |f_h(0, 0) - f(0, 0)| + C \epsilon_h^1(\mathcal{S}_f) \|(\alpha_1^0, \alpha_2^0)\| \\ &\leq |f_h(0, 0)| + C \epsilon_h^1(\mathcal{S}_f) \|(\alpha_1^0, \alpha_2^0)\|, \end{aligned}$$

we get

$$\begin{aligned} |f_h(\alpha_1, \alpha_2) - f(\alpha_1, \alpha_2)| &\leq |f_h(0, 0)| \\ &+ C \epsilon_h^1(\mathcal{S}_f) (\|(\alpha_1^0, \alpha_2^0)\| + \|(\alpha_1 - \alpha_1^0, \alpha_2 - \alpha_2^0)\|) \\ &+ \frac{1}{2} \epsilon_h^2 \|(\alpha_1 - \alpha_1^0, \alpha_2 - \alpha_2^0)\|^2 \end{aligned}$$

We follow then the arguments in Remark 21.2 by reducing  $f(\alpha_1, \alpha_2)$  to  $\alpha_1 \alpha_2$ , to check (21.14).  $\square$

The next result describes precisely the nature of the solution set  $\mathcal{S}_{f_h}$ .

**Theorem 21.2.** *Let us assume that (21.2) holds and in addition that*

$$(21.15) \quad \lim_{h \rightarrow 0} Df_h(0, 0) = 0,$$

$$(21.16) \quad \lim_{h \rightarrow 0} D^2 f_h(0, 0) = D^2 f(0, 0),$$

and there exists a constant  $L$  such that for all  $h \in (0, h_0]$  and  $(\alpha_1, \alpha_2) \in B(0, \delta)$

$$(21.17) \quad \|D^2 f_h(\alpha_1, \alpha_2) - D^2 f_h(0, 0)\| \leq L\|(\alpha_1, \alpha_2)\|.$$

Then there exist two positive constants  $C$  and  $h_0$ , a neighborhood  $\mathcal{V}$  of  $(0, 0)$ , and for  $h \leq h_0$  a point  $(\alpha_{1h}, \alpha_{2h}) \in \mathbb{R}^2$  and a  $C^1$  diffeomorphism  $\Phi_h$  from  $\mathcal{V}$  into  $\mathbb{R}^2$  such that

$$(21.18) \quad \|(\alpha_{1h}, \alpha_{2h})\| \leq C\|Df_h(0, 0)\|,$$

$$(21.19) \quad Df_h(\alpha_{1h}, \alpha_{2h}) = 0,$$

$$\Phi_h(\alpha_{1h}, \alpha_{2h}) = (\alpha_{1h}, \alpha_{2h}), \quad D\Phi_h(\alpha_{1h}, \alpha_{2h}) = I,$$

$\Phi_h(\mathcal{S}_{f_h} \cap \mathcal{V})$  is a part of a hyperbola.

*Proof.* Under the assumptions (21.15), (21.16), and (21.17), we can apply Theorem 6.1 with  $F_h = Df_h$  and  $L_h = L$  to prove the existence of  $(\alpha_{1h}, \alpha_{2h}) \in B(0, \delta)$ , for  $h \in (0, h_0]$ , only to restrict  $h_0$  if necessary, such that both relations (21.18) and (21.19) hold. Moreover for  $h \leq h_0$  we have

$$(21.20) \quad \det(D^2 f_h(\alpha_{1h}, \alpha_{2h})) < 0.$$

We are in a position to apply Theorem 19.1 to the mapping  $f_h$  at the point  $(\alpha_{1h}, \alpha_{2h})$ . There exists a  $C^{p-1}$  diffeomorphism  $(\alpha_1, \alpha_2) \in \mathcal{V} \rightarrow \Phi_h(\alpha_1, \alpha_2) \in \mathbb{R}^2$  such that

$$(21.21) \quad f_h(\alpha_1, \alpha_2) = f_h(\alpha_{1h}, \alpha_{2h}) + \frac{1}{2} \left( D^2 f_h(\alpha_{1h}, \alpha_{2h})(\Phi_h(\alpha_1, \alpha_2) - (\alpha_{1h}, \alpha_{2h}), \Phi_h(\alpha_1, \alpha_2) - (\alpha_{1h}, \alpha_{2h})) \right)$$

where  $\mathcal{V}$  is a neighborhood of  $(\alpha_{1h}, \alpha_{2h})$ , which a priori is depending on  $h$ . Moreover  $\Phi_h(\alpha_{1h}, \alpha_{2h}) = (\alpha_{1h}, \alpha_{2h})$  and  $D\Phi_h(\alpha_{1h}, \alpha_{2h}) = I$ . Consequently the set  $\mathcal{S}_{f_h}$  is diffeomorphic to a part of the hyperbola given in the  $\Phi$ -plane by

$$(21.22) \quad f_h(\alpha_{1h}, \alpha_{2h}) + \frac{1}{2} \left( D^2 f_h(\alpha_{1h}, \alpha_{2h})(\Phi - (\alpha_{1h}, \alpha_{2h}), \Phi - (\alpha_{1h}, \alpha_{2h})) \right) = 0.$$

Going back to the proof of Morse Lemma, see the appendix in CROUZEIX and RAPPAZ [1989], it is not difficult to check that under the assumptions (21.15), (21.16), and (21.17) and with  $p = 2$ , we can choose  $\mathcal{V}$  independently of  $h$  and containing the point  $(0, 0)$ , only to restrict  $h_0$  if necessary.  $\square$

*Remark 21.4.* With the relation (21.21), we deduce that solving the bifurcation equation (21.4) is equivalent to solving the equation (21.22). In the  $\Phi$ -plane, the equation

$$f_h(\alpha_{1h}, \alpha_{2h}) + \frac{1}{2} \left( D^2 f_h(\alpha_{1h}, \alpha_{2h})(\Phi - (\alpha_{1h}, \alpha_{2h}), \Phi - (\alpha_{1h}, \alpha_{2h})) \right) = 0$$

is the equation of a hyperbola centered at  $(\alpha_{1h}, \alpha_{2h})$ , since  $\det(D^2 f_h(\alpha_{1h}, \alpha_{2h})) < 0$ . The hyperbola represents the solution set to (21.4) since  $(\alpha_{1h}, \alpha_{2h})$  is invariant through

$\Phi_h$  and  $\bar{\Phi}_h$  is tangent to the identity at  $(\alpha_{1h}, \alpha_{2h})$ . The point  $(\alpha_{1h}, \alpha_{2h})$  is the center of the hyperbola and with both relations (21.18) and (21.15) we know that  $(\alpha_{1h}, \alpha_{2h})$  tends to 0.  $\square$

The solution set  $\mathcal{S}_{f_h}$  is described by the equation (21.22) and we can distinguish two cases.

If the value of  $f_h(\alpha_{1h}, \alpha_{2h})$  is not zero which is the common case,

$$(21.23) \quad f_h(\alpha_{1h}, \alpha_{2h}) \neq 0,$$

then  $\mathcal{S}_{f_h}$  is diffeomorphic to a part of a nondegenerate hyperbola.

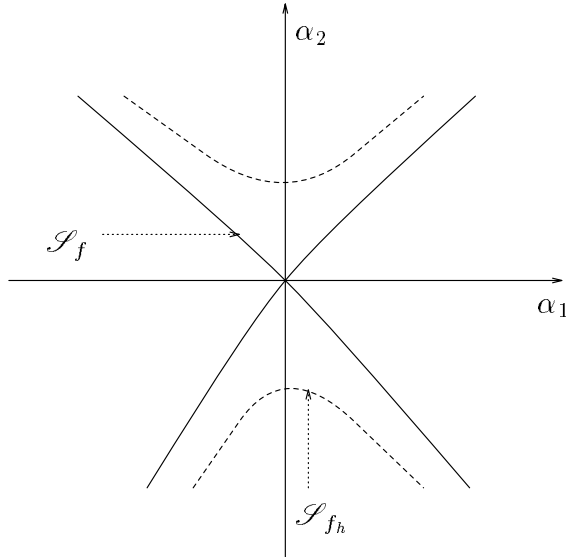


FIGURE 21.2: IMPERFECT NUMERICAL BIFURCATION.

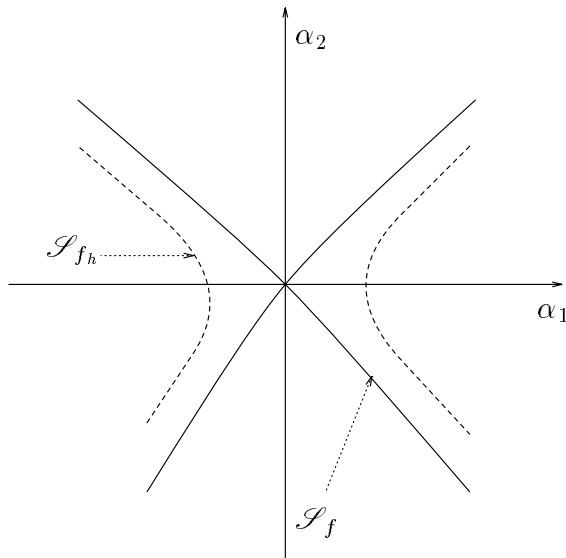


FIGURE 21.3: IMPERFECT NUMERICAL BIFURCATION.

That situation is represented on the figures 21.2 and 21.3, depending on the sign of  $f_h(\alpha_{1h}, \alpha_{2h})$ . We say that we have an imperfect numerical bifurcation.

The second case which can occur is

$$(21.24) \quad f_h(\alpha_{1h}, \alpha_{2h}) = 0.$$

Then  $\mathcal{S}_{f_h}$  is diffeomorphic to a part of a degenerate hyperbola.

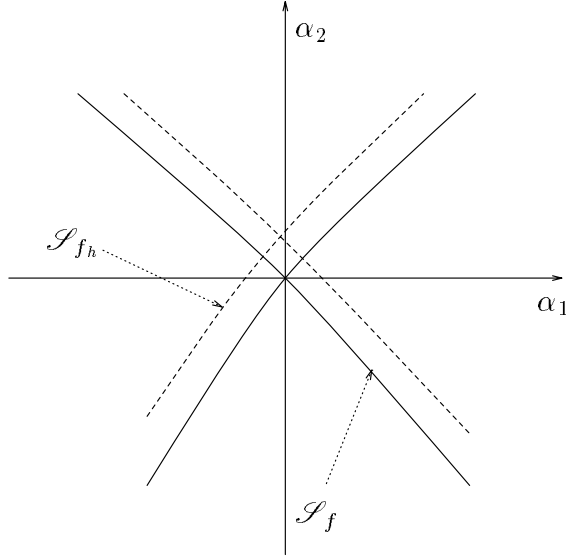


FIGURE 21.4: NUMERICAL BIFURCATION.

That situation is represented on the figure 21.4. We say that we have a numerical bifurcation. The condition (21.24) is satisfied for instance in the case of bifurcation from the trivial branch, when  $F(\lambda, 0) = 0$  for all  $\lambda \in \mathbb{R}$ , see the example developed in the next section. It can be also satisfied when we have symmetries, see SATTINGER [1979].

*Remark 21.5.* In the framework of Theorem 21.2, we assume furthermore that (21.5), (21.9) hold, and that  $f_h(\alpha_{1h}, \alpha_{2h}) = 0$ . Here  $(\alpha_{1h}, \alpha_{2h})$  is the center of the hyperbola given in Theorem 21.2. Then from the estimate (21.10) we get

$$(21.25) \quad \delta(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq C \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\|.$$

Indeed we notice that

$$|f_h(0, 0)| \leq C \|(\alpha_{1h}, \alpha_{2h})\|^2$$

since  $f_h(\alpha_{1h}, \alpha_{2h}) = 0$  and  $Df_h(\alpha_{1h}, \alpha_{2h}) = 0$ , see (21.19). From the estimate (21.18) we show

$$\begin{aligned} |f_h(0, 0)| &\leq C \|Df(0, 0) - Df_h(0, 0)\|^2 \\ &\leq C \left( \sup_{(\alpha_1, \alpha_2) \in B(0, \delta)} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\| \right)^2. \end{aligned}$$



The estimate (21.25) is then a direct consequence of (21.10).  $\square$

*Remark 21.6.* In the framework of Theorem 21.2, we assume furthermore that (21.5), (21.9), (21.13) hold, and that  $f_h(\alpha_{1h}, \alpha_{2h}) = 0$ . Like in Remark 21.5, we can derive from the estimate (21.14) the following one

$$(21.26) \quad \delta(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq C \sup_{(\alpha_1, \alpha_2) \in \mathcal{S}_f} \|Df(\alpha_1, \alpha_2) - Df_h(\alpha_1, \alpha_2)\|. \quad \square$$

*Remark 21.7.* It is possible to generalize the results of Section 21 to the case where the mappings  $f, f_h$  are defined in  $\mathbb{R}^3$ , that is

$$f : B(0, \delta) \subset \mathbb{R}^3 \rightarrow \mathbb{R}, \quad f_h : B(0, \delta) \subset \mathbb{R}^3 \rightarrow \mathbb{R}.$$

Then under analogous assumptions, we can prove that  $f$  is  $C^{p-1}$  diffeomorphic to a part of a cone and that  $f_h$  is  $C^{p-1}$  diffeomorphic to a part of a hyperboloid with one or two sheets.  $\square$

## 22. A model example

To illustrate the general results in Sections 18 through 21, we will develop a simple model example putting emphasis on the methodology.

Let  $\Omega \subset \mathbb{R}^2$  be a convex polygonal domain with the boundary  $\partial\Omega$ . We shall study the problem to find  $\lambda \in \mathbb{R}$  and  $u \in H_0^1(\Omega)$  satisfying

$$(22.1) \quad -\Delta u - \lambda u + u^2 = 0 \quad \text{in } \Omega$$

or in variational form, for all  $v \in H_0^1(\Omega)$

$$(22.2) \quad \int_{\Omega} \mathbf{grad} u \mathbf{grad} v \, dx - \lambda \int_{\Omega} uv \, dx + \int_{\Omega} u^2 v \, dx = 0.$$

If the mapping  $T \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$  denotes the inverse of the  $(-\Delta)$  operator with homogeneous boundary conditions, then problem (22.2) reads: find  $\lambda \in \mathbb{R}$  and  $u \in H_0^1(\Omega)$  such that

$$(22.3) \quad u + T(u^2 - \lambda u) = 0,$$

see (3.6). Clearly the trivial branch  $\{(\lambda, 0); \lambda \in \mathbb{R}\}$  belongs to the solution set of (22.3). Subsequently we look for bifurcation from the trivial branch and study a finite element approximation of the problem.

We are in the framework of Section 19 with the following notations:  $\mathcal{X} = H_0^1(\Omega)$ ,  $X = \mathbb{R} \times H_0^1(\Omega)$ ,  $Z = H_0^1(\Omega)$ , and the mapping  $F : \mathbb{R} \times H_0^1(\Omega) \rightarrow H_0^1(\Omega)$  given by, for  $(\lambda, v) \in \mathbb{R} \times H_0^1(\Omega)$

$$(22.4) \quad F(\lambda, v) = v + T(v^2 - \lambda v).$$

Clearly the mapping  $F$  is of class  $C^\infty$ . For  $\lambda \in \mathbb{R}$ , the derivative  $DF(\lambda, 0)$  is, for  $(\delta, w) \in \mathbb{R} \times H_0^1(\Omega)$

$$(22.5) \quad DF(\lambda, 0)(\delta, w) = w - \lambda T w.$$

We immediately deduce that

$$(22.6) \quad \text{Ker}(DF(\lambda, 0)) = \mathbb{R} \times \text{Ker}(I - \lambda T)$$

with  $I$  the identity operator in  $H_0^1(\Omega)$  and  $T$  the restriction to  $H_0^1(\Omega)$  of the above operator  $T \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$ . Let now  $\lambda_0$  be the first eigenvalue of the eigenproblem: find  $\lambda \in \mathbb{R}$ ,  $\xi \in H_0^1(\Omega)$ ,  $\xi \neq 0$  such that

$$-\Delta \xi = \lambda \xi \quad \text{in } \Omega.$$

It is well known that  $\lambda_0$  is a simple positive eigenvalue and that a corresponding eigenvector  $\xi$  can be chosen positive and with a norm equal to 1, that is

$$\xi > 0 \quad \text{in } \Omega \quad \text{and} \quad \int_{\Omega} |\mathbf{grad} \xi|^2 dx = 1.$$

Then

$$(22.7) \quad \text{Ker}(DF(\lambda_0, 0)) = \mathbb{R} \times \text{span}\{\xi\},$$

which is of dimension 2. With the expression (22.5), we get

$$(22.8) \quad \text{Range}(DF(\lambda_0, 0)) = \text{Range}(I - \lambda_0 T).$$

Since the operator  $T \in \mathcal{L}(H_0^1(\Omega); H_0^1(\Omega))$  is self adjoint when  $H_0^1(\Omega)$  is endowed with the scalar product, for  $v_1, v_2 \in H_0^1(\Omega)$

$$(v_1, v_2)_{1, \Omega} = \int_{\Omega} \mathbf{grad} v_1 \mathbf{grad} v_2 dx,$$

the range of  $DF(\lambda_0, 0)$  is the orthogonal complement of  $\xi$ , that is

$$(22.9) \quad \text{Range}(DF(\lambda_0, 0)) = \{v \in H_0^1(\Omega); (v, \xi)_{1, \Omega} = 0\},$$

which is of codimension 1.

The solution  $(\lambda_0, 0)$  satisfies the following hypotheses

$$(22.10) \quad \begin{cases} \text{(i)} & \text{Ker}(D_v F(\lambda_0, 0)) = \text{span}\{\xi\}, \\ \text{(ii)} & \text{the range of } D_v F(\lambda_0, 0) \text{ is closed, of codimension 1,} \\ \text{(iii)} & D_\lambda F(\lambda_0, 0) = 0. \end{cases}$$

Thus  $(\lambda_0, 0)$  is a singular solution of (22.3) and the assumptions (19.2) are verified. In our case a decomposition of both spaces  $\mathbb{R} \times H_0^1(\Omega)$  and  $H_0^1(\Omega)$  can be written explicitly. In (19.4) and (19.5), we choose

$$\begin{aligned} X_1 &= \{0\} \times \mathcal{X}_1 \equiv \{0\} \times \text{Range}(I - \lambda_0 T), \\ Z_1 &= \text{Ker}(I - \lambda_0 T) = \text{span}\{\xi\}. \end{aligned}$$

Then the two projectors  $Q_1 \in \mathcal{L}(H_0^1(\Omega); Z_1)$  and  $Q_2 \in \mathcal{L}(H_0^1(\Omega); \text{Range}(I - \lambda_0 T))$  are the orthogonal projectors with respect to the scalar product  $(\cdot, \cdot)_{1, \Omega}$  in  $H_0^1(\Omega)$ . The two projectors  $P_1 \in \mathcal{L}(\mathbb{R} \times H_0^1(\Omega); X_1)$  and  $P_2 \in \mathcal{L}(\mathbb{R} \times H_0^1(\Omega); \text{Ker}(DF(\lambda_0, 0)))$  can be written in the following way

$$P_1(\delta, v) = (0, Q_2 v), \quad P_2(\delta, v) = (\delta, Q_1 v).$$

There exist positive constants  $\gamma, \eta$ , and a mapping  $g : B(0, \gamma) \subset \mathbb{R}^2 \rightarrow \mathcal{X}_1$  of class  $C^\infty$  such that for all  $(\alpha_1, \alpha_2) \in B(0, \gamma)$

$$(Q_2 F(\lambda_0 + \alpha_1, \kappa_1 + \alpha_2 \xi) = 0 \text{ and } \kappa_1 \in B(0, \eta) \subset \mathcal{X}_1) \iff \kappa_1 = g(\alpha_1, \alpha_2).$$

In particular for  $(\alpha_1, \alpha_2) \in B(0, \gamma)$

$$(22.11) \quad Q_2 F(\lambda_0 + \alpha_1, g(\alpha_1, \alpha_2) + \alpha_2 \xi) = 0,$$

and  $g(\alpha_1, 0) = 0$  for  $(\alpha_1, 0) \in B(0, \gamma)$ .

If the functional  $B \in H^{-1}(\Omega)$  satisfies

$$B(\xi) = 1 \quad \text{and} \quad B(z) = 0 \quad \text{for all } z \in \text{Range}(I - \lambda_0 T),$$

then clearly for all  $z \in H_0^1(\Omega)$

$$(22.12) \quad B(z) = (z, \xi)_{1, \Omega}.$$

The bifurcation equation reads, for  $(\alpha_1, \alpha_2) \in B(0, \gamma)$

$$(22.13) \quad f(\alpha_1, \alpha_2) \equiv (F(\lambda_0 + \alpha_1, g(\alpha_1, \alpha_2) + \alpha_2 \xi), \xi)_{1, \Omega} = 0.$$

The function  $f$  can be written in a simpler form. With the definition (22.4) of  $F$  and the relations

$$(g(\alpha_1, \alpha_2), \xi)_{1, \Omega} = 0 = (Tg(\alpha_1, \alpha_2), \xi)_{1, \Omega}$$

we get

$$\begin{aligned} f(\alpha_1, \alpha_2) &= \left( \xi, \alpha_2 \xi + T((\alpha_2 \xi + g(\alpha_1, \alpha_2))^2 - (\lambda_0 + \alpha_1) \alpha_2 \xi) \right)_{1, \Omega} \\ (22.14) \quad &= \left( \xi, T((\alpha_2 \xi + g(\alpha_1, \alpha_2))^2 - \alpha_1 \alpha_2 \xi) \right)_{1, \Omega} \\ &= \int_{\Omega} \xi (\alpha_2 \xi + g(\alpha_1, \alpha_2))^2 dx - \frac{\alpha_1 \alpha_2}{\lambda_0}. \end{aligned}$$

From the relations (22.11) and (22.14), we check with some simple computations

$$(22.15) \quad g(0, 0) = 0 \quad \text{and} \quad f(0, 0) = 0,$$

$$(22.16) \quad \frac{\partial g}{\partial \alpha_1}(0, 0) = 0, \quad \frac{\partial g}{\partial \alpha_2}(0, 0) = 0 \quad \text{and} \quad \frac{\partial f}{\partial \alpha_1}(0, 0) = 0, \quad \frac{\partial f}{\partial \alpha_2}(0, 0) = 0,$$

$$(22.17) \quad \frac{\partial^2 g}{\partial \alpha_1^2}(0, 0) = 0, \quad \frac{\partial^2 g}{\partial \alpha_1 \partial \alpha_2}(0, 0) = 0, \quad \frac{\partial^2 g}{\partial \alpha_2^2}(0, 0) = -2(I - \lambda_0 T)^{-1} Q_2 T \xi^2,$$

$$(22.18) \quad \frac{\partial^2 f}{\partial \alpha_1^2}(0, 0) = 0, \quad \frac{\partial^2 f}{\partial \alpha_1 \partial \alpha_2}(0, 0) = -1/\lambda_0, \quad \frac{\partial^2 f}{\partial \alpha_2^2}(0, 0) = 2 \int_{\Omega} \xi^3 dx;$$

notice that here the operator  $I - \lambda_0 T$  is an isomorphism onto  $\mathcal{X}_1$ . Since

$$\det(D^2 f(0, 0)) = -\left(\frac{1}{\lambda_0}\right)^2,$$

we are precisely in the case where  $(\lambda_0, 0)$  is a simple bifurcation point. Moreover we notice that for all  $(\alpha_1, 0) \in B(0, \gamma)$

$$(22.19) \quad f(\alpha_1, 0) = 0,$$

since  $g(\alpha_1, 0) = 0$ .

Once studied the solution set of (22.3) in a neighborhood of the singular solution  $(\lambda_0, 0)$ , we turn our attention to finite element approximations of (22.3). For ease of exposition we have assumed that  $\Omega$  is a convex polygonal domain in  $\mathbb{R}^2$ , so we can consider a regular family  $\{\mathcal{T}_h\}_{0 < h \leq 1}$  of triangulations of  $\bar{\Omega}$ . To each triangulation we associate the finite element subspace

$$(22.20) \quad V_h = \{v \in C^0(\bar{\Omega}); v|_T \in \mathcal{P}_1(T), \text{ for all } T \in \mathcal{T}_h\} \cap H_0^1(\Omega),$$

see (3.12). A finite element approximation of (22.2) consists in finding  $\lambda \in \mathbb{R}$  and  $u_h \in V_h$  such that for all  $v_h \in V_h$

$$(22.21) \quad \int_{\Omega} \mathbf{grad} u_h \mathbf{grad} v_h dx - \lambda \int_{\Omega} u_h v_h dx + \int_{\Omega} u_h^2 v_h dx = 0.$$

We define the mapping  $T_h \in \mathcal{L}(H^{-1}(\Omega); V)$ , the discrete analogue of  $T$ ; given  $f \in H^{-1}(\Omega)$ ,  $T_h f \equiv w_h \in V_h$  is the unique solution of

$$(22.22) \quad \text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad} w_h \mathbf{grad} v_h dx = \langle f, v_h \rangle_{H^{-1}(\Omega) H_0^1(\Omega)}.$$

Then problem (22.21) can be written in the equivalent way: find  $\lambda \in \mathbb{R}$  and  $u_h \in H_0^1(\Omega)$  such that

$$(22.23) \quad u_h + T_h(u_h^2 - \lambda u_h) = 0.$$

We notice that if  $u_h \in H_0^1(\Omega)$  satisfies (22.23), then  $u_h$  belongs to the range of  $T_h$ , which is precisely  $V_h$ , so  $u_h$  is a solution to (22.21). The discrete analogue of the mapping  $F$  in (22.4) is  $F_h : \mathbb{R} \times H_0^1(\Omega) \rightarrow H_0^1(\Omega)$  defined by, for  $(\lambda, v) \in \mathbb{R} \times H_0^1(\Omega)$

$$(22.24) \quad F_h(\lambda, v) = v + T_h(v^2 - \lambda v).$$

To study the approximate problem (22.23), we shall apply the results of Sections 20 and 21, with  $F$  defined in (22.4) and  $F_h$  in (22.24). We check first the assumption (20.4). For  $(\lambda, v) \in \mathbb{R} \times H_0^1(\Omega)$

$$(22.25) \quad F_h(\lambda, v) - F(\lambda, v) = (T_h - T)(v^2 - \lambda v),$$

for  $(\delta, w) \in \mathbb{R} \times H_0^1(\Omega)$

$$(22.26) \quad (DF_h(\lambda, v) - DF(\lambda, v))(\delta, w) = -\delta(T_h - T)v + (T_h - T)(2vw - \lambda w),$$

and for  $(\delta_1, w_1), (\delta_2, w_2) \in \mathbb{R} \times H_0^1(\Omega)$

$$(22.27) \quad \begin{aligned} & (D^2 F_h(\lambda, v) - D^2 F(\lambda, v))((\delta_1, w_1), (\delta_2, w_2)) \\ & = -(T_h - T)(\delta_1 w_2 + \delta_2 w_1) + 2(T_h - T)(w_1 w_2). \end{aligned}$$

Moreover for  $\ell > 2$

$$(22.28) \quad D^\ell F(\lambda, v) = D^\ell F_h(\lambda, v) = 0.$$

From the classical finite element estimates, for all  $f \in L^2(\Omega)$  we have

$$(22.29) \quad |(T - T_h)f|_{1,\Omega} \leq Ch \|f\|_{0,\Omega},$$

where  $T$  and  $T_h$  are given in (3.4) and (22.22) and so

$$(22.30) \quad \lim_{h \rightarrow 0} \|T - T_h\|_{L^2(\Omega); H_0^1(\Omega)} = 0.$$

Then from the formulae (22.25)-(22.28) and the continuity of the embedding  $H_0^1(\Omega) \subset L^p(\Omega)$ , we get

$$(22.31) \quad \lim_{h \rightarrow 0} \sup_{(\lambda, w) \in B} \|D^\ell F_h(\lambda, w) - D^\ell F(\lambda, w)\| = 0$$

for all  $\ell = 0, 1, 2, \dots$  and  $B \subset H_0^1(\Omega)$  bounded. The assumption (20.4) is verified and we can apply the theory developed in Section 20 at the singular solution  $(\lambda_0, 0)$ .

There exist positive constants  $h_0 \leq 1$ ,  $\eta$ , and a  $C^\infty$  mapping  $g_h : B(0, \gamma) \subset \mathbb{R}^2 \rightarrow \mathcal{X}_1$  such that for all  $(\alpha_1, \alpha_2) \in B(0, \gamma)$ ,  $0 < h \leq h_0$ ,

$$(Q_2 F_h(\lambda_0 + \alpha_1, \kappa_1 + \alpha_2 \xi) = 0 \text{ and } \kappa_1 \in B(0, \eta) \subset \mathcal{X}_1) \iff \kappa_1 = g_h(\alpha_1, \alpha_2).$$

In particular for  $(\alpha_1, \alpha_2) \in B(0, \gamma)$

$$(22.32) \quad Q_2 F_h(\lambda_0 + \alpha_1, g_h(\alpha_1, \alpha_2) + \alpha_2 \xi) = 0,$$

and  $g_h(\alpha_1, 0) = 0$  for  $(\alpha_1, 0) \in B(0, \gamma)$ . Moreover we can derive the following error estimates, see (20.7), (20.8), (12.14), (12.15), (22.25)–(22.29),

$$(22.33) \quad \sup_{(\alpha_1, \alpha_2) \in B(0, \gamma)} |D^\ell g(\alpha_1, \alpha_2) - D^\ell g_h(\alpha_1, \alpha_2)|_{1, \Omega} \leq C_\ell h,$$

where the constants  $C_\ell$ ,  $0 \leq \ell < \infty$ , are independent of  $h$ .

The approximate bifurcation equation reads, for  $(\alpha_1, \alpha_2) \in B(0, \gamma)$

$$(22.34) \quad f_h(\alpha_1, \alpha_2) \equiv (F_h(\lambda_0 + \alpha_1, g_h(\alpha_1, \alpha_2) + \alpha_2 \xi), \xi)_{1, \Omega} = 0.$$

We can explicit with the definitions (22.13) of  $f$  and (22.34) of  $f_h$ , the difference  $f - f_h$ ,

$$\begin{aligned} f(\alpha_1, \alpha_2) - f_h(\alpha_1, \alpha_2) &= (T(g^2(\alpha_1, \alpha_2) - g_h^2(\alpha_1, \alpha_2) + (2\alpha_2 \xi - \alpha_1)(g(\alpha_1, \alpha_2) - g_h(\alpha_1, \alpha_2))), \xi)_{1, \Omega} \\ &\quad + ((T - T_h)((g_h(\alpha_1, \alpha_2) + \alpha_2 \xi)^2 - (\lambda_0 + \alpha_1)(g_h(\alpha_1, \alpha_2) + \alpha_2 \xi)), \xi)_{1, \Omega}. \end{aligned}$$

With the estimates (22.29) and (22.33) we deduce the following

$$(22.35) \quad \sup_{(\alpha_1, \alpha_2) \in B(0, \gamma)} |D^\ell f(\alpha_1, \alpha_2) - D^\ell f_h(\alpha_1, \alpha_2)| \leq C'_\ell h,$$

where the constants  $C'_\ell$ ,  $0 \leq \ell < \infty$ , are independent of  $h$ . Since for all  $(\alpha_1, 0) \in B(0, \gamma)$  we have  $g_h(\alpha_1, 0) = 0$ , it is not difficult to check with the definition of  $f_h$  in (22.34) that for all  $(\alpha_1, 0) \in B(0, \gamma)$

$$(22.36) \quad f_h(\alpha_1, 0) = 0,$$

which is the analogue of (22.19).

To study the approximate bifurcation equation, we are exactly in the framework presented in Section 21. The functions  $f$ ,  $f_h$  are of class  $C^\infty$  defined on  $B(0, \gamma) \subset \mathbb{R}^2$  with values in  $\mathbb{R}$ , see (22.13) and (22.34). We have checked that

$$\det(D^2 f(0, 0)) = -(1/\lambda_0)^2 < 0,$$

so  $(\lambda_0, 0)$  is a simple bifurcation point. Moreover the difference  $D^\ell f - D^\ell f_h$  is estimated in (22.35).

We can apply Theorem 21.2. For  $h$  small enough and only to restrict  $\gamma$ , there exists a  $C^\infty$  diffeomorphism from  $\mathcal{S}_{f_h} = \{(\alpha_1, \alpha_2) \in B(0, \gamma) \subset \mathbb{R}^2; f_h(\alpha_1, \alpha_2) = 0\}$  onto a

part of a hyperbola centred at  $(\alpha_{1h}, \alpha_{2h})$ . Going back to definition of  $f_h$ , we can check that the center of the hyperbola is in fact  $(\alpha_{1h}, 0)$  with

$$f_h(\alpha_{1h}, 0) = 0 \quad \text{and} \quad Df_h(\alpha_{1h}, 0) = 0,$$

which means that the hyperbola is degenerate and we are in the case of Figure 21.4.

From the estimates (21.18) and (22.35), we have

$$(22.37) \quad |\alpha_{1h}| \leq Ch.$$

The estimates (21.25) and (22.35) lead to

$$(22.38) \quad \delta(\mathcal{S}_f, \mathcal{S}_{f_h}) \leq Ch,$$

and the estimates (20.11), (22.29), and (22.38) imply

$$(22.39) \quad \delta(\mathcal{S}, \mathcal{S}_h) \leq Ch,$$

where  $\mathcal{S}$  and  $\mathcal{S}_h$  are the solution sets for  $F$  and  $F_h$  introduced in Section 20.

*Remark 22.1.* In fact we can check that there exists a  $\xi_h \in V_h$ ,  $\xi_h \neq 0$  such that

$$\text{for all } v_h \in V_h \quad \int_{\Omega} \mathbf{grad} \xi_h \mathbf{grad} v_h \, dx = (\lambda_0 + \alpha_{1h}) \int_{\Omega} \xi_h v_h \, dx.$$

It is known from the theory of approximation of eigenproblems, see BABUŠKA and OSBORN [1991], that the following estimate holds

$$|\alpha_{1h}| \leq Ch^2.$$

The estimate (22.37) is not optimal. Optimal estimates are derived with a different method in CROUZEIX and RAPPAZ [1989].  $\square$

*Remark 22.2.* If we take polynomials of degree  $k > 1$  on each triangle in the definition (22.20) of  $V_h$ , then we cannot improve the estimate (22.33) but the estimate (22.39) is modified into

$$\delta(\mathcal{S}, \mathcal{S}_h) \leq Ch^k. \quad \square$$

### 23. Bibliographical comments

The brief comments below are intended only to outline the wealthy literature which is related to our work. A lot of additional results could have been included or quoted.

In Sections 18 through 22 we describe the different steps to study approximations of singular solutions. Only the main features of the method are developed and the example of Section 22 is elementary but it illustrates the general theory. An example of a simple bifurcation point not on the trivial branch is studied into details in CAUSSIGNAC,

DESCLOUX, and RAPPAZ [1987] and an example of bifurcation for a problem on an unbounded domain is presented in DESCLOUX and RAPPAZ [1982], see also CROUZEIX and RAPPAZ [1989].

Bifurcation phenomena often occur in parameter dependent nonlinear problems. In KELLER and ANTMAN [1969] different examples are presented. General results on bifurcation problems can be found in the books of CHOW and HALE [1982], GOLUBITSKY and SCHAEFFER [1985], GOLUBITSKY, STEWART, and SCHAEFFER [1988], IOSS and JOSEPH [1981], NIRENBERG [1974], and SATTINGER [1973] [1979], and in the bibliographies therein.

The approximation of parameter dependent problems has been widely studied in the literature. For a general approach we can mention the works of BREZZI, RAPPAZ, and RAVIART [1980] [1981a] [1981b], CROUZEIX and RAPPAZ [1989], FUJI and YAMAGUTI [1980], KELLER [1977], KIKUCHI [1976], RABIER [1985], RHEINBOLDT [1986] and SIMPSON [1972].

Further developments concerning turning point can be found for instance in GRIEWANK and REDDIEN [1984], KIKUCHI [1979], and PAUMIER [1981], and concerning bifurcation can be found in RAPPAZ and RAUGEL [1981] for bifurcation from multiple eigenvalues, in BERNARDI [1982] and DESCLOUX [1984] for Hopf bifurcation. A generalization in the direction of non differentiable problems was mentioned in Remark 2.4.

When the exact problem and its approximation are covariant with respect to representations of a finite group, see SATTINGER [1979], GOLUBITSKY and SCHAEFFER [1985], GOLUBITSKY, STEWART, and SCHAEFFER [1988], then the study can be sharpened, see CROUZEIX and RAUGEL [1988], DELLNITZ and WERNER [1989], RAUGEL [1984] for instance.

The applications of the general approximation theory are numerous and quite diverse. For instance we can find an account on the analysis of discretization of the Navier-Stokes equations in GIRAULT and RAVIART [1986]. Approximations of the von Kármán equations, see CIARLET and RABIER [1980], can be found in BREZZI, RAPPAZ, and RAVIART [1980] [1981b], KESAVAN [1979] [1980], REINHARDT [1982].

In our work we have developed a general method based on the inverse and the implicit function theorems to study approximations of regular solutions of nonlinear problems and approximations of regular and singular solutions of parametrized nonlinear problems. There are other methods to study nonlinear problems, which are not even mentioned here. For instance the ones based on the maximum principle or on monotonicity are well-suited to variational inequalities, see for instance BARRETT and ELLIOTT [1991], GLOWINSKI [1984], CONRAD and CORTEY-DUMONT [1987], ŽENÍŠEK [1990].

There is a wide variety of algorithms to solve nonlinear approximate problems. Some fundamental results are due to Keller, see KELLER [1970], KEENER and KELLER [1974], and KELLER [1987]. For an overview we refer to the conference proceedings edited by KUEPPER, MITTELMANN, and WEBER [1984], KUEPPER, SEYDEL, and TROGER [1987], MITTELMANN and WEBER [1980]. There is an extensive bibliography on the topic of numerical continuation methods in the book of ALLGOWER and GEORG [1990].



## REFERENCES

- Adams, R. (1975), *Sobolev spaces*, Academic Press, New York.
- Allgower, E.L. and Georg K. (1990), *Numerical continuation methods. An introduction*, Springer Series in Computational Mathematics, Springer-Verlag, Berlin.
- Amann, H. (1976), *Fixed point equations and nonlinear eigenvalue problems in ordered Banach space*, SIAM Review **18**, 620–709.
- Babuška, I. and Aziz, A.K. (1972), *Survey lectures on the mathematical foundations of the finite element method*, in The mathematical foundations of the finite element method with application to partial differential equations (Aziz, A.K., eds.), Academic Press, New York, pp. 3–359.
- Babuška, I. and Osborn, J. (1991), *Eigenvalue problems*, in Handbook of Numerical Analysis, vol. II, North-Holland, Amsterdam, pp. 641–787.
- Babuška, I. and Rheinboldt, W.C. (1982), *Computational error estimates and adaptive processes for some nonlinear structural problems*, Comp. Methods Appl. Mech. Eng. **34**, 895–937.
- Baranger, J. and El Amri H. (1991), *Estimateur a posteriori d'erreur pour le calcul adaptatif d'écoulements quasi-newtoniens*, M.<sup>2</sup>A.N. **25**, 31–48.
- Barrett, J.W. and Elliott, C.M. (1989), *Finite element approximation of a plasma equilibrium problem*, IMA J. Num. Anal. **9**, 443–464.
- Barrett, J.W. and Elliott, C.M. (1991), *Finite element approximation of a free boundary problem arising in the theory of liquid drops and plasma physics*, M.<sup>2</sup>A.N. **25**, 213–252.
- Bernardi, C. (1982), *Approximation of Hopf bifurcation*, Numer. Math. **39**, 14–37.
- Bonic, R.A. (1969), *Linear functional analysis*, Notes on mathematics and its applications, Gordon and Breach, New York.
- Brezis, H. (1983), *Analyse fonctionnelle. Théorie et applications*, Collection mathématiques appliquées pour la maîtrise, Masson, Paris.
- Brezzi, F., Rappaz, J., and Raviart, P.A. (1980), *Finite dimensional approximation of nonlinear problems. Branches of nonsingular solutions*, Numer. Math. **36**, 1–36.
- Brezzi, F., Rappaz, J., and Raviart, P.A. (1981a), *Finite dimensional approximation of nonlinear problems. Limit points*, Numer. Math. **37**, 1–28.
- Brezzi, F., Rappaz, J., and Raviart, P.A. (1981b), *Finite dimensional approximation of nonlinear problems. Simple bifurcation points*, Numer. Math. **38**, 1–30.
- Caloz, G. (1987), *Simulation numérique des équilibres d'un plasma dans un tokamak: modélisation et études mathématiques*, Thèse **650**, École Polytechnique Fédérale de Lausanne.
- Caloz, G. (1991), *Approximation by finite element method of the model plasma problem*, M.<sup>2</sup>A.N. **25**, 213–252.
- Caloz, G. (1994), *Stability on a regular branch of solutions*, Tech. Rpt. 94, I.R.M.A.R., Université de Rennes I.
- Cartan, H. (1967), *Cours de calcul différentiel*, Collection méthodes, Hermann, Paris.
- Caussignac, Ph., Descloux, J., and Rappaz, J. (1987), *Study of an elliptic problem with nonlinear boundary conditions*, Math. Met. in the Appl. Sci. **9**, 261–275.

- Chow, S.N. and Hale, J. (1982), *Methods of bifurcation theory*, Grundlehren 251, Springer-Verlag, Berlin.
- Ciarlet, P.G. (1978), *The finite element method for elliptic problems*, Studies in mathematics and its applications, North-Holland, Amsterdam.
- Ciarlet, P.G. (1991), *Basic error estimates for elliptic problems*, in Handbook of Numerical Analysis, vol. II, North-Holland, Amsterdam, pp. 17–351.
- Ciarlet, P.G. and Rabier, P. (1980), *Les équations de von Kármán*, Lecture notes in mathematics, 826, Springer-Verlag, Berlin.
- Clément, P. (1975), *Approximation by finite element functions using local regularization*, R.A.I.R.O. Analyse Numérique **2**, 77–84.
- Conrad, F. and Cortey-Dumont, P. (1987), *Nonlinear eigenvalue problems in elliptic variational inequalities: some results for the maximal branch*, Numer. Funct. Anal. and Optimiz. **9**, 1059–1114.
- Crandall, M.G. and Rabinowitz, P.H. (1973), *Bifurcation, perturbation of simple eigenvalues and linearized stability*, Arch. Rational Mech. Anal. **52**, 161–180.
- Crandall, M.G. and Rabinowitz, P.H. (1975), *Some continuation and variational methods for positive solutions of nonlinear eigenvalue problems*, Arch. Rational Mech. Anal. **58**, 241–269.
- Crouzeix, M. and Rappaz, J. (1989), *On numerical approximation in bifurcation theory*, RMA 13, Masson, Paris.
- Crouzeix, M. and Raugel, G. (1988), *Invariance under the dihedral group and application to bifurcation problems*, Nonlinear Analysis **12**, 75–99.
- Dauge, M. (1992), *Neumann and mixed problems on curvilinear polyhedra*, Integ. Equat. Op. Th. **15**, 227–261.
- Dautray, R. and Lions, J.L. (1987), *Analyse mathématique et calcul numérique pour les sciences et les techniques*, Masson, Paris.
- Dellnitz, M. and Werner, B. (1989), *Computational methods for bifurcation problems with symmetries*, J. Comput. Appl. Math. **26**, 97–123.
- Descloux, J. (1984), *On Hopf and subharmonic bifurcations*, in Numerical methods for bifurcation problems, Birkhauser-Verlag, Basel, pp. 145–161.
- Descloux, J. and Rappaz, J. (1982), *Approximation of solution branches of nonlinear equations*, R.A.I.R.O. Analyse Numérique **16**, 319–349.
- Dieudonné, J. (1968), *Éléments d'analyse*, Tome I, Gauthier-Villars, Paris.
- Dunford, N. and Schwartz, J.T. (1958), *Linear operators*, Interscience, New York.
- Fink, J.P. and Rheinboldt, W.C. (1983a), *On the discretization error of parametrized nonlinear equations*, SIAM J. Numer. Anal. **20**, 732–746.
- Fink, J.P. and Rheinboldt, W.C. (1983b), *Solution manifolds and submanifolds of parametrized equations and their discretization errors*, Tech. Rpt. ICMA-83-59, Institute for computational mathematics and applications, University of Pittsburgh (1984), Numer. Math. **45**, 323–343.
- Fink, J.P. and Rheinboldt, W.C. (1986), *Folds on the solution manifold of a parametrized equation*, SIAM J. Numer. Anal. **23**, 693–706.
- Fink, J.P. and Rheinboldt, W.C. (1987), *A geometric framework for the numerical study of singular points*, SIAM J. Numer. Anal. **24**, 618–633.
- Fuji, H. and Yamaguti, M. (1980), *Structure of singularities and its numerical realization in nonlinear elasticity*, J. Math. Kyoto Univ. **20**, 489–590.
- Gilbarg, D. and Trudinger, N.S. (1977), *Elliptic partial differential equations of second order*, Springer-Verlag, Berlin.
- Girault, V. and Raviart, P.A. (1982), *An analysis of upwind schemes for the Navier-Stokes equations*, SIAM J. Numer. Anal. **19**, 312–333.
- Girault, V. and Raviart, P.A. (1986), *Finite element methods for Navier-Stokes equations*, Springer-Verlag, Berlin.
- Golubitsky, M. and Schaeffer, D.G. (1985), *Singularities and groups in bifurcation theory*, Volume 1, Springer-Verlag, Berlin.

- Golubitsky, M., Stewart, I., and Schaeffer, D.G. (1988), *Singularities and groups in bifurcation theory*, Volume 2, Springer-Verlag, Berlin.
- Griewank, A. and Reddien, G.W. (1984), *Characterization and computation of generalized turning points*, SIAM J. Numer. Anal. **21**, 176–185.
- Grisvard, P. (1985), *Elliptic problems in nonsmooth domains*, Pitman, Boston.
- Ioss, G. and Joseph, D.D. (1981), *Elementary stability and bifurcation theory*, Springer-Verlag, Berlin.
- Joseph, D.D. (1965), *Non-linear heat generation and stability of the temperature distribution in conduction solids*, Int. J. Heat Mass Transfer **8**, 281–288.
- Keener, J.P. and Keller, H.B. (1974), *Perturbed bifurcation theory*, Arch. Rational Mech. Anal. **50**, 159–175.
- Keller, H.B. (1970), *Nonlinear bifurcation*, J. Diff. Equ. **7**, 417–434.
- Keller, H.B. (1977), *Numerical solution of bifurcation and nonlinear eigenvalue problems*, in Application of bifurcation theory (P.H. Rabinowitz, eds.), Academic Press, pp. 359–384.
- Keller, H.B. (1987), *Lectures on numerical methods in bifurcation problems*, Springer-Verlag, Berlin.
- Keller, H.B. and Antman, S., Eds (1969), *Bifurcation theory and nonlinear eigenvalue problems*, W.A. Benjamin, Inc., New York.
- Keller, H.B. and Cohen, D.S. (1967), *Some positive problems suggested by nonlinear heat generation*, J. Math. Mech. **16**, 1361–1376.
- Kesavan, S. (1979), *La méthode de Kikuchi appliquée aux équations de von Kármán*, Numer. Math. **32**, 209–232.
- Kesavan, S. (1980), *Une méthode d'éléments finis mixtes pour les équations de von Kármán*, R.A.I.R.O. Analyse Numérique **14**, 149–173.
- Kikuchi, F. (1976), *An iterative finite element scheme for bifurcation analysis of semilinear elliptic equations*, Report **542**, Inst. Space Aero. Sc., Tokyo University.
- Kikuchi, F. (1979), *Finite element approximations to bifurcation problems of turning point type*, Theoretical and Applied Mechanics **27**, 99–144.
- Kuepper, T., Mittelman, H.D., and Weber, H., Eds (1984), *Numerical methods for bifurcation problems*, ISNM, 70, Birkhauser-Verlag, Basel.
- Kuepper, T., Seydel, R., and Troger, H., Eds (1987), *Bifurcation: analysis, algorithms, applications*, ISNM, 79, Birkhauser-Verlag, Basel.
- Leray, J. and Schauder, J. (1934), *Topologie et équations fonctionnelles*, Ann. Sci. École Norm. Sup. **51**, 45–78.
- Lions, J.L. and Magenes, E. (1968), *Problèmes aux limites non homogènes et applications*, Travaux et recherches mathématiques, Dunod, Paris.
- Lions, P.L. (1982), *On the existence of positive solutions of semilinear elliptic equations*, SIAM Review **24**, 441–467.
- Maz'ya, V.G. (1985), *Sobolev spaces*, Springer, Berlin.
- Mignot, F. and Puel, J.P. (1980), *Sur une classe de problèmes non linéaires avec non linéarité positive, croissante, convexe*, Comm. P.D.E. **5**, 791–836.
- Mittelman, H.D. and Weber, H., Eds (1980), *Bifurcation problems and their numerical solution*, ISNM, 54, Birkhauser-Verlag, Basel.
- Nečas, J. (1967), *Les méthodes directes en théorie des équations elliptiques*, Masson, Paris.
- Nijenhuis, A. (1974), *Strong derivatives and inverse mappings*, Amer. Math. Monthly **81**, 969–980.
- Nirenberg, L. (1974), *Topics in nonlinear functional analysis*, Courant Institute of Mathematical Science, New York University, New York.
- Paumier, J.C. (1981), *Une méthode numérique pour le calcul de points de retournement. Application à un problème aux limites non-linéaires*, Numer. Math. **37**, 433–452.
- Picasso, M. (1992), *Simulation numérique des traitements de surface par laser*, Thèse **1011**, École Polytechnique Fédérale de Lausanne.
- Pousin, J. and Rappaz, J. (1991), *Consistance, stabilité, erreurs a priori et a posteriori pour des problèmes non linéaires*, C.R. Acad. Sc. Paris, Série I **312**, 699–703.

- Pousin, J. and Rappaz, J. (1992), *Consistency, stability, a priori and a posteriori errors for Petrov-Galerkin methods applied to nonlinear problems*, Tech. Rpt. 04-92, D.M.A., École Polytechnique Fédérale de Lausanne; to appear in Numer. Math..
- Rabier, P. (1985), *Topics in one-parameter nonlinear problems*, Tata Institute Lecture Notes, Springer-Verlag, Berlin.
- Rannacher, R. and Scott, R. (1982), *Some optimal error estimates for piecewise linear finite element approximations*, Math. Comp. **38**, 437–445.
- Rappaz, J. (1983), *Estimations d'erreur dans différentes normes pour l'approximation de problèmes de bifurcation*, C.R. Acad. Sc. Paris, Série I **296**, 179–182.
- Rappaz, J. (1984), *Approximation of a nondifferentiable nonlinear problem related to MHD equilibria*, Numer. Math. **45**, 117–133.
- Rappaz, J. and Raugel, G. (1982), *Approximation of double bifurcation points for nonlinear eigenvalue problems*, MAFELAP 1981, J. Whiteman editor, Academic Press, New York.
- Raugel, G. (1984), *Approximation numérique de problèmes non linéaires*, Thèse, Université de Rennes I.
- Reinhart, L (1982), *On the numerical analysis of the von Kármán equations*, Numer. Math. **39**, 371–404.
- Rheinboldt, W.C. (1981), *Numerical analysis of continuation methods for nonlinear structural problems*, Comput. and Structures **13**, 103–116.
- Rheinboldt, W.C. (1985), *Error estimates for nonlinear finite element computations*, Comput. and Structures **20**, 91–98.
- Rheinboldt, W.C. (1986), *Numerical analysis of parametrized nonlinear equations*, Wiley-Interscience, John Wiley and Sons.
- Sattinger, D.H. (1973), *Topics in stability and bifurcation theory*, Lecture notes in mathematics, 309, Springer-Verlag, Berlin.
- Sattinger, D.H. (1979), *Group theoretic methods in bifurcation theory*, Lecture notes in mathematics, 762, Springer-Verlag, Berlin.
- Simader, C.G. (1972), *On Dirichlet's boundary value problem*, Lecture notes in mathematics, 268, Springer-Verlag, Berlin.
- Simpson, R.B. (1972), *Existence and error estimates for solutions of a discrete analog of nonlinear eigenvalue problems*, Math. Comp. **26**, 359–375.
- Temam, R. (1977), *Theory and numerical analysis of the Navier-Stokes equations*, North-Holland, Amsterdam.
- Verfürth, R. (1993), *A review of a posteriori error estimation and adaptive mesh-refinement techniques*, Tech. Rpt., Institut für Angewandte Mathematik, Universität Zürich.
- Wahlbin, L.R. (1978), *Maximum norm error estimates in the finite element method with isoparametric quadratic elements and numerical integration*, R.A.I.R.O. Analyse Numérique **12**, 173–202.
- Ženíšek, A. (1990), *Nonlinear elliptic and evolution problems and their finite element approximations*, Computational mathematics and applications, Academic Press Limited, London.