

The goal of this mini-course is to provide an introduction into finite volume strategies and convergence analysis of the nonlinear diffusion operators of the p -laplacian kind.

- Contents
- Leray-Lions kind operators, key properties
 - Elliptic problems with Leray-Lions operators: well-posedness, variational nature, regularity of solutions
 - meshes of 2D domains, dual meshes, boundary cells discrete spaces, discrete divergence operator
 - strategies for gradient reconstruction, 2D lemmas, discrete duality (covolume scheme, DDFV)
 - finite volume schemes with Leray-Lions operators, discrete weak formulation
 - discrete solutions: uniqueness, a priori estimates, existence (Brouwer, topological degree), maximum principle
 - discrete energy functional, minimization, practical resolution of the discrete equations
 - projection of test functions, consistency properties
 - discrete compactness (H^1 weak, L^1 strong)
 - discrete Poincaré inequality and embedding inequalities
 - passage to the limit: Minty trick, strong convergence
 - passage to the limit: use of the Young measures
 - basic strategy for error estimates for $C^2/W^{2,p}$ solutions
 - consistency with Sobolev functions
 - superconvergence estimates for $C^4/W^{4,1}$ solutions
 - refinements using Barrett-Liu arguments
 - discrete Besov framework on uniform cartesian meshes
 - Chow error estimates in the Besov context, optimality
 - some words on different numerical strategies

I Def (Leray-Lions operators)

Let $\vec{a}: \mathbb{R}^d \rightarrow \mathbb{R}^d$, continuous, st for some $p \in (1, \infty)$

• $\vec{a}(\vec{\xi}) \cdot \vec{\xi} \geq \frac{1}{C} |\vec{\xi}|^p$ (coercivity)

• $|\vec{a}(\vec{\xi})|^{p'} \leq C(1 + |\vec{\xi}|)^p$ (growth)

• $(\vec{a}(\vec{\xi}) - \vec{a}(\vec{\eta})) \cdot (\vec{\xi} - \vec{\eta}) \geq 0$ (monotonicity)

Then the operator $A: W^{1,p}(\Omega) \rightarrow W^{-1,p'}(\Omega)$
 $u \mapsto -\text{div } \vec{a}(\nabla u) + BC$

is called Leray-Lions operator.

The prototype of Leray-Lions ops is the p-laplacian: $\vec{a}(\vec{\xi}) = |\vec{\xi}|^{p-2} \vec{\xi}$.

This is also an example of operator identity from a scalar potential:

namely, $\vec{a}(\vec{\xi}) = \nabla_{\vec{\xi}} \Phi(\vec{\xi})$, where $\Phi(\vec{\xi}) = \frac{1}{p} |\vec{\xi}|^p$. In this case, \vec{a} is cyclically monotone.

The p-laplacian corresponds to a strictly monotone \vec{a} ;

moreover, \vec{a} and \vec{a}^{-1} are (locally) Hölder continuous (we'll see later the precise meaning).

Reference Lions "Quis pbs aux limites non linéaires",
 Chap. "Compacité + monotonie".

Model problem $\Omega \subset \mathbb{R}^d$ bdd domain

(MPL) $\begin{cases} u - \text{div } \vec{a}(\nabla u) = f & \text{in } \Omega \\ + \text{homogeneous Neumann (zero-flux) BC: } \vec{a}(\nabla u) \cdot \vec{n} = 0 \text{ on } \partial\Omega \end{cases}$

Mathematical interpretations:

weak sol for $f \in L^{p'}(\Omega)$ (more generally, $f \in W^{-1,p'}(\Omega)$)

$u \in W^{1,p}(\Omega) \cap L^2(\Omega)$

$\forall \varphi \in \tilde{W}^{1,p}(\Omega)$ ($C^\infty(\bar{\Omega})$ is enough)

(W) $\int_{\Omega} u \varphi + \vec{a}(\nabla u) \cdot \nabla \varphi = \int_{\Omega} f \varphi$

Optimization

$u = \text{arg min}_{v \in W^{1,p}(\Omega)} \mathcal{J}[v], \quad \mathcal{J}: v \mapsto \frac{1}{2} \int_{\Omega} (u-f)^2 + \int_{\Omega} \Phi(\nabla v)$

Actually, \mathcal{J} is a strictly convex lower semi-continuous coercive functional (thus the minimization problem has a unique solution), and (W) is the associated Euler-Lagrange equation for critical points.

It is clear that the minimum is unique; to show uniqueness for (W) take u, \hat{u} two solutions, and take $\varphi = u - \hat{u}$ for the test function.

Then $\int |u - \hat{u}|^2 \leq \int |u - \hat{u}|^2 + (\vec{a}(\nabla u) - \vec{a}(\nabla \hat{u})) \cdot (\nabla u - \nabla \hat{u}) = \int (f - f)(u - \hat{u}) = 0$.

Existence of a solution will be the byproduct of the numerical method; stability can be obtained from the above calculation, but L^1 stability is more interesting.

One can note the following properties:

- $\int_{\Omega} u^2 + \int_{\Omega} |u|^p \leq C(\|f\|_{L^p})$ (test function $\varphi = u$ + the Poincaré inequality: $\|u\|_{L^p} \leq C(p, \Omega)(\|u\|_{L^1} + \|u\|_{L^2})$)
- $\int |u| \leq \int |f|$ (test functions approximating $\text{sign}(u)$), $\int |u|^q \leq \int |f|^q$ (test fcts approx. $|u|^{q-2} u$, then Young's ineq.)
- $\int (u - \hat{u})^{\pm} \leq \int (f - \hat{f})^{\pm}$ (idem, with $\text{sign}^{\pm}(u - \hat{u})$) \leadsto Δ contraction and comparison principle
- maximum principle: $f \leq M$ a.e. $\Rightarrow u \leq M$ a.e. (idem, with $\text{sign}^+(u - M)$; comparison with constant solution)

The accurate proof of the three latter properties uses regularization, but in the discrete context things are simpler.

The issue of regularity of solutions, essential for error estimates, is much more delicate than in the linear case. Available results are:

- Hölder regularity (Lieberman, '1988)
- $W^{2,p}$ regularity for $p \leq 2$, and some Besov space $(B_{\infty}^{1-\frac{1}{p-1}, p})$ regularity for $p > 2$.

II

In numerical analysis, one seeks for "good" ways to approach solutions (here, of PDE). "Good" may mean many things:

- error of approximation as small as possible
- computational complexity as low as possible
- qualitative properties of solutions (e.g. maximum principle, comparison and contraction principles, ^{like with optimization problem} preserved at the discrete level) *qualitative properties of differential operators (discrete duality)*
- possibility to work on "general meshes" (complex geometries, mesh refinement)

For a theoretical analysis, a minimal requirement is: consistency + stability of the method. Convergence is sometimes more difficult to obtain.

Here, we will explore methods that are

- convergent
- preserve many of the structure properties of operators
- work on general meshes, and preserve comparison properties on "orthogonal" meshes.

The discrete problem includes a finite number of real values (called degrees of freedom). They may have different interpretation (coefficients of a finite element basis, spectral coefficients); in finite volumes context the degrees of freedom u_K are associated to volumes K (eventually, via the "centers" x_K of the volumes).

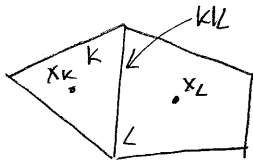
Therefore the discretization starts with a mesh of Ω . *Dim. 1 is too easy, dim. 3 is too hard, take dimension 2.*

We assume that $\Omega \subset \mathbb{R}^2$ is polygonal, and look on polygonal meshes. (4)

Description

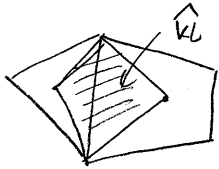
$K \in \mathcal{M}$: volumes

x_K : center of K (barycenter: inside K if K convex
circumcenter: may lie outside K)



KL : interface between neighbours ($L \in \mathcal{N}(K)$), $d_{KL} = |x_K - x_L|$

\hat{KL} : "diamond"



the mesh is orthogonal if $\overrightarrow{x_K x_L} \perp KL$

Important: $\overrightarrow{x_K x_L}$ should form an acute angle with \vec{n}_K .

Approx. of $u - \text{div } \vec{F} = f$

We need one unknown u_K per volume K , which is interpreted (\equiv which approximate) both as the mean value of u on K and as the point value of u at x_K .

One "integrates" the PDE on K and uses the Green-Gauss formula!

$$m_K u_K - \int_{\partial K} \vec{F} \cdot \vec{n}_K = m_K f_K \quad \text{where } f_K = \int_K f \quad (\text{or } f_K = f(x_K), \text{ anyway quadratures are needed to evaluate } f_K)$$

Thus we need to approximate

$$\int_{\partial K} \vec{F} \cdot \vec{n} = \sum_{L \in \mathcal{N}(K)} m_{KL} (f_{KL} \vec{F} \cdot \vec{n}_{KL}) \quad ; \quad \text{one needs the values } F_{KL} = \int_{KL} \vec{F} \cdot \vec{n}_{KL}$$

In our context, $\vec{F} = \vec{\sigma}(\nabla u)$; for a discrete counterpart, we need to relate F_{KL} to $(u_K)_K$. There are different strategies for doing this:

- (mixed formulations) put more equations keeping F_{KL} as unknowns
- (finite volume schemes) reconstruct F_{KL} "by hands" from the values u_K, u_L and may be some neighbour values.

Case $\vec{\sigma} = \text{Id}$, orthogonal mesh

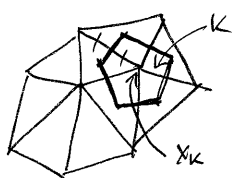
$$F_{KL} = \frac{u_K - u_L}{d_{KL}}$$

Case of the p-Laplacian Even on an orthogonal mesh, one seems to need

the whole vector \vec{F}_{KL} on $KL \dots$

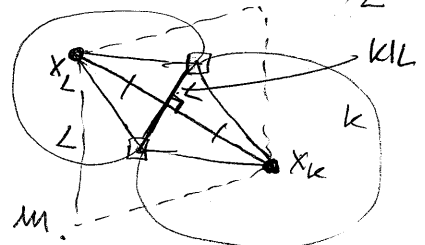
III) Covolume scheme on triangle mesh

Quite often, we are given a partition of Ω into triangles satisfying the so-called Delaunay condition. With $x_{K^*} = \text{circumcenter of } K^*$, we get an orthogonal mesh... Instead of using this for FV discretization, we will put the unknowns into the vertices of this mesh, and consider \mathcal{T} is the dual to the mesh \mathcal{M} defined as follows (the Voronoi mesh); for x_K a vertex of \mathcal{T} , $K \in \mathcal{M}$ is the set of $x \in \Omega$ st $|x - x_K| = \min_L |x - x_L|$.



This mesh is also orthogonal, making zoom on a diamond we have:

- : vertices of $\mathcal{T} \equiv$ centers of \mathcal{M}
- : (circum)centers of $\mathcal{T} \equiv$ vertices of \mathcal{M} .



Thus our unknowns are located at the vertices of the triangles constituting the mesh \mathcal{T} ; then, as in the finite element context, it is tempting to use affine / triangle interpolation of these values to produce a function on Ω . Its gradient is constant per triangle; in particular, on each interface KL it takes two values.

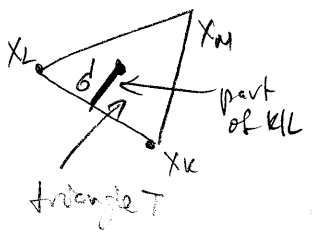
Then we can use this to produce the following scheme:

$$DEq) \forall K \in \mathcal{M}, \quad m_K u_K - \sum_{\delta \in \mathcal{K}(K)} m_\delta \vec{a}(\nabla_{T(\delta)} u_h) \cdot \vec{n}_\delta = m_K f_K$$

NB Formally, at $\delta \in \mathcal{K}(K)$ we set the flux to zero; thus the BC is taken into account!

Here δ 's are half-interfaces included in ∂K , and $T = T(\delta)$ is the triangle containing δ ; $u_h = (u_K)_K$ is the set of all unknowns; $m_K = \text{meas}_2(K)$, $m_\delta = \text{meas}_1(\delta)$; and

$\nabla_T u_h$ = the gradient of the unique affine function obtained by interpolation of the values u_K, u_L, u_M at the vertices x_K, x_L, x_M of the triangle $T(\delta)$, for $\delta \in \mathcal{K}(K)$.



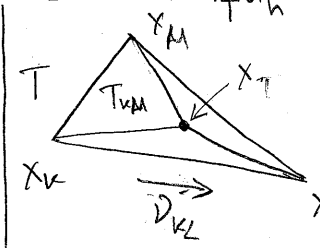
This gives a set of $\#K$ eqns on $\#K$ unknowns, which looks reasonable:

- u_K interpreted as values at (and near) x_K
- starting from $(u_K)_K$, the approximation of ∇u (thus of $\vec{a}(\vec{\sigma}u)$) is produced in an "intuitively consistent" way
- starting from the approximate fluxes replacing $\vec{a}(\vec{\sigma}u)$ the approximation of $\text{div } \vec{a}(\vec{\sigma}u)$ is consistent both mathematically and physically (balance law respected by construction); in particular, the scheme is "locally conservative".

To write the system explicitly, one needs

- geometric quantities from the mesh (m_K, m_δ, \vec{n}_K)
- pointwise values of gradients of f (to get f_K)
- an explicit formula for $\nabla_T u_h$ in terms of u_K, u_L, u_M and x_K, x_L, x_M .

Lemma Prop
$$\nabla_T u_h = \frac{2}{m_T} \left(m_{TKL} \frac{u_L - u_K}{d_{KL}} \vec{v}_{KL} + m_{TLM} \frac{u_M - u_L}{d_{LM}} \vec{v}_{LM} + m_{TMK} \frac{u_K - u_M}{d_{MK}} \vec{v}_{MK} \right)$$



where T_{KL}, T_{LM}, T_{MK} are triangles with vertex $x_T =$ circumcenter of T and the opposite side KL, LM, MK , respectively, and m_{TKL}, \dots are the signed measures of these triangles.

Proof: can be obtained by elementary tools (same theorem, etc.)

• generalization and proof will be given later.

Heuristic explanation: If $\vec{r} = \vec{\nabla}_T u_h$, then $\frac{u_h - u_k}{d_{kk}} \vec{r}_{kk} = \text{Proj}_{\vec{r}_{kk}} \vec{r}$.

The formula reads $\text{Id} = 2 \left(\alpha \text{Proj}_{\vec{r}_{kL}} + \dots + \gamma \text{Proj}_{\vec{r}_{kR}} \right)$, $\alpha + \beta + \gamma = 1$. (*)

Because projection takes "in the mean" one direction out of two, the role of the coefficient 2 is clear (say, if \vec{r} is random, then $E[\vec{r}] = 2 E[\alpha \text{Proj}_{\vec{r}_{kL}} \vec{r} + \gamma \text{Proj}_{\vec{r}_{kR}} \vec{r}]$. The fact that we have the identity (*) is due to the choice of the point $x_T =$ the circumcenter of $T \dots$ yes!! we'll see.

IV) Now we have written the scheme; to analyse it, we need much more. For the analysis of convergence, we have to place u (a function on Ω , in $W^1P(\Omega) \cap \dots$) and $u_h = (u_k)_k \in \mathbb{R}^{\#K}$ at the same level; that is, either "project" u on the mesh, or else "lift" $(u_k)_k$ to a function on Ω . (cf. the way $(P_k)_k$ were obtained from f) $(\nabla_T u_h)_T \dashv \dots \dashv \dots$. If we project u , we still need some M -independent (or rather M -uniform) way to measure the difference $u_h - P_h u$ (actually, this approach is used for error estimates); thus we still need to interpret discrete functions in some M -independent way. Therefore lifting operators seem necessary ...

Def L • Given $u_h = (u_k)_{k \in M}$, $\mathcal{L}_h u_h = \begin{bmatrix} \Omega \rightarrow \mathbb{R} \\ x \mapsto \sum_k u_k \mathbb{1}_k(x) \end{bmatrix}$
 In the sequel, $(\mathcal{L}_h u_h)(x)$ will be simply (=abusively) denoted by $u_h(x)$.
 • Given $\vec{F}_h = (\vec{F}_T)_{T \in \mathcal{T}}$, $\mathcal{L}_h \vec{F}_h = \begin{bmatrix} \Omega \rightarrow \mathbb{R}^2 \\ x \mapsto \sum_T \vec{F}_T \mathbb{1}_T(x) \end{bmatrix}$
 $(\mathcal{L}_h \vec{F}_h)(x)$ will be denoted $\vec{F}_h(x)$.
 In particular, if $\vec{F}_T = \vec{\nabla}_T u_h$, we write $(\vec{\nabla}_h u_h)(x)$.

The analysis of the continuous equation been based on multiplication of the equation by test functions (in some duality framework), we'll probably need the same in the discrete setting. But because everything is finite-dimensional here, we can always use the " L^2 - L^2 " multiplication. Thus, what we need is a discrete scalar product related to $\int_{\Omega} uv$ or to $\int_{\Omega} \vec{F} \cdot \vec{G}$. As we have seen, in our discrete framework scalars and vectors "live" on different meshes: $u_h = (u_k)_k \in \mathbb{R}^{\#K}$ and $\vec{F}_h = (\vec{F}_T)_T \in (\mathbb{R}^2)^{\#T}$. Then we can use the above definition to set


Def SP
 $\left[\begin{array}{l} \langle u_h, v_h \rangle_h = \sum_k m_k u_k v_k \quad (\equiv \int_{\Omega} u_h(x) v_h(x) dx) \\ \langle \vec{F}_h, \vec{G}_h \rangle_h = \sum_T m_T \vec{F}_T \cdot \vec{G}_T \quad (\equiv \int_{\Omega} \vec{F}_h(x) \cdot \vec{G}_h(x) dx) \end{array} \right.$

Now, recall our discrete equations: $\forall k \quad u_k v_k - \sum_{\substack{\sigma \subset \partial k \\ \sigma \neq \partial \Omega}} m_\sigma \vec{a}(\vec{\nabla}_{T(\sigma)} u_h) \cdot \vec{n}_k = u_k f_k$
 Take some $v_h = (v_k)_k$ and multiply the eqns term by term:

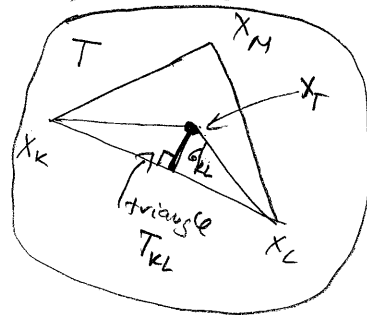
$$\sum m_k u_k v_k + \text{truc} = \sum m_k f_k v_k, \text{ i.e., } [u_h, v_h]_h + \text{truc} = [f_h, v_h]_h.$$

Nice! What have we done with "truc" in the continuous framework? we integrated by parts! There is an analogous procedure in the discrete setting: it is called "gathering by edges" (here, by half-edges σ).

Indeed, each ^{half}edge σ is included in some interface kL ; thus in "truc" it appears twice (the boundary edges are excluded from summation, they never appear). Moreover, while $\vec{a}(\vec{\nabla}_{T(\sigma)} u_h)$ only depends on T (i.e. it is the same for k and for L), the normals are reversed: $\vec{n}_L = -\vec{n}_k$. Therefore (with $\vec{F}_{T(\sigma)} = \vec{a}(\vec{\nabla}_{T(\sigma)} u_h)$)

$$\begin{aligned} \text{truc} &:= - \sum_k \left(\sum_{\sigma \subset \partial k \cap \partial \Omega} m_\sigma \vec{F}_{T(\sigma)} \cdot \vec{n}_k \right) v_k = - \sum_{\sigma} m_\sigma \vec{F}_{T(\sigma)} \cdot (\vec{v}_L - \vec{v}_k) \cdot \vec{n}_{kL} \\ &= + \sum_{\sigma} m_\sigma (\vec{F}_{T(\sigma)} \cdot \vec{n}_{kL}) (v_L - v_k) \text{ where } \sigma \subset kL. \end{aligned}$$


Here comes a wonderful thing... let us continue the calculation, gathering by T :

$$\begin{aligned} \text{truc} &= \sum_T \vec{F}_T \cdot \left\{ m_{\sigma_{kL}} (v_L - v_k) \vec{n}_{kL} + \dots + m_{\sigma_{hK}} (v_k - v_h) \vec{n}_{hK} \right\} \\ &= \sum_T \vec{F}_T \cdot \left\{ \underbrace{\frac{1}{2} m_{\sigma_{kL}} \frac{d_{kL}}{d_{kL}}}_{m_{T_{kL}}} \frac{v_L - v_k}{d_{kL}} \vec{n}_{kL} + \dots + \underbrace{\frac{1}{2} m_{\sigma_{hK}} \frac{d_{hK}}{d_{hK}}}_{m_{T_{hK}}} \frac{v_k - v_h}{d_{hK}} \vec{n}_{hK} \right\} \end{aligned}$$


with the notation on $T \rightarrow$

! $= \sum_T \vec{F}_T \cdot \vec{\nabla}_T v_h$ miracle! Pq Thanks to the orthogonality of the mesh, we have $\vec{n}_{kL} = \vec{n}_{hK}$.

We have found that $\text{truc} = \int \vec{a}(\vec{\nabla}_h u_h), \vec{\nabla}_h v_h$. We have shown one implication in the following proposition:

Proposition WF The scheme is equivalent to its "weak" formulation: find $u_h = (u_k) \in \mathbb{R}^{\#k}$ such that

$$\forall v_h = (v_k) \in \mathbb{R}^{\#k} \quad [u_h, v_h]_h + \int \vec{a}(\vec{\nabla}_h u_h), \vec{\nabla}_h v_h = [f_h, v_h]_h. \quad (\text{WDEq})$$

(The proof of the implication " \Leftarrow " is obtained by making the same calculation in the inverse sense, and taking $v_{k_0} = 1, v_L = 0$ for all $L \neq k_0$).

The miracle of the above calculation should be examined more closely (actually, we will need the same calculation, with different \vec{F}_h , for the convergence analysis). Let's formalize the calculation in terms of discrete operators.

Def. 00. $\nabla_h: \mathbb{R}^{\#k} \rightarrow (\mathbb{R}^2)^{\#T}$
 $u_h = (u_k)_k \mapsto (\vec{\nabla}_T u_h)_T$ with $\vec{\nabla}_T u_h = \frac{2}{m_T} \left\{ m_{T_{kL}} \frac{u_L - u_k}{d_{kL}} \vec{d}_{kL} + \dots \right\}$
 $\text{div}_h: (\mathbb{R}^2)^{\#T} \rightarrow \mathbb{R}^{\#k}$
 $\vec{F}_h = (\vec{F}_T)_T \mapsto (\text{div}_k \vec{F}_h)_k$ with $\text{div}_k \vec{F}_h = \frac{1}{m_k} \sum_{\sigma \subset \partial k \cap \partial \Omega} m'_\sigma \vec{F}_{T(\sigma)} \cdot \vec{n}_k$ Pq This is "div + zero flux BC"

Proposition 11 (The discrete duality property)

The operators $-div_h$ and ∇_h are dual for the products $\{\cdot, \cdot\}_h, \{\cdot, \cdot\}_h$, i.e.

$$\forall v_h \in \mathbb{R}^{\#k} \quad \left[-div_h \vec{F}_h, v_h \right]_h = \left\{ \vec{F}_h, \vec{\nabla}_h v_h \right\}_h$$

Proof: was given in the "true" calculation.

V) Properties of the scheme, straight forward from (WDEq):

• uniqueness: as in the continuous setting, take $v_h = u_h - \hat{u}_h$ with u_h, \hat{u}_h two solutions.

Then (WDEq) yields $\left[u_h - \hat{u}_h, u_h - \hat{u}_h \right]_h + \left\{ \vec{a}(\vec{\nabla}_h u_h) - \vec{a}(\vec{\nabla}_h \hat{u}_h), \vec{\nabla}_h u_h - \vec{\nabla}_h \hat{u}_h \right\}_h = 0$.

Using lifting operators $\int_{\Omega} (u_h - \hat{u}_h)^2(x) dx \stackrel{\parallel}{=} \int_{\Omega} (\vec{a}(\vec{F}) - \vec{a}(\vec{\hat{F}})) \cdot (\vec{F} - \vec{\hat{F}})(x) dx \geq 0$

Hence $u_h(\cdot) \equiv \hat{u}_h(\cdot)$ a.e, which means $u_k = \hat{u}_k$ for all k , i.e. $u_h = \hat{u}_h$.

• Δ^q estimates, $1 \leq q \leq \infty$.

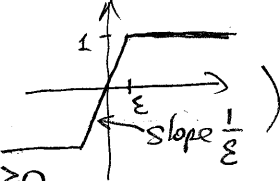
- $q=1$: Take $v_k := \text{sign } u_k$ (component per component, i.e., $v_k := \text{sign } u_k \forall k$)

NB No regularization is needed! At this point, discrete setting is a bit simpler than the continuous setting.

Then $\left[u_h, \text{sign } u_h \right]_h + \left\{ \vec{a}(\vec{\nabla}_h u_h), \vec{\nabla}_h(\text{sign } u_h) \right\}_h \leq \left[f_h, \text{sign } u_h \right]_h$

$$\int_{\Omega} u_h(x) \text{sign } u_h(x) dx = \int_{\Omega} |u_h(x)| dx \stackrel{\geq 0}{=} \int_{\Omega} f_h(x) \text{sign } u_h(x) dx \leq \int_{\Omega} |f_h(x)| dx = \left[|f_h|, 1 \right]_h$$

NB In general, let S be a Lipschitz non-decreasing function (e.g. a regularization of $z \mapsto \text{sign } z$, e.g. He:



Then in the continuous setting, $\vec{a}(\vec{\nabla} u) \cdot \vec{\nabla} S(u) \stackrel{\text{a.e.}}{=} S'(u) \vec{a}(\vec{\nabla} u) \cdot \vec{\nabla} u \geq 0$.

But in the discrete setting, it is easy to see that the chain rule is lost: chain rule for Lipow'IP functions!

$$\vec{\nabla}_T S(u_h) = \frac{2}{m_T} \left(m_{T_{kl}} \frac{S(u_k) - S(u_l)}{d_{kl}} \vec{\nu}_{kl} + \dots \right) = \frac{2}{m_T} \left(S'(\xi_{kl}) m_{T_{kl}} \frac{u_k - u_l}{d_{kl}} \vec{\nu}_{kl} + \dots + S'(\xi_{mk}) m_{T_{mk}} \frac{u_m - u_k}{d_{mk}} \vec{\nu}_{mk} \right)$$

The values $S'(\xi_{kl}), S'(\xi_{lm}), S'(\xi_{mk})$ are, in general, all different, and we cannot write $\vec{\nabla}_T S(u_h) = \Theta \vec{\nabla}_T u_h$; in particular, the orientation of $\vec{\nabla}_T S(u_h)$ is not necessarily the one of $\vec{\nabla}_T u_h$ and it is delicate (and sometimes false!) to conclude that $\vec{a}(\vec{\nabla}_h u_h) \cdot \vec{\nabla}_h S(u_h) \geq 0$. At this point, discrete setting is much more difficult than the continuous one.

We make two additional assumptions:

- Restrict. Str
- structure of \vec{a} : $\vec{a}(\vec{\xi}) = k(|\vec{\xi}|) \vec{\xi}$, $k \geq 0$ (verified by the p-laplacian)
 - structure of the triangulation T : all triangles have angles $\leq 90^\circ$.
- NB In practice, there is no algorithm for creating such partitions of general Ω !

Lemma Under the restrictions Str, for all non-decreasing (even discontinuous) $S(\cdot)$

[we have for all u_h , for all T $\vec{a}(\vec{\nabla}_T u_h) \cdot \vec{\nabla}_T S(u_h) \geq 0$.

Proof Let u_K, u_L, u_M be the values at the vertices x_K, x_L, x_M of T .

Assume u_M is the greatest value. Then also $S(u_M)$ is the greatest among $S(u_K), S(u_L), S(u_M)$.

Take $x_{KL} \in [x_K, x_L]$ such that $x_{KL}x_M$ is aligned with \vec{p}_{KL}^\perp (this is possible because the angles of T at x_K and at x_L are acute)

Because $\vec{\nabla}_T u_h$ is the gradient of the affine function interpolating the vertices, we have in particular $\vec{\nabla}_T u_h \cdot \vec{p}_{KL}^\perp = \frac{u_M - (\alpha u_K + (1-\alpha)u_L)}{|x_M - x_{KL}|}$

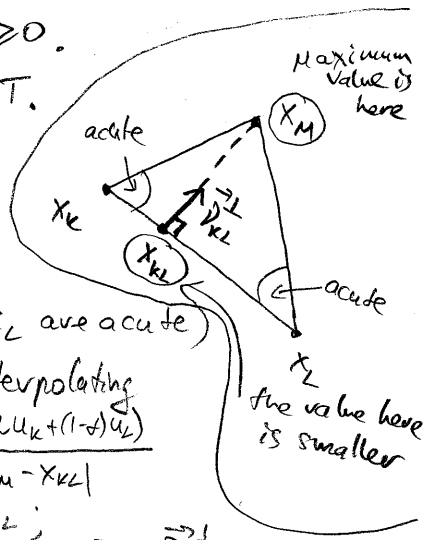
where $\alpha u_K + (1-\alpha)u_L$ is the value of the interpolation at x_{KL} ; it is crucial that $\alpha \in [0, 1]$, thus $u_M \geq \alpha u_K + (1-\alpha)u_L$ and $\vec{\nabla}_T u_h \cdot \vec{p}_{KL}^\perp \geq 0$.

In the same way, we see that $\vec{\nabla}_T S(u_h) \cdot \vec{p}_{KL}^\perp \geq 0$. Thus each of the vectors $\vec{\nabla}_T u_h$ and $\vec{\nabla}_T S(u_h)$ forms an acute angle with \vec{p}_{KL}^\perp ; and $\vec{a}(\vec{\nabla}_T u_h)$ is aligned with $\vec{\nabla}_T u_h$ by restrictions Str.

Now we write both vectors in the basis $(\vec{p}_{KL}, \vec{p}_{KL}^\perp)$; we have

$$\vec{a}(\vec{\nabla}_T u_h) \cdot \vec{\nabla}_T S(u_h) = k(|\vec{\nabla}_T u_h|) \left\{ \underbrace{\frac{u_L - u_K}{d_{KL}} \cdot \frac{S(u_L) - S(u_K)}{d_{KL}}}_{\geq 0} + \underbrace{(\vec{\nabla}_T u_h \cdot \vec{p}_{KL}^\perp)}_{\geq 0} \underbrace{(\vec{\nabla}_T S(u_h) \cdot \vec{p}_{KL}^\perp)}_{\geq 0} \right\} \geq 0$$

because $(t-s)(S(t)-S(s)) \geq 0$.
 the components of the projection on \vec{p}_{KL}^\perp are known! as we have just shown



Remark What we have not seen in this proof, is that the orthogonality of the mesh is very important.

Conclusion: we do have $\| |u_h|, 1 \|_h \leq \| |f_h|, 1 \|_h$

- $q \in (1, \infty)$: same reasoning, but use $S(z) = \begin{cases} |z|^{q-2} z, & z \neq 0 \\ 0, & z = 0 \end{cases}$ we get $\| |u_h|^q, 1 \|_h \leq \| |f_h|^q, 1 \|_h$

- $q = \infty$: see the previous relation as $(\int_\Omega |u_h|^q(x))^{1/q} \leq (\int_\Omega |f_h|^q(x))^{1/q}$, and pass to the limit $q \rightarrow \infty$, we get $\| u_h(\cdot) \|_\infty \leq \| f_h(\cdot) \|_\infty$, thus $\max_k |u_k| \leq \max_k |f_k|$.

• maximum principle:
 almost the same proof; e.g. if $M = \max_k f_k$, subtract M to all the values of u_h, f_h and use the test function $S(u_h) = \text{sign}^+(u_h - M)$; we get

$$\| u_h - M, \text{sign}^+(u_h - M) \|_h + (\text{term} \geq 0) = \| \underbrace{f_h - M}_{\leq 0}, \underbrace{\text{sign}^+(u_h - M)}_{\geq 0} \|_h,$$

thus $\int_\Omega (u_h - M)^+(x) dx \leq 0$, i.e. $u_h(x) \leq M$ a.e., i.e. $\forall k, u_k \leq M$, i.e. $\max_k u_k \leq \max_k f_k$.

• maximum principle bis:
 let us give another proof, using the standard technique for classical solutions of continuous problem.
NB At this point, discrete setting is much simpler than the continuous one.

Namely, assume $f_h \leq M$ and assume, by contradiction, that $u_k > M$ for some k . Then we can take the k for which $u_k - M$ attains its (positive) maximum.

Look at the equation for the volume K :

$$\underbrace{m_K u_K}_{\text{term 1}} - \sum_{\sigma \in \mathcal{K}(K)} \underbrace{m_\sigma k(\vec{\nabla}_T u_h)}_{\text{term 2}} \underbrace{\frac{u_L - u_K}{d_{KL}}}_{\text{term 3}} = \underbrace{m_K f_K}_{\text{term 3}}$$

Term 1 is $> m_k M$, term 2 is ≥ 0 because $-(u_k - \hat{u}_k) \geq 0$ and $k(\cdot) \geq 0$.
 and term 3 is $\leq m_k M$. We get $m_k M < \text{term 1} + \text{term 2} = \text{term 3} \leq m_k M$, contradiction.

• Can we prove the comparison principle in the same way?

Let's try! Assume $f_k \leq \hat{f}_k$ for all k , and (by contradiction) take the volume K where $u_k - \hat{u}_k$ attains its positive maximum. The corresponding discrete equation yields:

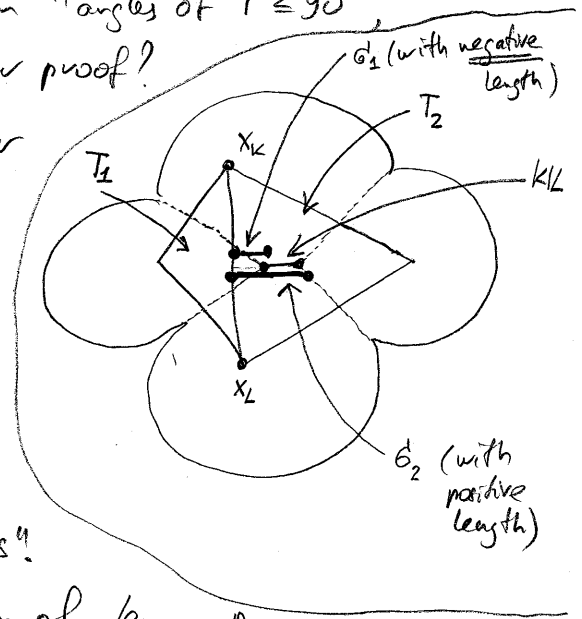
$$\underbrace{m_k (u_k - \hat{u}_k)}_{\geq 0} - \sum_{\delta \in \partial \text{dev} K} m_\delta \left\{ k(|\vec{T}_\delta| u_k) \frac{u_k - u_k}{d_{kL}} - k(|\vec{T}_\delta| \hat{u}_k) \frac{\hat{u}_k - \hat{u}_k}{d_{kL}} \right\} = m_k (f_k - \hat{f}_k) \leq 0$$

is it ≥ 0 ??

In the linear case ($k \equiv \text{const}$) we get $-(u_k - \hat{u}_k) - (u_k - \hat{u}_k) \geq 0$ by the choice of K .
 Looks ok? But where have we used the restriction "angles of $T \leq 90^\circ$ " that seemed to be useful for another proof?

Remark In the above scheme, one should consider "negative half-edges δ " except for the case "angles of $T \leq 90^\circ$ ". Example:

here we have not
 $k_{1L} = \delta_1 \cup \delta_2$, but
 $k_{1L} = \delta_2 \setminus \delta_1$.



Thus we assign m_δ as usual, but we take for m_δ the length of δ_1 with the sign "minus".

NB This convention is compatible with the convention of Lemma Rec (there we have taken the signed triangle areas m_{TKL} , etc).
 Proposition DD is true if we stick to the sign convention for both δ 's and T_k 's.

Conclusion: the above proof of comparison principle for $k \equiv \text{const}$ works under the same restriction str on the triangulation.

What's for general $k(\cdot)$?

Guess Not true in general.



True if angles of T are all $\leq 90^\circ - \Theta$ where $\Theta = \Theta(p)$, with $\Theta(2) = 0$ and $\lim_{p \rightarrow 1} \Theta(p) = 90^\circ = \lim_{p \rightarrow \infty} \Theta(p)$.

• Should we look at the contraction principle?

No! It is enough to use the celebrated Crandall-Tartar lemma.

Lemma CT Assume that $C \subset L^1(\Omega)$ verifies $[f, g \in C \Rightarrow \max\{f, g\} \in C]$

[Proc. AMS '80]

and $T: C \rightarrow L^1(\Omega)$ satisfies $\int_\Omega T(f) \leq \int_\Omega f$.

Then (for $f, g \in C$) the following properties are equivalent

- $f \leq g \Rightarrow T(f) \leq T(g)$ (order preservation)
- $\int_\Omega |T(f) - T(g)| \leq \int_\Omega |f - g|$ (L^1 contraction)
- $\int_\Omega (T(f) - T(g))^+ \leq \int_\Omega (f - g)^+$ (L^1 T-contraction)

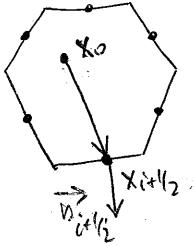
Application: $C = \{ \sum_k u_k \mathbb{1}_{K_k}(\cdot), u_k \in \mathbb{R} \}$; $T: f_k \mapsto u_k$ the solution of (DE_q) (existence first).
 we have $\int_\Omega T(f_k) = \int_\Omega u_k = \int u_k \mathbb{1}_{K_k} = \int u_k \mathbb{1}_{K_k} + \int 0 \mathbb{1}_{\Omega \setminus K_k}$, $\vec{v}_k \cdot \vec{\beta}_k = \beta_k^+ \mathbb{1}_{K_k} - \beta_k^- \mathbb{1}_{\Omega \setminus K_k}$ (use (WDE_q) with $v_k = 1$).

Thus whenever we have the comparison principle, we get contraction on Δ .
 Yet in general, both are false. 111

VI) Yet many other important properties of the continuous problem are preserved at the discrete level, thanks to the discrete duality property.

So, let's generalize and prove Lemma Rec 1.

Lemma Rec 2 let T be a polygon, $x_0 \in \mathbb{R}^2$. Let $(x_i)_{i=1}^l$ be the vertices (we set $x_{l+1} = x_1$) and let $x_{i,i+1/2} := \frac{x_i + x_{i+1}}{2}$ (the barycenter of the edge). let $d_{i,i+1} = |x_i - x_{i+1}|$ (unbound counter-clockwise)



Then for all $\vec{r} \in \mathbb{R}^2$, $\vec{r} = \frac{1}{m_T} \sum_{i=1}^l d_{i,i+1} (\vec{r} \cdot \overrightarrow{x_0 x_{i,i+1/2}}) \vec{n}_{i,i+1/2}$

Proof Consider $u: \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $u(x) = \vec{r} \cdot \overrightarrow{x_0 x}$; we have $\nabla u = \vec{r}$.

Using the Green-Gauss formula, $m_T \vec{r} = \int_T \nabla_2 u = \int_T \text{div}(\vec{u}, \vec{0}) = \int_{\partial T} (\vec{u}, \vec{0}) \cdot \vec{n} = \int_{\partial T} u n_1$.

In the same way, $m_T \vec{r} = \int_{\partial T} u n_2$, thus $\vec{r} = \frac{1}{m_T} \sum_{i=1}^l d_{i,i+1} \int_{[x_i, x_{i+1}]} (\vec{r} \cdot \overrightarrow{x_0 x}) \vec{n}_{i,i+1/2}$

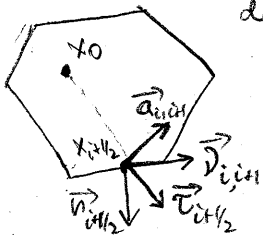
Because u is affine and $x_{i,i+1/2}$ is the barycenter of $[x_i, x_{i+1}]$, the result follows. □

Lemma Rec 3 With the notation of Lemma Rec 2,

denote in addition $\vec{v}_{i,i+1/2}$ = unit vector from x_0 to $x_{i,i+1/2}$

$\vec{v}_{i,i+1}$ = unit vector from x_i to x_{i+1}

$\vec{a}_{i,i+1}$ = unit vector orthogonal to $\vec{v}_{i,i+1}$, oriented in the direct sense.



Then for all $\vec{r} \in \mathbb{R}^2$, $\vec{r} = \frac{2}{m} \sum_{i=1}^l \frac{m_{i,i+1}}{\cos d_{i,i+1}} (\vec{r} \cdot \vec{v}_{i,i+1}) \vec{a}_{i,i+1}$

with $m_{i,i+1}$ = (signed) area of the triangle $x_0 x_i x_{i+1}$

$d_{i,i+1}$ = angle between $\vec{n}_{i,i+1/2}$ and $\vec{v}_{i,i+1}$, $\cos d_{i,i+1} = (\vec{n}_{i,i+1/2}, \vec{v}_{i,i+1}) = (\vec{v}_{i,i+1}, \vec{a}_{i,i+1})$

Proof: Several lines starting from Lemma Rec 2

and the identity $\vec{r} = \frac{1}{\cos d_{i,i+1}} \left\{ (\vec{r} \cdot \vec{v}_{i,i+1/2}) \vec{n}_{i,i+1/2} + (\vec{r} \cdot \vec{v}_{i,i+1}) \vec{a}_{i,i+1} \right\}$

(the identity is justified taking scalar product with $\vec{v}_{i,i+1/2}$, $\vec{v}_{i,i+1}$ and using orthogonality. □)

NB Lemma Rec 1 is a particular case of Lemma Rec 3 (corresponding to $T = \text{triangle}$, $x_0 = \text{circumcenter}$)

Generalization of the scheme \mathcal{T}

Take an arbitrary partition \mathcal{V} of Ω into polygons,

with one point per polygon: then create the

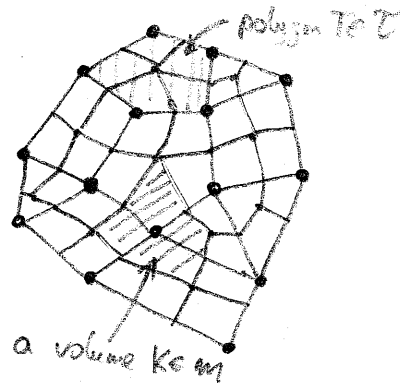
"median dual mesh" (= "Donald mesh") \mathcal{M}

as shown on the picture,

Define div_h on the mesh \mathcal{M} as previously,

Define ∇_h piecewise constant per $T \in \mathcal{T}$,

using the formula of Lemma Rec 3 ($\vec{r} \cdot \vec{v}_{i,i+1}$ substituted by $\frac{u_{i+1} - u_i}{d_{i,i+1}}$).



Proposition DD bis The discrete duality property holds for the so defined scheme.

III) Now, let's embark on analysis of the equations (DEq) (more exactly under the form (WDEq))

A priori estimates Remark: (WDEq) $[[u_h, v_h]]_h + \{ \vec{a}(\vec{\nabla}_h u_h), \vec{\nabla}_h v_h \}_h = [[f_h, u_h]]_h$ for all v_h

Take $v_h = u_h$; we get $\int_{\Omega} u_h^2(x) dx + \int_{\Omega} \vec{a}(\vec{\nabla}_h u_h)(x) \cdot (\vec{\nabla}_h u_h)(x) dx = [[u_h, u_h]]_h + \{ \}$, $[[u_h, u_h]]_h = \int_{\Omega} f_h(x) u_h(x) dx \leq \frac{1}{2} \int_{\Omega} f_h^2(x) dx + \frac{1}{2} \int_{\Omega} u_h^2(x) dx$

thus $\bullet \|u_h(\cdot)\|_{L^2} \leq \|f_h(\cdot)\|_{L^2}$
 \bullet (by coercivity of $\vec{a}(\cdot)$) $\|(\vec{\nabla}_h u_h)(\cdot)\|_{L^p} \leq C(\|f_h(\cdot)\|_{L^2})$ (with tacit convention that the dependence is non-decreasing!)

Moreover, if $f_k := \int_k f(x) dx$, then from the Jensen inequality, $f_k^2 \leq \int_k f^2(x) dx$;

then $\int_{\Omega} f_h^2(x) dx = \sum m_k f_k^2 \leq \sum m_k \int_k f^2(x) dx = \int_{\Omega} f^2(x) dx$,

thus $\|f_h(\cdot)\|_{L^2} \leq \|f(\cdot)\|_{L^2}$, for all mesh.

We have shown Proposition AEst Let f be given. Then for any mesh, any solution of (DEq)

verifies $\|u_h(\cdot)\|_{L^2} \leq C$, with some constant C that only depends on f, Ω, \vec{a} .
 $\|(\vec{\nabla}_h u_h)(\cdot)\|_{L^p} \leq C$

Variational interpretation

Here we assume $\vec{a} = \nabla \Phi$ (this is the case of the p-laplacian $\Phi(\xi) = \frac{1}{p} |\xi|^p$) with strictly convex coercive functional $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}$

Set $\mathcal{J}_h: \mathbb{R}^{*k} \rightarrow \mathbb{R}$
 $v_h \mapsto \frac{1}{2} \int_{\Omega} (u_h - f_h)^2(x) dx + \int_{\Omega} \Phi(\vec{\nabla}_h u_h)(x) dx = \frac{1}{2} \sum_k m_k (u_k - f_k)^2 + \sum_T m_T \Phi(\vec{\nabla}_T u_h)$

Proposition 4.6.11 (i) \mathcal{J}_h admits a unique minimizer on \mathbb{R}^{*k}

(ii) This minimiser is a solution of (WDEq) (the unique solution)

Proof (i) One easily checks that \mathcal{J}_h is a strictly convex continuous (in any norm of \mathbb{R}^{*k}) functional; it is coercive e.g. for $\|u_h\|_{\mathbb{R}^{*k}} := \sqrt{\sum_k m_k u_k^2} + \sqrt[2]{\sum_T m_T |\vec{\nabla}_T u_h|^2}$.

Thus it admits a unique minimizer

(ii) It is sufficient to write the Euler-Lagrange equation: $(d\mathcal{J}_h|_{u_h})[v_h] = 0 \Leftrightarrow$

$\Leftrightarrow \lim_{t \rightarrow 0} \frac{\mathcal{J}_h[u_h + t v_h] - \mathcal{J}_h[u_h]}{t} = \lim_{t \rightarrow 0} \left(\int_{\Omega} \frac{(u_h + t v_h - f_h)^2 - (u_h - f_h)^2}{2t} + \int_{\Omega} \frac{\Phi(\vec{\nabla}_h u_h + t \vec{\nabla}_h v_h) - \Phi(\vec{\nabla}_h u_h)}{t} \right)$
 $= \lim_{t \rightarrow 0} \left(\int_{\Omega} \frac{t v_h (2u_h - 2f_h + t v_h)}{2t} + \int_{\Omega} \frac{\vec{a}(\vec{\nabla}_h u_h) \cdot (t \vec{\nabla}_h v_h) + \bar{o}(t)}{t} \right)$
 $= \int_{\Omega} (u_h - f_h) v_h + \int_{\Omega} \vec{a}(\vec{\nabla}_h u_h) \cdot \vec{\nabla}_h v_h = [[u_h - f_h, v_h]]_h + \{ \vec{a}(\vec{\nabla}_h u_h), \vec{\nabla}_h v_h \}_h$

Thus (WDEq) is the Euler-Lagrange equation for the minimization of \mathcal{J}_h . □

Existence, general case

If \vec{a} does not derive from a scalar potential Φ , we have to use topological existence results (a version of the Brouwer theorem or the topological degree arguments) in order to get existence.

For the application of the Brouwer theorem: see [Lions '69, lemma 4.3].

Let us give the topological degree argument. Reference: J. Droniou (visit his personal homepage)

Idea: Let X be a finite-dimensional space.
Let $H: [0,1] \times X \rightarrow X$ be a continuous map such that

- $H(0, \cdot) = 0$ admits a solution
- for all solution x of $H(t, \cdot) = 0$, one has $\|x\|_X \leq R$

Then for all t , $H(t, \cdot) = 0$ has at least one solution, and in particular $H(1, \cdot) = 0$ admits a solution.

(tools: homotopy, notion of topological degree and its invariance per homotopy)

Application: $X = \mathbb{R}^{*K}$ (any norm can be taken on X)
 $H: ([0,1] \times \mathbb{R}^{*K}) \rightarrow \mathbb{R}^{*K}$
 $(t, u_h) \mapsto (m_k(u_k - t f_k) - \operatorname{div}_{\vec{a}}(\vec{f}_h u_h))_k$

- The continuity of H is evident ($H(t, (u_k)_k)$ is produced by a finite number of arithmetic operations and a composition by the continuous map $\vec{a}(\cdot)$)
- For $t=0$, $u_h=0$ is the unique solution
- Applying Proposition AEst with f_h replaced with $t f_h$, we get the uniform bound on possible solutions $u_h^{(t)}$ of $H(t, \cdot) = 0$ (\Leftrightarrow the scheme (DEq) with source $t f_h$), eg. in the norm $\|u_h\|_X := \sqrt{\|u_h, u_h\|_h}$.

Conclusion. There exists a unique solution to the scheme (DEq).

- In the variational case ($\vec{a} = \nabla \Phi$), it can be obtained by minimization of the discrete energy J_h .

In practice, this allows for the use of "descent" iteration methods: conjugate gradient, Polak-Ribière...

- In the non-variational and in the variational case, one could use the Newton method to solve the nonlinear system; but in practice, it is difficult to pick the initial guess for which the algorithm would converge. Newton works better for the evolution problem associated with $-\operatorname{div} \vec{a}(0, \cdot)$.

- In the non-variational and in the variational case, the "coordination-decomposition" algorithm of Glowinski and Maurocco (see the recent book of Glowinski) appears to be quite efficient.

We do not discuss more the efficiency of the solution of the nonlinear system (DEq).

Now, the questions are:

- in which sense the discrete solutions could converge to the continuous ones?
- how to prove such convergence?
- at what minimal rate the convergence takes place (a priori error estimates)?
- what rate of convergence do we observe in practice (numerical tests)?
- * can we use the observed numerical solutions to improve the convergence rate estimate (a posteriori error estimates)?
- * can we use the observed numerical solution in order to refine the mesh adaptively and thus improve the ratio approximation/computational effort?

* : will not be answered.

VIII

Because we have lifted u_h and $\nabla_h u_h$ to functions on Ω , it is natural to state convergence in the following terms:

$$u_h(\cdot) \rightarrow u(\cdot)$$

$$(\nabla_h u_h)(\cdot) \rightarrow \nabla u(\cdot),$$

but in which sense?

We have the following "asymptotic compactness" result

Proposition Cmp

T_h should be triangulations

Let T_h, M_h be a family of dual/primal meshes and u_h , the associated discrete functions satisfying ^(the estimates of) Prop. AEst. Then there exists a

sequence $(h_m)_{m \in \mathbb{N}}$ [in the sequel we drop the subscript m] and a function $u \in W^{1,p}(\Omega) \cap L^2(\Omega)$ such that

- $u_h(\cdot) \rightarrow u(\cdot)$ in $L^2(\Omega)$ (weakly)
- $(\nabla_h u_h)(\cdot) \rightarrow \nabla u(\cdot)$ in $L^p(\Omega)$ (weakly)
- moreover, $u_h \rightarrow u(\cdot)$ in $L^1(\Omega)$ (strongly).

Remark This is the simplest statement; using discrete Poincaré embeddings (that may require some uniform regularity of the meshes) one can improve the last statement to a strong convergence in $L^q(\Omega)$ for all $q < p^*$ (the critical Sobolev exponent)

The proof is split into several steps

- Basic weak compactness: the a priori estimates of Proposition AEst yield: there exists $u \in L^2(\Omega)$ such that (for a subsequence) $u_h(\cdot) \rightarrow u(\cdot)$ in $L^2(\Omega)$ and $\bar{w} \in L^2(\Omega; \mathbb{R}^d)$ such that $(\nabla_h u_h)(\cdot) \rightarrow \bar{w}(\cdot)$ in $L^p(\Omega)$ (both convergences are weak)

Estimate of L^1 translates:

- It is enough to show the $L^1_{loc}(\Omega)$ compactness. Indeed, if $(\omega_n)_n$ is a sequence of compacts such that $\omega_n \uparrow \Omega$, if $(u_h)_n$ is relatively compact in $L^1(\omega_n)$ for all n , then using the diagonal procedure we get a subsequence $(u_h)_n$ and a $u \in L^1_{loc}(\Omega)$ such that $u_h \rightarrow u$ in $L^1(\omega_n)$ for all n . The limit $u(\cdot)$ is the same as in the previous step, because weak $L^2(\Omega)$ convergence and strong $L^1(\omega_n)$ convergence both imply $\mathcal{D}'(\omega_n)$ convergence, and the \mathcal{D}' limit is unique [we work with some fixed convergent extracted subsequences]

Moreover, $\|u_h(\cdot) - u(\cdot)\|_{L^1(\Omega)} \leq \|u_h(\cdot) - u(\cdot)\|_{L^1(\omega_n)} + \sqrt{|\Omega \setminus \omega_n|} \|u_h(\cdot) - u(\cdot)\|_{L^2(\Omega)} \rightarrow 0 \leq 2 \cdot \text{const by Prop. AEst.}$

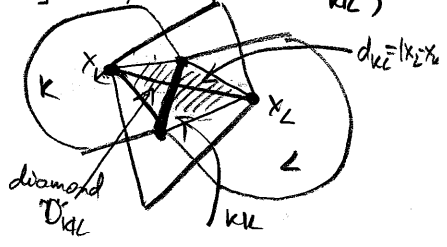
Fix $\omega \Subset \Omega$, take $\Delta \in \mathbb{R}^2$ such that $|\Delta| < \text{dist}(\omega, \partial\Omega)$. For all $x \in \Omega$ and all interface $K|L$, define $\Theta_{K|L}(x) := \begin{cases} 1, & \text{if } [x, x+\Delta] \cap (K|L) \neq \emptyset \\ 0, & \text{otherwise} \end{cases}$

Consider $I(\Delta) = \int_{\omega} |u_h(x+\Delta) - u_h(x)| dx$; we notice that

$$|u_h(x+\Delta) - u_h(x)| \leq \sum_{K|L: [x, x+\Delta] \cap (K|L) \neq \emptyset} |u_L - u_K| = \sum_{K|L} \Theta_{K|L}(x) |u_L - u_K|$$

Then $I(\Delta) \leq \sum_{K|L} \left(\int_{\omega} \Theta_{K|L}(x) dx \right) |u_L - u_K|$. But $\text{meas} \{x \mid [x, x+\Delta] \cap (K|L) \neq \emptyset\} \leq \Delta \cdot m_{K|L}$,

thus $I(\Delta) \leq \Delta \sum_{K|L} m_{K|L} d_{K|L} \frac{|u_L - u_K|}{d_{K|L}}$. Further, introducing "diamonds" $D_{K|L}$ as in the picture, we have $I(\Delta) \leq 2\Delta \sum_{K|L} m_{D_{K|L}} \frac{|u_L - u_K|}{d_{K|L}} \leq 2\Delta \int_{\Omega} |\nabla_h u_h(x)| dx \leq C\Delta \leq C \|\nabla_h u_h\|_{L^p(\Omega)}$



Therefore the space translates of the functions $u_h(\cdot)$ on ω are uniformly bounded; by the Kolmogorov compactness theorem, one can extract an $L^1(\omega)$ -convergent subsequence.

Question Where have we used the assumption that T are triangles??

Answer Only in this case, the divided difference $\frac{|u_L - u_R|}{d_{KL}}$ is $\leq |\nabla_T u_h|$ where T is an element of \mathcal{T} with vertices x_L, x_R .

Rq Actually we have shown that the translates of $u_h(\cdot)$ are estimated by $\text{const} \cdot \Delta$, which means (see e.g. [Brezis]) that $u \in BV(\omega)$; moreover, $u_h(\cdot) \rightarrow u(\cdot)$ on $BV(\omega)$. For this result, an L^1 estimate of $\nabla_h u_h(\cdot)$ would be enough.

• Consistency for test functions It remains to show that $\vec{w} = \vec{\partial} u$, on the $W^1 P(\Omega)$ case. This will be done on indirect way (through a weak formulation) First we need a first consistency result.

Def. Proj For a scalar or vector function on $L^1(\Omega)$, one defines the projection operators

• $P_h : f \mapsto P_h f = (\int_K f(x) dx)_K \in \mathbb{R}^{\#K}$
 • $\vec{P}_h : \vec{\Psi} \mapsto \vec{P}_h \vec{\Psi} = (\int_T \vec{\Psi}(x) dx)_T \in (\mathbb{R}^2)^{\#T}$

Proposition Consist 1 Let $\vec{\Psi} \in C^1_0(\Omega)$. Assume that all $T \in \mathcal{T}_h$ are triangles, and moreover, $\cos d_{i,j} \geq \delta > 0$ uniformly (see lemma Rec 3).

Then for all $v \in \mathbb{R}^{\#K}$,

$$\left| \left[P_h(\text{div } \vec{\Psi}) - \text{div}_h(\vec{P}_h \vec{\Psi}), v \right]_h \right| \leq C(\Psi) h \|\vec{\nabla}_h v\|_{L^1(\Omega)},$$

where $h \approx \max_{T \in \mathcal{T}_h} (\text{diam } T)$ (the parameter h has the sense of the size of the mesh).

Rq In the framework of volume schemes, the operators do not commute and moreover, the consistency property should be stated in the weak sense above.

Proof (sketched)

We introduce for all half-edge σ , the value $\vec{\Psi}_\sigma = \int_\sigma \vec{\Psi}(x) dx$ (in addition to the values $\vec{\Psi}_T = \int_T \vec{\Psi}(x) dx$ taken by $\vec{P}_h \vec{\Psi}$).

We can write $\text{div}_K(\vec{P}_h \vec{\Psi}) = \frac{1}{m_K} \sum_{\sigma \in \mathcal{E}_K} m_\sigma \vec{\Psi}_{T(\sigma)} \cdot \vec{n}_K$
 and $P_K(\text{div } \vec{\Psi}) = \frac{1}{m_K} \sum_{\sigma \in \mathcal{E}_K} m_\sigma \vec{\Psi}_\sigma \cdot \vec{n}_K$.

Gathering by triangles T and proceeding as in the "discrete duality calculation", we find the expression for $\left[P_h \text{div } \vec{\Psi} - \text{div}_h(\vec{P}_h \vec{\Psi}), v \right]_h$:

$$\sum_T m_T \frac{2}{m_T} \left(\frac{m_{TKL}}{\cos d_{KL}} \frac{V_L - V_R}{d_{KL}} (\vec{\Psi}_{d_{KL}} - \vec{\Psi}_T) + \dots \right)$$

But $|\vec{\Psi}_{d_{KL}} - \vec{\Psi}_T| \leq C \|\vec{\nabla} \Psi\|_\infty h$ and $\left| \frac{V_L - V_R}{d_{KL}} \right| \leq |\vec{\nabla}_T v|$ (because T is a triangle).

Using the lower bound on $\cos d_{i,j}$, we end up with the bound $\sum_T m_T C h |\vec{\nabla}_T v|$, which ends the proof. \square

Proposition Const 2 Let $\vec{\Psi} \in C(\Omega)$. Then $\|(\mathbb{P}_h \vec{\Psi})(\cdot) - \vec{\Psi}(\cdot)\|_{L^\infty(\Omega)} \leq Ch$.

Proof: Exercise.

Use of discrete duality Now we are ready to conclude. Let $\vec{\Psi} \in C^\infty(\Omega; \mathbb{R}^2)$. Set $\vec{\Psi}_h = \mathbb{P}_h \vec{\Psi}$; then by Proposition DD,

$$\| -\text{div}_h \vec{\Psi}_h, u_h \|_h = \{ \vec{\Psi}_h, \vec{\nabla}_h u_h \}_h$$

Moreover, by the commutation property of Prop. Const 1,

$$\| -\mathbb{P}_h(\text{div} \vec{\Psi}), u_h \|_h = \{ \vec{\Psi}_h, \vec{\nabla}_h u_h \}_h + \bar{o}(1) \quad (\text{because } \|\vec{\nabla}_h u_h\|_{L^2(\Omega)} \leq \text{const}!)$$

Rewriting this under the form of integrals, we get

$$= \sum_k m_k \frac{1}{m_k} \left(\int_k (\text{div} \vec{\Psi})(x) dx \right) u_k = \sum_T m_T \frac{1}{m_T} \left(\int_T \vec{\Psi}(x) dx \right) \cdot \vec{\nabla}_T u_h + \bar{o}(1) \quad \text{thus}$$

$$= \int_\Omega u_h(x) (\text{div} \vec{\Psi})(x) dx = \int_\Omega \vec{\Psi}(x) \cdot (\vec{\nabla}_h u_h)(x) dx + \bar{o}(1)$$

By the definition of $u(\cdot)$ and $\vec{w}(\cdot)$, we get at the limit

$$- \int_\Omega u \text{div} \vec{\Psi} = \int_\Omega \vec{\Psi} \cdot \vec{w}$$

for all $\vec{\Psi} \in C^\infty(\Omega)$. Recall that $u \in L^2(\Omega)$, $\vec{w} \in L^p(\Omega)$. This means exactly that $u \in W^{1,p}(\Omega) \cap L^2(\Omega)$ and $\vec{w} = \vec{\nabla} u$ in the same $W^{1,p}(\Omega)$. \square

IX - Now we want to prove that (for any subsequence of $(u_h)_h$ convergent in the sense of Proposition Cmp, the limit u is a solution to (MP_1) . This gives existence! Since uniqueness for (MP_1) is known, this implies that the whole sequence $(u_h)_h$ converges to the unique solution of (MP_1) [argument: by contradiction, extract a subsequence that remains ε -far from u ; by Proposition Cmp, it admits a subsequence convergent and if we show that the limit solves (MP_1) , we know that the limit is u ; this is absurd].

Because our notion of solution is weak sol., we try to pass to the limit onto the weak formulation (WEqD) of the discrete system (EqD).

Theorem V Let $\mathcal{T}_h, \mathcal{M}_h$ be a family of meshes: \mathcal{T}_h a triangulation of maximal size h , \mathcal{M}_h the associated mesh of the co-volume method,

Assume the centers x_T of triangles $T \in \mathcal{T}_h$ are chosen so that $\text{cost}_k \geq \delta > 0$ uniformly in h [this kind of uniform regularity assumption on meshes often appears].

Then for all h , there exists a unique solution u_h of the co-volume scheme (DEq) for problem (MP_1) . Moreover, the associated lifted functions $u_h(\cdot)$ and the associated lifted gradients $(\vec{\nabla}_h u_h)(\cdot)$ converge to $u, \vec{\nabla} u$ in the following sense:

- $u_h(\cdot) \rightarrow u(\cdot)$ in $L^2(\Omega)$ strongly [even in $L^q(\Omega)$ strongly for $q < p^*$; if we can use optimal discrete Sobolev embedding]
- $(\vec{\nabla}_h u_h)(\cdot) \rightarrow \vec{\nabla} u(\cdot)$ in $L^p(\Omega)$ weakly, and also strongly if \vec{a} is strictly monotone

here u is the unique solution of (MP_1) .

Remark The strong convergences are a byproduct of the proof: they do not follow from Prop. Cmp. In fact, Prop. Cmp only uses boundedness of u_h in discrete L^2 and $W^{1,p}$ spaces; but Theorem V uses in addition the fact that u_h (resp., u) satisfy the discrete (resp., continuous) "differential relations"; thus we have an effect of the "compensated compactness" kind.

Proof We take $\varphi \in C^1(\bar{\Omega})$ and use $\varphi_h = (\varphi(x_k))_k$ (not $P_h \varphi$!) as the test function in (WDEq). We have $[(u_h - f_h, \varphi_h)]_h + \{ \vec{a}(\vec{\nabla}_h u_h), \vec{\nabla}_h \varphi_h \}_h = 0$. We will write this as integrals...

As in Prop. Consist 2, we can say that $\varphi_h(\cdot) \rightarrow \varphi(\cdot)$ in $L^2(\Omega)$ strongly.
 We need three more convergences: $f_h(\cdot) \rightarrow f(\cdot)$ in $L^2(\Omega)$ weakly
 $\vec{\nabla}_h(P_h^C \varphi) \rightarrow \vec{\nabla} \varphi(\cdot)$ in $L^p(\Omega)$ strongly.

this part will be extremely delicate, because \vec{a} is nonlinear!
 $\rightarrow [\cdot \vec{a}(\vec{\nabla}_h u_h)(\cdot) \rightarrow \vec{a}(\vec{\nabla} u)(\cdot)]$ in $L^p(\Omega)$ weakly

Proposition Consist 3 Let $f \in L^2(\Omega)$ and $f_h = P_h f$. Then $f_h(\cdot) \xrightarrow{h \rightarrow 0} f(\cdot)$ in $L^2(\Omega)$ strongly.

Proof From the Jensen inequality, we know that $\|f_h(\cdot)\|_{L^2} \leq \|f(\cdot)\|_{L^2}$; in the abstract way, we can write that the linear operator $P_h \circ P_h : L^2(\Omega) \rightarrow L^2(\Omega)$ is of norm ≤ 1 .

Now, it is sufficient to approximate f by a sequence $(f^m)_m$ of $C(\Omega)$ functions.

By Prop. Consist 2, we find $\|f_h^m(\cdot) - f^m(\cdot)\|_{L^2} \rightarrow 0$ as $h \rightarrow 0$.

Then $\|f_h(\cdot) - f(\cdot)\|_{L^2} \leq \|f_h(\cdot) - f_h^m(\cdot)\|_{L^2} + \|f^m(\cdot) - f(\cdot)\|_{L^2} + \|f_h^m(\cdot) - f^m(\cdot)\|_{L^2} \leq 2 \underbrace{\bar{O}(h)}_{\text{uniform in } h!} + \underbrace{\bar{O}^m(h)}_{h \rightarrow 0}$

⚠ Idea: we used abstract tools to avoid doing estimates "by hands". The price to pay is that we ignore the convergence rate

$\|P_h \circ P_h(f - f^m)\|_{L^2} \leq \|f - f^m\|_{L^2}$

Taking $\varepsilon > 0$, then $m = m(\varepsilon)$, then $h = h(\varepsilon, m(\varepsilon))$, we conclude. \square

$(\varphi_h = P_h^C \varphi)$ a new projection

Proposition Consist 4 Let $\varphi \in C_0^2(\bar{\Omega})$ and $\varphi_h = (\varphi(x_k))_k$. Then $\|\vec{\nabla}_h \varphi_h - \vec{\nabla} \varphi(\cdot)\|_{L^\infty} \leq C(\varphi, \delta) h$ where δ is the mesh regularity constant: $\cos \delta_{kl} \geq \delta > 0, \frac{d_{kl}}{\text{diam}(T)} \geq \delta > 0$.

Proof We argue separately in each triangle T . Let w be the affine Taylor polynomial of φ at some fixed point of T . Then from the C^2 regularity of φ , because $\text{diam}(T) \leq h$, we get $\|\vec{\nabla} \varphi - \vec{\nabla} w\|_{L^\infty(T)} \leq C(\varphi) h$ and $|\varphi_k - w_k|, |\varphi_l - w_l|, |\varphi_m - w_m| \leq C(\varphi) (\text{diam}(T))^2$

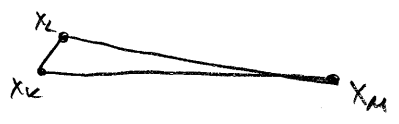
Now, from Lemma Rec 3 we know that $\vec{\nabla}_h$ is exact for affine functions, thus if we reconstruct $\vec{\nabla}_T w$ from w_k, w_l, w_m , we get $\vec{\nabla}_T w = \vec{\nabla} w$ on T .

It remains to estimate $\|\vec{\nabla}_T \varphi_h - \vec{\nabla}_T w\|_{L^\infty}$; looking at the reconstruction formula, $\vec{\nabla}_T \varphi_h - \vec{\nabla}_T w = \vec{\nabla}_T(\varphi_h - w)$, because $\vec{\nabla}_T$ is linear

$$\vec{\nabla}_T(\varphi_h - w) = \frac{2}{m_T} \left\{ \frac{m_{Tkl}}{\cos \delta_{kl}} \frac{(\varphi_k - w_k) - (\varphi_l - w_l)}{d_{kl}} + \dots \right\} \leq \frac{1}{\delta} \leq C(\varphi) \frac{(\text{diam}(T))^2}{d_{kl}} \leq C(\varphi) \frac{1}{\delta} \text{diam}(T) \leq C(\varphi, \delta) h$$

⚠ here it is important that we take $\varphi_k = \varphi(x_k)$ and not $\int_k \varphi(x) dx$, except if x_k is the barycenter of k

Remark Here we see that further proportionality assumptions on the meshes appear: the points x_k, x_l should not be too close with respect to the "typical length" in T . This only forbids "too flat" triangles of the kind (this restriction is much weaker than the "angle condition: $\leq 90^\circ$ " that was used for the maximum principle on Voronoi mesh)



Now, from the uniform bound of $\|\vec{\nabla}_h u_h(\cdot)\|_{L^p}$ and the growth assumption on $\vec{a}(\cdot)$, we get the uniform bound on $\|\vec{a}(\vec{\nabla}_h u_h(\cdot))\|_{L^{p_1}}$, and [extracting a further subsequence] we can say that $\vec{a}(\vec{\nabla}_h u_h(\cdot)) \rightarrow \vec{\chi}(\cdot)$ in $L^{p_1}(\Omega)$ as $h \rightarrow 0$.

Main question of the theory: why would we have $\vec{\chi} = \vec{a}(\vec{\nabla}u)$???

From the function-analytic arguments, nothing follows (or almost nothing): using the Young measures to describe weak convergence, we can establish some relation between $\vec{\chi}$ and ∇u , because both are calculated via the Young measure... we will see it later, for the $p(x)$ -Laplacian. Therefore the only hope is that the PDE itself brings additional information. Here is the information:

Lemma Anti-weak With the notation above, we have $\int \vec{\chi} \cdot \vec{\nabla}u \geq \liminf_{h \rightarrow 0} \int \vec{a}(\vec{\nabla}_h u_h) \cdot \vec{\nabla}_h u_h$

- UB. We have the inequality with \lim , but \liminf would be enough
- In the sequel we will show that the limit exists, and the equality holds (this is related to upgrading the weak convergence of $\vec{\nabla}_h u_h$ to the strong one)
- If $\vec{a} = \text{Id}$, the inequality yields $\|\vec{\nabla}u\|_{L^2} \geq \liminf_{h \rightarrow 0} \|(\vec{\nabla}_h u_h)\|_{L^2}$, which is known as a sufficient condition to upgrade the weak convergence " $\vec{\nabla}_h u_h \rightarrow \vec{\nabla}u$ in L^2 " into the strong L^2 convergence; in this sense, the inequality of the lemma "contradicts" weak (=non strong) convergence.

Proof The idea is: - take u for the test fct. in the limit of the scheme (involving $\vec{\chi}$)
- take u_h in (WDEq)
- pass to the limit where you can do so, and compare what remains.

So, we have proved so far that

$\forall \varphi \in C^\infty(\bar{\Omega}) \quad \int_{\Omega} u \varphi + \int_{\Omega} \vec{\chi} \cdot \vec{\nabla} \varphi = \int_{\Omega} f \varphi$; by the density of $C^\infty(\bar{\Omega})$ in $W^{1,p}(\Omega) \cap C(\bar{\Omega})$, we can take $\varphi = u$ for the test function;

then $\int_{\Omega} \vec{\chi} \cdot \vec{\nabla} u = \int_{\Omega} f u - \int_{\Omega} u^2$.

At the discrete level, we denote $\vec{\chi}_h := \vec{a}(\vec{\nabla}_h u_h)(\cdot)$ and from (WDEq) with $v_h = u_h$,

we get $\int_{\Omega} \vec{\chi}_h \cdot \vec{\nabla}_h u_h \, dx = \int_{\Omega} f_h(x) u_h(x) \, dx - \int_{\Omega} u_h^2(x) \, dx$.

We can "almost" pass to the limit in the right-hand side:

- $\int f u = \lim_{h \rightarrow 0} \int f_h u_h$ by the properties $f_h \rightarrow f$ in L^2 strongly (Prop. Consist 3)
- $\int u^2 \leq \liminf_{h \rightarrow 0} \int u_h^2$ by the weak convergence (or by the Fatou lemma); recall that we have shown $u_h \rightarrow u$ in $L^1(\Omega)$, which also gives "a.e." upon extraction of a further subseq.)

This is "fortunately" the good sign!

Comparing the two relations, we get the result. \square

Now, it's time to exploit the monotonicity of $\vec{a}(\cdot)$ which is the key property for identification of $\vec{\chi}$ to $\vec{a}(\vec{\nabla}u)$.

We have the following "miracle" lemma (known as "Minty trick") 19

Lemma Minty Assume $(v_n)_n \subset V$ and $(w_n)_n \subset V^*$ (V a Banach space, V^* its dual, $\langle \cdot, \cdot \rangle$ the duality pairing) such that

- $v_n \rightharpoonup v$ in V weakly
- $w_n \xrightarrow{*} w$ in V^* weakly-*
- $w_n = Tv_n$ where $T: V \rightarrow V^*$ is a monotone [nonlinear] operator

Then $w = Tv$ provided $\forall v, \hat{v} \quad \langle T(v) - T(\hat{v}), v - \hat{v} \rangle \geq 0$ and hemicontinuous: $T(v + \lambda \hat{v}) \xrightarrow{*} T(v)$ in V^* weakly as $\lambda \rightarrow 0, \forall v, \hat{v} \in V$

we have $\langle w, v \rangle \geq \lim_{n \rightarrow \infty} \langle w_n, v_n \rangle$.

Proof For $\lambda \in \mathbb{R}$ and $z \in V$, consider $v := v - \lambda z$. We have

$$\langle w, v - \varphi \rangle \geq \liminf \langle w_n, v_n - \varphi \rangle = \liminf \langle T(\varphi) + (T(v_n) - T(\varphi)), v_n - \varphi \rangle$$

$$\stackrel{\text{assumption}}{\geq} \liminf \langle T(\varphi), v_n - \varphi \rangle = \liminf \langle T(v - \lambda z), (v_n - v) + \lambda z \rangle = \lambda \langle T(v - \lambda z), z \rangle$$

Now, $\lambda \langle w, z \rangle$.
 • take $\lambda > 0$, divide by λ , send λ to zero $\Rightarrow \langle w, z \rangle \geq \lim_{\lambda \downarrow 0} \langle T(v - \lambda z), z \rangle \stackrel{\text{hemicontinuity of } T}{=} \langle T(v), z \rangle$
 • take $\lambda < 0$, divide by λ — change the sign " \geq " to " \leq "! — and get $\langle w, z \rangle \leq \langle T(v), z \rangle$.

Hence $\forall z \in V \quad \langle w, z \rangle = \langle T(v), z \rangle$, which yields $w = T(v)$. □

Announcement €100 to the person who convinces me that the proof is a natural one...

Application: $V = L^p(\Omega), V^* = L^q(\Omega), T: V \rightarrow V^*$ maps V into V^* thanks to the growth assumption

Hemicontinuity comes from the dominated convergence theorem + growth assumption (exercise).
 In conclusion, we establish $\vec{\nabla} \bar{x} = \vec{\nabla}(\vec{\nabla} u)$.

It remains to prove strong convergence claims.

Upgrading to strong convergence (if $\forall \vec{z}, \vec{y} \quad (\vec{\nabla}(\vec{z}) - \vec{\nabla}(\vec{y})) \cdot (\vec{z} - \vec{y}) > 0$ whenever $\vec{z} \neq \vec{y}$)

Coming back to the calculation of Lemma Anti-weak, we get

$$\left(\int u^2 - \int u_n^2 \right) + \left(\int \vec{\nabla}(\vec{\nabla} u) \cdot \vec{\nabla} u - \int \vec{\nabla}(\vec{\nabla}_n u_n) \cdot \vec{\nabla}_n u_n \right) \xrightarrow{u \rightarrow 0} 0$$

$$\text{Notice that } \vec{\nabla}(\vec{\nabla} u) \cdot \vec{\nabla} u - \vec{\nabla}(\vec{\nabla}_n u_n) \cdot \vec{\nabla}_n u_n = (\vec{\nabla}(\vec{\nabla} u) - \vec{\nabla}(\vec{\nabla}_n u_n)) \cdot (\vec{\nabla} u - \vec{\nabla}_n u_n)$$

$$\text{we have } \int \vec{\nabla}_n \cdot \vec{\nabla} u \rightarrow \int \vec{\nabla} \cdot \vec{\nabla} u \quad + \underbrace{\vec{\nabla}(\vec{\nabla}_n u_n) \cdot \vec{\nabla} u}_{\vec{\nabla}_n} - \underbrace{\vec{\nabla}(\vec{\nabla} u) \cdot \vec{\nabla}_n u_n}_{\vec{\nabla}}$$

$$\int \vec{\nabla}(\vec{\nabla} u) \cdot \vec{\nabla}_n u_n \rightarrow \int \vec{\nabla}(\vec{\nabla} u) \cdot \vec{\nabla} u = \int \vec{\nabla} \cdot \vec{\nabla} u$$

Hence $\lim_{n \rightarrow \infty} \left(\int u^2 - \int u_n^2 \right) + \lim_{n \rightarrow \infty} \left(\vec{\nabla}(\vec{\nabla} u) - \vec{\nabla}(\vec{\nabla}_n u_n) \right) \cdot (\vec{\nabla} u - \vec{\nabla}_n u_n) \geq 0$
 ≥ 0 , by the weak convergence ≥ 0 , by the monotonicity of $\vec{\nabla}(\cdot)$.

From the weak convergence of u_n to u on L^2 we know that $\|u\|_{L^2} \leq \liminf_{n \rightarrow \infty} \|u_n\|_{L^2}$, thus we deduce $\|u\|_{L^2} = \lim_{n \rightarrow \infty} \|u_n\|_{L^2}$ (and the limit exists). Uniform convexity \Rightarrow the cv. is strong.

We also have from the convergence to zero in $L^1(\Omega)$ of the nonnegative function

$$(\vec{a}(\vec{\sigma}u) - \vec{a}(\vec{\sigma}_n u_n)) \cdot (\vec{\sigma}u - \vec{\sigma}_n u_n) \quad (\text{up to a subsequence}).$$

Thus we also have the a.e. convergence to zero of this quantity.

Now we can reason pointwise: set $\vec{\xi} = \vec{\sigma}u(x)$ and $\vec{\eta}_n = (\vec{\sigma}_n u_n)(x)$.

If $\vec{\eta}_n \not\rightarrow \vec{\xi}$, we have $(\vec{a}(\vec{\xi}) - \vec{a}(\vec{\eta}_n)) \cdot (\vec{\xi} - \vec{\eta}_n) \not\rightarrow 0$ by the strong monotonicity of $\vec{a}(\cdot)$.

Thus $\vec{\sigma}_n u_n \rightarrow \vec{\sigma}u$ a.e. on Ω .

Then we have the following situation:

$$\text{setting } \begin{cases} F := \vec{a}(\vec{\sigma}u) \cdot \vec{\sigma}u \\ F_n := \vec{a}(\vec{\sigma}_n u_n) \cdot \vec{\sigma}_n u_n \end{cases}, \quad \left. \begin{cases} F_n \rightarrow F \text{ a.e. on } \Omega \\ F_n \geq 0 \end{cases} \right\} \text{Fatou lemma: } \int F \leq \liminf \int F_n$$

(and moreover, $\int F = \lim \int F_n$)

Lemma Sch (known as the Schaeffe lemma) (i.e. we are in the critical case of the Fatou lemma)

If on the assumptions of the Fatou lemma we have in addition $\int F = \lim \int F_n$, then $F \rightarrow F_n$ in L^1 strongly.

Proof $\lim \int |F - F_n| = \lim (\int (F - F_n)^+ + \int (F - F_n)^-)$, provided the limits exist.

The nice thing about $(F - F_n)^+$ is that it is dominated by F , because $F_n \geq 0$; therefore since $F \in L^1$, the dominated convergence can be used; since $F_n - F \rightarrow 0$, we get $\lim \int (F - F_n)^+ = 0$.

This trick cannot be used on $\int (F - F_n)^-$ (not dominated), but we use the additional assumption:

$$0 = \lim \int (F - F_n) = \lim \int (F - F_n)^+ - \lim \int (F - F_n)^- = 0 - \lim \int (F - F_n)^-; \text{ the proof is finished } \square$$

Now we know that F_n converge in L^1 , therefore it's

equi-integrable: $\forall E \subset \Omega$ measurable $\forall n \int_E |F_n| \leq \omega(\text{meas}(E))$ where $\omega(t) \xrightarrow{t \rightarrow 0} 0$.

We use the following refinement of the dominated convergence theorem:

Theorem (Vitali) Assume G_n are equi-integrable and $G_n \rightarrow G$ a.e. Then $G \in L^1(\Omega)$ and $G_n \rightarrow G$ in $L^1(\Omega)$ strongly.

We use this theorem on the quantity $|\vec{\sigma}_n u_n|^p =: G_n$; it is equi-integrable because $\vec{a}(\vec{\xi}) \cdot \vec{\xi} \geq c |\vec{\xi}|^p$ (coercivity) and $\vec{a}(\vec{\sigma}_n u_n) \cdot \vec{\sigma}_n u_n = F_n$ is equi-integrable.

$$\text{Thus we have } \|(\vec{\sigma}_n u_n)(\cdot)\|_p \rightarrow \|\vec{\sigma}u\|_p.$$

The norm convergence permits to upgrade the weak to the strong convergence.

NB The combination of the Minty trick with the Fatou-Schaeffe-equi-int.-Vitaly arguments is called "the Minty-Browder argument".

In the sequel, we'll give a "more straight forward" version of the argument, using Young measures to characterize weak convergence.

Pg It follows that the energies converge: $\int \eta_n[u_n] \rightarrow \int \eta[u]$.
Actually, Γ -convergence arguments could have been used for the variational case. (This ends the convergence proof. \square)

(X) Now let us look at a slightly different equation and at a different finite volume scheme.

$p(x)$ Laplacian with Dirichlet BC

let $p: \Omega \mapsto [p_-, p_+] \subset (\frac{1}{2}, \infty)$

be a sufficiently regular (continuous, with log-Hölder variable exponent modulus of continuity)

Consider

(MP₂) $\begin{cases} -\operatorname{div}(\underbrace{|\nabla u|^{p(x)-2}}_{\text{denoted } \vec{a}(x, \nabla u)} \nabla u) = f & \text{with } f \in (L^{p_-})^* \text{ (for simplicity)} \\ u|_{\partial\Omega} = 0 \end{cases}$

Weak sol

A weak solution of (MP₂) is $u \in W_0^{1,p(\cdot)}(\Omega)$ such that for all $v \in W_0^{1,p(\cdot)}(\Omega)$ ($C_0^\infty(\Omega)$ is enough) one has

(w₂) $\int_{\Omega} \vec{a}(x, \nabla u) \cdot \nabla v = \int_{\Omega} f v$

Optimization

$u = \operatorname{argmin}_{v \in W_0^{1,p(\cdot)}(\Omega)} \mathcal{J}[v], \quad \mathcal{J}: v \mapsto \int_{\Omega} \frac{1}{p(x)} |\nabla v|^{p(x)} dx - \int_{\Omega} f v$

There exists a unique weak solution, which is the unique minimizer of the energy \mathcal{J} .

DFV meshing and operators

Take a possibly non-conformal unstructured mesh with centers located inside cells (barycenters, etc)



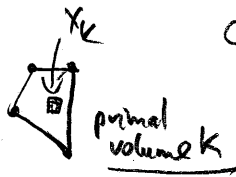
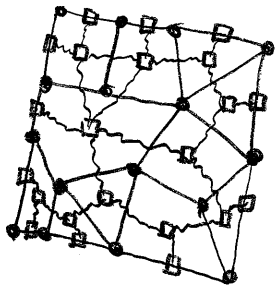
non-conformity

Create the dual mesh

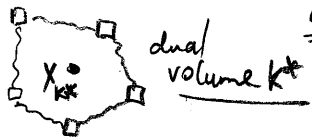
(primal vertices \equiv dual centers)
(primal centers \equiv dual vertices)

NB At the boundary: put "degenerate primal volumes" that are flat

They are called boundary volumes, and no unknown is associated with them (Dirichlet BC is used instead). We also have boundary dual volumes, they have centers located on $\partial\Omega$; the Dirichlet BC is also used in these volumes.



primal volume K



dual volume K^*

Unknowns:

u_K at points x_K that are not on $\partial\Omega$

u_{K^*} at points x_{K^*} that are not on $\partial\Omega$.

Diamonds: isolate interaction between two couples of neighbours and dual neighbours:

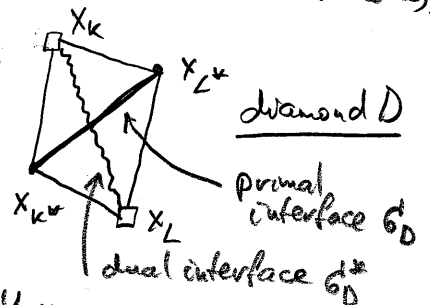
Discrete spaces: $((u_K)_K, (u_{K^*})_{K^*}) \in \mathbb{R}^{\#K} \times \mathbb{R}^{\#K^*}$ (denoted u_h)

$(\vec{F}_D)_D \in (\mathbb{R}^2)^{\#D}$ (denoted \vec{F}_h)

Scalar products: $\|u_h, v_h\|_h := \frac{1}{2} \sum_K m_K u_K v_K + \frac{1}{2} \sum_{K^*} m_{K^*} u_{K^*} v_{K^*}$

$\|\vec{F}_h, \vec{G}_h\|_h := \sum_D m_D \vec{F}_D \cdot \vec{G}_D$

Discrete divergence: Standard FV approach. Eg. $\operatorname{div}_{K^*} \vec{F}_h := \frac{1}{m_{K^*}} \sum_{D: D \cap K^* \neq \emptyset} m_D \vec{F}_D \cdot \vec{n}_{K^*}^D$



Diamond D

primal interface ∂_D

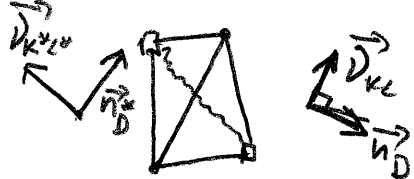
dual interface ∂_D^*

Discrete gradient This is the key idea (Kermeline; Domelevo, Omnès).

We reconstruct the gradient per diamond (according to the definition of the space of discrete vectors) as follows:

$$\vec{\nabla}_D u = \text{the unique vector of } \mathbb{R}^2 \text{ st } \begin{cases} \text{Proj}_{\vec{D}_{kL}}(\vec{\nabla}_D u) = \frac{u_L - u_k}{d_{kL}} \vec{D}_{kL} \\ \text{Proj}_{\vec{D}_{k^*L^*}}(\vec{\nabla}_D u) = \frac{u_{L^*} - u_{k^*}}{d_{k^*L^*}} \vec{D}_{k^*L^*} \end{cases}$$

Explicit formula: this was calculated in the proof of Lemma Rec 3:

$$(OGr) \quad \vec{\nabla}_D u = \frac{1}{\cos d_D} \left(\frac{u_L - u_k}{d_{kL}} \vec{n}_D + \frac{u_{L^*} - u_{k^*}}{d_{k^*L^*}} \vec{n}_D \right)$$


d_D is the angle between \vec{D}_{kL} and $\vec{D}_{k^*L^*}$.

Important particular case: $d_D = 90^\circ$ (orthogonal mesh)

In this case, we can prove the maximum principle for the scheme.

Discrete duality holds true!! $\llbracket -\text{div } \vec{F}_h, v_h \rrbracket_h = \{ \vec{F}_h, \vec{\nabla}_h v_h \}$
 for all \vec{F}_h and for all v_h such that the boundary values of v_h are zero.

Scheme: pick a value p_D per diamond D

(e.g. $p_D = \int_D p(x) dx$, or $p_D = p(x_D)$, $x_D \in D$)

or $p_D = \sup_D p(\cdot)$ or $p_D = \inf_D p(\cdot)$

Define $\vec{a}_D(\vec{s}) := |\vec{s}|^{p_D-2} \vec{s}$, write $\vec{a}_h(\cdot)$ for $(\vec{a}_D(\cdot))_D$,

$$\begin{cases} -\text{div}_h [\vec{a}_h(\vec{\nabla}_h u_h)] = f_h & f_h = p_h f \\ \text{with } u_h = 0 \text{ on boundary volumes.} \end{cases}$$

We can again write, thanks to the discrete duality,

(WDEq₂) for all v_h which is zero on the boundary volumes, $\{ \vec{a}_h(\vec{\nabla}_h u_h), \vec{\nabla}_h v_h \}_h = \llbracket p_h, v_h \rrbracket_h$.

Uniqueness immediate (same proof as for (WDEq))

Variational nature immediate, for $J[v_h] = \sum_D m_D \frac{|\vec{\nabla}_h v_h|^{p_D}}{p_D}$.

Existence From the connection to the minimization problem (a priori estimates + topological degree are also ok)

A priori estimate: requires the discrete Poincaré inequality.

Here we use the fixed $W_0^1 P$ -discrete Poincaré (proof: Eyraud, Gallouët, Herbin IMANNA, Sushi scheme) and the assumption $f \in L^p(\mathbb{R}^2)$

The proof may use translations (cf. the proof of Proposition Cmp) in L^p , but the Gallouët idea is much nicer (L^1 translates + the Nirenberg idea) as in the continuous case

The estimate obtained is: $\sum_D m_D |\vec{\nabla}_h u_h|_{p_D} \leq C(\rho)$ uniformly in h ,
 i.e., $\int_{\Omega} |(\vec{\nabla}_h u_h)(x)|^{p_h(x)} dx \leq C$.

Discrete compactness Uses again discrete Poincaré.

But it is essentially Prop. Comp, provided we do it on $W_0^{1,p}(\Omega)$.

- Thus:
- $u_h(\cdot) \rightarrow u(\cdot)$ strongly in L^1 (and in L^p , from discrete compact Sobolev embedding)
 - $(\vec{\nabla}_h u_h)(\cdot) \rightarrow \vec{\nabla} u(\cdot)$ weakly in $L^p(\Omega)$.

Consistency lemmas Roughly speaking, the same as for the co-volume scheme

Proof of convergence The approach remains the same:

- write (WDEq₂) with a test function $\varphi_h = ((\varphi(x_k))_k, (\varphi(x_{k+1}))_{k+1})$
- pass to the limit using compactness; where φ is a test function on Ω
- identify $\vec{\gamma}(\cdot)$ to $\vec{a}(\cdot, \vec{\nabla} u(\cdot))$.

Two arguments: • creation of $\vec{\gamma}$ become more difficult.
 • its identification

Indeed, there seems to be no universal space in which all the discrete norms could be calculated! The exponent $p_h(\cdot)$ depends on h .

We can manage to put all the functions $(\vec{\nabla}_h u_h)(\cdot)$ into $L^{p(\cdot)}(\Omega)$ if we pick $p_D = \sup p(\cdot)$.
 But then $\vec{a}_h(\vec{\nabla}_h u)(\cdot)$ is not in $L^{p(\cdot)}(\Omega)$! So we loose duality, more exactly, the duality framework should move with h .

(XII) A way out of the trouble is "to pull everything down to L^1 convergence".
 The key fact is the fundamental theorem for Young measures:

Theorem FTYM (J. Ball, -, N. Hungerbühler)

Assume Ω is a bounded domain, and $(v_n)_n$ is an equi-integrable sequence of functions with values in \mathbb{R}^d . (\Leftrightarrow weakly compact in $L^1(\Omega)$ by the Dunford-Pettis Criterion, see e.g. [Brézis])

Then, upon extraction of a subsequence, there exists a family $(\nu_x(\cdot))_{x \in \Omega}$ of probability measures on \mathbb{R}^d , measurable on the ad hoc sense, such that for all Carathéodory function F on $\Omega \times \mathbb{R}^d$

int 1/41) $[F(\cdot, v_n(\cdot)) \text{ equi-integrable}] \Rightarrow [\exists \lim_{n \rightarrow \infty} \int_{\Omega} F(x, v_n(x)) dx = \int_{\Omega} \left(\int_{\mathbb{R}^d} F(x, \vec{\lambda}) d\nu_x(\vec{\lambda}) \right) dx]$

Moreover, we have

$[\nu_x(\cdot) = \delta(\cdot - v(x)) \text{ a.e.}] \Leftrightarrow [v_n \text{ converge to } v \text{ in measure on } \Omega]$,

Finally, if $v_n = (w_n, z_n)$ and $(\nu_x(\cdot))_x$ is the Young measure for v_n

- $z_n \rightarrow z$ in measure (so that $(\delta(\cdot - z(x)))_x$ is the Young measure for z_n)

then $\nu_x = \mu_x \otimes \delta(\cdot - z(x))$ where $(\mu_x(\cdot))_x$ is the Young measure for w_n .

Here is the way we use Theorem FITYM.

- $(\vec{\nabla}_h u_n)_n$ is a sequence bounded in $L^p(\Omega)$, thus it is equi-integrable on Ω . (Exercise)
- We pass to a subsequence and get a Young measure $(\nu_x^i)_x$
- We first use (IntYM) on $F(x, \vec{\xi}) = \vec{\xi}^T g(x)$ with $g \in L^\infty(\Omega)$. Equi-integrability is clear.

The outcome is:

$$\int_{\Omega} (\vec{\nabla}_h u_n)(x) g(x) dx = \int_{\Omega} g(x) \left(\int_{\mathbb{R}^2} \vec{\lambda} d\nu_x(\vec{\lambda}) \right) dx, \forall g \in L^\infty(\Omega).$$

This means that $\vec{\nabla}_h u_n$ weakly converge to $\int_{\mathbb{R}^2} \vec{\lambda} d\nu_x(\vec{\lambda})$ on $L^1(\Omega)$; because we already know that $\vec{\nabla}_h u_n \rightarrow \vec{\nabla} u$ in $L^p(\Omega)$ (from compactness proposition), we get

$$\boxed{(\vec{\nabla} u)(x) = \int_{\mathbb{R}^2} \vec{\lambda} d\nu_x(\vec{\lambda}) \text{ a.e. on } \Omega} \quad (\text{YEq 1})$$

- Then we use Theorem FITYM on the couples $(\vec{\nabla}_h u_n, p_n)$ and use the fact that $p_n \rightarrow p$ on Ω (convergence in measure is enough) [If $p(\cdot)$ is continuous, any choice of p_0 leads to the L^∞ convergence ...]

Then the associated Young measure is $\nu_x(\cdot) \otimes \delta(\cdot - p(x))$.

We would like to use (IntYM) with

$$F(x, \vec{\xi}, p) = |\vec{\xi}|^{p(x)-2} \vec{\xi}^T g(x), \quad g \in L^\infty(\Omega);$$

but we have to check the equi-integrability of $F(\cdot, (\vec{\nabla}_h u_n)(\cdot), p_n(\cdot))$.

This follows from the estimate $\int_{\Omega} |\vec{\nabla}_h u_n(x)|^{p_n(x)} dx \leq \text{const}$

and the variable exponent Hölder inequality. The proof is a bit technical and left to a reader with motivation $\geq \delta > 0$.

Then the outcome is, making g vary in $L^\infty(\Omega)$:

$$\boxed{\vec{a}_h(\vec{\nabla}_h u_n)(\cdot) \rightarrow \vec{\chi}(\cdot) := \int_{\mathbb{R}^2} |\vec{\lambda}|^{p(x)-2} d\nu_x(\vec{\lambda}) \text{ in } L^1(\Omega) \text{ weakly}} \quad (\text{YEq 2})$$

- Two things remain:
- identify $\vec{\chi}(\cdot)$ to $\vec{a}(\cdot, \vec{\nabla} u)$
 - prove that $u \in W_0^{1,p(\cdot)}(\Omega)$ (for the time being, we only have $u \in W_0^{1,p}(\Omega)$)

• We prove that $\int_{\Omega} |\vec{\nabla} u|^{p(x)} dx < \infty$;

then one can use the fact [p-log-Hölder continuous] \Rightarrow a fortiori, $W_0^{1,p}(\Omega)$ is ok here

! log-Hölder condition is important here [$u \in W_0^{1,p}(\Omega)$ iff $u \in W_0^{1,1}(\Omega)$ and $\vec{\nabla} u \in L^{p(\cdot)}(\Omega)$].

The estimate follows

by using $F_k(\vec{\xi}, p) = \begin{cases} |\vec{\xi}|^p, & |\vec{\xi}| \leq k \\ k^p, & \text{otherwise.} \end{cases}$ It is L^∞ , thus $F_k(\vec{\nabla}_h u_n(\cdot), p_n(\cdot))$ are equi-integrable.

Hence $\int_{\Omega} \int_{\mathbb{R}^2} F_k(\lambda, p(x)) d\nu_x(\vec{\lambda}) = \lim_{k \rightarrow \infty} \int_{\Omega} F_k(\vec{\nabla}_h u_n, p_n) \leq \int_{\Omega} |\vec{\nabla}_h u_n|^{p_n} \leq \text{const.}$

But F_k is increasing in k ; and the limit is $F(\vec{\xi}, p) = |\vec{\xi}|^p$. Thus by Beppo Levi theorem, $\int_{\Omega} \int_{\mathbb{R}^2} |\vec{\lambda}|^{p(x)} d\nu_x(\vec{\lambda}) \leq \text{const.}$

It remains to use the Jensen inequality for the probability measure $\nu_x(\cdot)$:

$$\text{for a.e. } x, \quad |\text{div}(\vec{x})|^{p(x)} = \left| \int_{\mathbb{R}^2} \vec{x} d\nu_x(\vec{\lambda}) \right|^{p(x)} \leq \int_{\mathbb{R}^2} |\vec{\lambda}|^{p(x)} d\nu_x(\vec{\lambda}).$$

Hence $\int_{\Omega} |\text{div}(\vec{x})|^{p(x)} dx \leq \text{const} < \infty$; we have therefore established that $u \in W_0^{1,p(x)}(\Omega)$.

Now we would like to prove the following inequality:

$$\text{DivRot} \int_{\Omega} \left(\int_{\mathbb{R}^2} |\vec{\lambda}|^{p(x)-2} \vec{\lambda} d\nu_x(\vec{\lambda}) \right) \left(\int_{\mathbb{R}^2} \vec{\lambda} d\nu_x(\vec{\lambda}) \right) dx \geq \int_{\Omega} \left(\int_{\mathbb{R}^2} |\vec{\lambda}|^{p(x)} d\nu_x(\vec{\lambda}) \right) dx$$

$\cong \vec{x} \qquad \qquad \qquad \cong \vec{\nabla} u$

This is reminiscent of the inequality that was the starting point of the tricky Minty-Browder argument. $\int_{\Omega} \vec{x} \cdot \vec{\nabla} u \geq \lim \int_{\Omega} \vec{a}(\vec{\nabla}_h u) \cdot \vec{\nabla}_h u$

We have $\int_{\Omega} \vec{x} \cdot \vec{\nabla} u = \int_{\Omega} f u$ because we were able to pass to the limit in the (WDEq₂) and then put $v=u$ (recall that $u \in W_0^{1,p(x)}(\Omega)$ and $C_c^\infty(\Omega)$ is dense in $W_0^{1,p(x)}(\Omega)$).

Moreover, $\int_{\Omega} f u = \lim_{n \rightarrow \infty} \int_{\Omega} f u_n = \lim \int_{\Omega} a_n(\vec{\nabla}_h u_n) \cdot \vec{\nabla}_h u_n \cong \lim \int_{\Omega} |\vec{\nabla}_h u_n|^{p_n(x)} dx$.

Thus it would be sufficient to show that

$$\int_{\Omega} \int_{\mathbb{R}^2} |\vec{\lambda}|^{p(x)} d\nu_x(\vec{\lambda}) \leq \lim \int_{\Omega} |\vec{\nabla}_h u_n|^{p_n(x)} dx.$$

But we have just shown this (in a weaker form) in the previous argument!

Rq Why have we used truncations?

It's because we could not take directly $F(\vec{\xi}, p) = |\vec{\xi}|^p$:

The equi-integrability of $|\vec{\nabla}_h u_n(\cdot)|^{p_n(\cdot)}$ was missing, only an L^1 bound was available!

To conclude, we show that ("DivRot") is a compactification relation; this is because $\vec{a}(x, \cdot) = |\cdot|^{p(x)-2} \cdot$ is monotone!

We calculate

$$\begin{aligned} \Delta &:= \int_{\Omega} \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} (\vec{a}(x, \vec{\lambda}) - \vec{a}(x, \vec{\mu})) \cdot (\vec{\lambda} - \vec{\mu}) d\nu_x(\vec{\lambda}) d\nu_x(\vec{\mu}) = \\ &= \int_{\Omega} \underbrace{\left(\int_{\mathbb{R}^d} d\nu_x(\vec{\mu}) \right)}_{=1} \left(\int_{\mathbb{R}^d} \vec{a}(x, \vec{\lambda}) \cdot \vec{\lambda} d\nu_x(\vec{\lambda}) \right) + \text{identical term with } \lambda, \mu \text{ exchanged} \\ &\quad - \int_{\Omega} \left(\int_{\mathbb{R}^d} \vec{a}(x, \vec{\lambda}) d\nu_x(\vec{\lambda}) \right) \left(\int_{\mathbb{R}^d} \vec{\mu} d\nu_x(\vec{\mu}) \right) + \text{identical term with } \lambda, \mu \text{ exchanged} \\ &\leq 0 \text{ from ("DivRot").} \end{aligned}$$

\uparrow here we can now substitute $\vec{\mu}$ by $\vec{\lambda}$

But $\Delta \geq 0$ because the integrand is ≥ 0 , by monotonicity of \vec{a} ! Thus $\Delta = 0$, moreover, the integrand is zero $d\nu_x(\vec{\lambda}) d\nu_x(\vec{\mu}) dx$ - a.e.

By the strict monotonicity of $\vec{a}^\rightarrow(x, \cdot)$, this means that the support of $d\mu_x(\cdot)$ should be reduced to one point (exercise).

This means, for a.e. $x \in \Omega$ $\nu_x(\cdot)$ is a Dirac measure concentrated at some $\vec{v}(x)$: By Theorem RTYM, we deduce that $\vec{\nabla}_h u_h(\cdot) \rightarrow \vec{v}(\cdot)$ on measure, and also weakly ^{in $L^1(\Omega)$} . From the Proposition on compactness, we already know that $\vec{\nabla}_h u_h(\cdot) \rightarrow \vec{\nabla} u(\cdot)$ in $L^1(\Omega)$.

We conclude that $\nu_x(\cdot) = \delta(\cdot - (\vec{\nabla} u)(x))$ a.e. on Ω .

Now from the representation formula for $\vec{\chi}^\rightarrow$,

$$\vec{\chi}(x) = \int_{\mathbb{R}^2} |\vec{\chi}|^{p(x)-2} \vec{\chi} d\delta(\vec{\chi} - \vec{\nabla} u(x)) = |\vec{\nabla} u(x)|^{p(x)-2} \vec{\nabla} u(x).$$

The identification is complete! We conclude that $u = \lim u_h$ solves (MP_2) . \square

In addition, upon extraction of a subsequence, $\vec{\nabla}_h u_h(\cdot) \rightarrow \vec{\nabla} u$ a.e.

Then we can upgrade the weak convergence a bit; we get in particular

$$\int_{\Omega} |\vec{\nabla}_h u_h|^{p_h} \rightarrow \int_{\Omega} |\vec{\nabla} u|^p \quad (\text{convergence of the energy } J_h[u_h] \text{ to } J[u] \text{ follows}).$$

NB We cannot state that $\vec{\nabla}_h u_h$ converge to $\vec{\nabla} u$ in $L^{p(\cdot)}(\Omega)$ because $p_h(\cdot)$ is different from $p(\cdot)$; but, interpolating the uniform bound on $\int |\vec{\nabla}_h u_h|^{p_h}$ with the a.e. convergence of $\vec{\nabla}_h u_h(\cdot)$ to $\vec{\nabla} u(\cdot)$, we get strong convergence in $L^q(\Omega)$ with $q < p$ (recall that $p_h \rightarrow p$ on the L^∞ norm).

XII Let's give essential ideas for a priori error estimates.

Properties of $\vec{a}^\rightarrow: \vec{\xi}^\rightarrow \rightarrow |\vec{\xi}^\rightarrow|^{p-2} \vec{\xi}^\rightarrow$ One has to distinguish between $p \geq 2$ and $1 < p < 2$.

- ① $\begin{cases} |\vec{a}^\rightarrow(\vec{\xi}^\rightarrow) - \vec{a}^\rightarrow(\vec{\eta}^\rightarrow)| \leq C |\vec{\xi}^\rightarrow - \vec{\eta}^\rightarrow|^{p-1}, & 1 < p \leq 2 \\ |\vec{a}^\rightarrow(\vec{\xi}^\rightarrow) - \vec{a}^\rightarrow(\vec{\eta}^\rightarrow)| \leq C |\vec{\xi}^\rightarrow - \vec{\eta}^\rightarrow| (|\vec{\xi}^\rightarrow|^p + |\vec{\eta}^\rightarrow|^p)^{\frac{p-2}{p}}, & p \geq 2 \end{cases}$: $\vec{a}^\rightarrow(\cdot)$ is $C^{0,\alpha}$ Hölder, $\alpha = p-1$ locally Lipschitz

Remark: If we denote $\vec{\Sigma}^\rightarrow = |\vec{\xi}^\rightarrow|^{p-2} \vec{\xi}^\rightarrow$, then $\vec{\xi}^\rightarrow = |\vec{\Sigma}^\rightarrow|^{\frac{1}{p-2}} \vec{\Sigma}^\rightarrow$ (exercise)

Therefore the inverse of $\vec{a}_p^\rightarrow(\cdot)$ is $\vec{a}_p^\rightarrow(\cdot)$.

Consequence: $1 < p \leq 2$, $\vec{a}^\rightarrow(\cdot)$ is locally Lipschitz
 $p \geq 2$, $\vec{a}^\rightarrow(\cdot)$ is $C^{0,\beta}$ Hölder, $\beta = p-1 = \frac{1}{p-1}$

The Hölder continuity results are the basis of most straight forward error estimates:

roughly speaking, from the order of consistency (e.g., h)

we first extract an estimate on consistency of fluxes $\vec{F} = \vec{a}^\rightarrow(\cdot)$;

then it becomes of order h^α , $\alpha = \begin{cases} p-1, & p \leq 2 \\ 1, & p \geq 2 \end{cases}$.

We get an estimate on $|\vec{a}^\rightarrow(\vec{\nabla}_h u_h) - \vec{a}^\rightarrow(\vec{\nabla}_h P_h u)|$ in h^α , and inverting $\vec{a}^\rightarrow(\cdot)$, we deduce for $|\vec{\nabla}_h u_h - \vec{\nabla}_h P_h u|$ in $(h^\alpha)^\beta$ where $\beta = \begin{cases} \frac{1}{p-1}, & p \leq 2 \\ \frac{1}{p-1}, & p \geq 2 \end{cases}$. Finally, the convergence order is $h^{\min\{p-1, \frac{1}{p-1}\}}$.

② To be precise, instead of the Hölder continuity of $\vec{a}^{-1}(\cdot)$ we rather use

$$\begin{cases} (\vec{a}(\vec{z}) - \vec{a}(\vec{y})) \cdot (\vec{z} - \vec{y}) \geq C |\vec{z} - \vec{y}|^2 / (|\vec{z}|^p + |\vec{y}|^p)^{\frac{p-2}{p}}, & 1 < p \leq 2 \\ (\vec{a}(\vec{z}) - \vec{a}(\vec{y})) \cdot (\vec{z} - \vec{y}) \geq C |\vec{z} - \vec{y}|^p, & p \geq 2. \end{cases}$$

③ Everything comes from the following general inequalities:

$\forall p \geq 1 \quad \forall \delta > 0 \quad \exists C_1, C_2$ depending on p, δ s.t.

$$|\vec{a}(\vec{z}) - \vec{a}(\vec{y})| \leq C_1 |\vec{z} - \vec{y}|^{1-\delta} (|\vec{z}| + |\vec{y}|)^{p-2+\delta}$$

$$(\vec{a}(\vec{z}) - \vec{a}(\vec{y})) \cdot (\vec{z} - \vec{y}) \geq C_2 |\vec{z} - \vec{y}|^{2+\delta} (|\vec{z}| + |\vec{y}|)^{p-2-\delta}$$

These inequalities allow to improve the Hölder continuity exponents when the solutions have gradients far enough from zero or infinity,

Reference Barrett, Liu 1993.

thus to improve error estimates for solutions for which some additional information is available.

How to get consistency orders?

A typical estimate that is needed looks as follows:

(on covolume mesh) compare $\int_{\Omega} \vec{a}(\vec{\nabla} u_e) - \vec{a}(\vec{\nabla}_{T(\Omega)} P_h u_e)$, where u_e is the exact solution.

Using Hölder continuity of \vec{a} , we reduce the question to comparison of $(\vec{\nabla} u_e)(x)$ for fixed x and of $\vec{\nabla}_{T(\Omega)} P_h u_e$. Recall that $P_h u_e = (\int_k u_e)_k$ or $(u_e|_k)_k$

• If one assumes $u_e \in C^{st_{tag}}$, one can use Taylor expansions to prove that for C^2 solutions, the difference is order h .

Moreover, on meshes possessing symmetry properties and in case u_e is smooth and \vec{a} is smooth ($p \geq st_{tag} \geq 2$ or $p=2$), summing on edges we observe compensations that may lead to better convergence orders (cancellation in Taylor expansions).

• If one assumes $u_e \in W^{2,p}$, which seems more realistic, then the essential tool is the so-called Hadamard formula: say, if $\begin{cases} x \text{ runs over } k \\ y \text{ runs over } L \end{cases}$

$$\begin{aligned} u_L - u_k &= \int_L u_e(x) dx - \int_k u_e(y) dy = \frac{1}{m_k m_L} \iint_{k \times L} (u_e(x) - u_e(y)) dx dy \\ &= \frac{1}{m_k m_L} \iint_{k \times L} \left(\int_0^1 \vec{\nabla} u_e(t x + (1-t)y) dt \right) \cdot (x-y) dx dy \end{aligned}$$

Then a change of variables

$$(x, y, t) \mapsto (x, z, t)$$

brings an estimate in terms of $\int_{k \cup L} \vec{\nabla} u_e(z) dz$.

$$\text{We end up with } \|\vec{\nabla}_h u_h - \vec{\nabla} u\|_p \leq C h^{\min\{p-1, \frac{1}{p}\}}$$

Rq • If we only assume the minimal $(W^{1,p}(\Omega))$ regularity of u_e , no convergence order can be obtained: we have shown convergence, but the rate of convergence can be as slow as desired.

• Assuming $u_e \in W^{2,p}(\Omega)$ is nice (standard in the context of the Laplacian) but we can prove it only for $p \leq 2$ (with $f \in L^p(\Omega)$, which is more restrictive than $f \in L^1$ for which we have shown convergence).

Rg. For $p > 2$, no condition on f is known to ensure $u \in W^{2,p}(\Omega)$,
So, the result is not satisfactory.

- Optimal regularity for $p > 2$ and $f \in L^p$ was proved by J. Simon '81: this is the so-called Besov $B_{\infty}^{1+\frac{1}{p-1}, p}(\Omega)$ regularity, in case the solution admits an extension outside Ω (on a rectangle with Dirichlet homogeneous BC, we achieve such extension by reflexions and periodization).
- How to use this regularity? The approach described above would lead to $((h^{\frac{1}{p-1}})^d)^p$ convergence order, i.e., $h^{\frac{1}{(p-1)2}}$ (for $p > 2$).
- There is a very smart approach due to Tyukhtin '82, popularized by S.S. Chow '89: it uses the minimization properties for both discrete and continuous energies, and the outcome is the $h^{\frac{2}{(p-1)2}}$ convergence order.
- Numerically, we observe that for $u \in$ which has precisely the $B_{\infty}^{1+\frac{1}{p-1}, p}$ regularity (and not more), the order $h^{\frac{2}{(p-1)2}}$ is what happens in practice. Thus, this estimate is optimal!

Some details on error estimates are given in the "Annexes" pages (in French) at the end of this document.

References

- Existence for p-laplacian & Co!
 - J.L. Lions, 1969 "Quelques problèmes aux limites non linéaires", Chapter II. Dunod, Paris
- Finite element approximation of the p-laplacian:
 - S.S. Chow, 1989, Numerische Mathematik
 - J.W. Barrett, W.B. Liu, 1993, J. Math. Anal. Appl.
 - || — || — || — || 1993, Math. Comp.
- Covolume Finite Volume schemes
 - Walkington 1996, SIAM J. Numer. Anal.
 - Afif, Amaziane 2002, Comput. Math. Appl. Mech. Eng.
 - Handlovičová, Mikula, Sgallari 2003, Numer. Math.
 - Handlovičová, Mikula 2008, Appl. Math.
- DDFV Finite Volume schemes
 - Kernevale 1998, CRAS; 2000, J. Comput. Phys.; 2009, Comp. Math. Appl. Mech. Eng.
 - Domelevo, Omnes 2005, M2AN
 - Andreianov, Boyer, Kubert 2007, Num. Math. PDE
 - Andreianov, Bendahmane, Kubert, preprint to appear on HAL
- Schemes on cartesian meshes and error estimates
 - Andreianov, Boyer, Kubert, 2004, M2AN; 2005, Numer. Math.; 2006, IMA J. Num. Anal.

- Other FV schemes
 - Eymard, Gallouët, Herbin "Finite volume methods" 2000, personal page of R. Herbin
 - Coindre, Vila, Villedieu 1999, M2AN
 - Andreianov, Gutnic, Wittbold 2004, SIAM J. Num. Anal.
 - Droniou, 2006, M2AN.
 - Eymard, Gallouët, Herbin 2009, IMA J. Num. Anal.
 - Eymard, Gallouët, Herbin HAL preprint, January 2009
- Momentic schemes
 - Shashkov, Lipnikov, Brezzi, Manzini 2005-...

Méthodes de volumes finis pour les problèmes elliptiques non linéaires

B. Andreianov (Besançon) & R. Boyer, F. Hubert (Marseille)

Problème à discrétiser $\left\{ \begin{array}{l} -\operatorname{div} \varphi(z, \nabla u_e) = f(z), \quad z = (x, y) \in \Omega \subset \mathbb{R}^2 \\ + \text{CB sur } u_e \text{ (Dirichlet homogène / inhomogène)} \end{array} \right. \quad u_e: \text{ solution exacte}$

Prototype: le p -laplacien, avec $\varphi(z, \xi) = |\xi|^{p-2} \xi, \quad 1 < p < +\infty$.

Propriétés structurelles de φ :

Pour existence, unicité, convergence: • monotone de φ en ξ
• coercivité: $\varphi(\xi) \cdot \xi \geq c|\xi|^p$

• hypothèse de croissance tq $u \mapsto -\operatorname{div} \varphi(\cdot, \nabla u)$ agisse de $W_0^{1,p}$ dans $W^{-1,p'}$

• (cas particulier) structure variationnelle: $\varphi(\xi) = \nabla \Phi(\xi)$

pour l'analyse quantitative de l'erreur:

• $\varphi: \xi \rightarrow \varphi(\xi)$ est (localement) Hölderienne d'ordre $\begin{cases} 1, & p \geq 2 \\ p-1, & 1 < p < 2 \end{cases}$

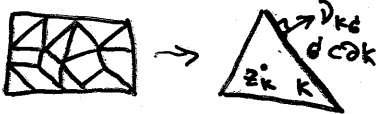
• φ^{-1} est (localement) Hölderienne d'ordre $\begin{cases} \frac{1}{p-1}, & p \geq 2 \\ 1, & 1 < p < 2 \end{cases}$.

Méthodes de volumes finis

Pour le problème linéaire sur maillages conformes:

cf. Eymard, Gallouët, Herbin "Finite Volume Methods" '05

Principe: • Partitionner Ω en volumes de contrôle, notés K ($K \in \mathcal{T}$)



• On cherche une "solution discrète", notée $u^{\mathcal{T}} = (u_K)_{K \in \mathcal{T}}$, constante par volume

• On "intègre" l'équation sur chaque volume:

$$\forall K \in \mathcal{T}, \quad a_K [u^{\mathcal{T}}] = \sum_{\text{arêtes } d_{k,l}} \int_{\sigma} (\varphi \cdot \nu_{k,l}) [u^{\mathcal{T}}] = m_{K,l} f_K$$

• On fournit $(\varphi \cdot \nu_{k,l})^{\mathcal{T}} [u^{\mathcal{T}}]$, qui représente une reconstruction du flux normal discret via l'arête d , à partir des valeurs $(u_K)_{K \in \mathcal{T}}$.

Cas linéaire: Si on peut munir chaque volume K d'un centre z_K tq $z_K z_L \perp d$, ($\varphi(z, \xi) = k(z) \xi$) où $d = \partial K \cap \partial L$ (i.e. maillage conforme), alors $\frac{u_L - u_K}{|z_L - z_K|}$ fournit la reconstruction "naturelle" de $(\varphi \cdot \nu_{k,l})^{\mathcal{T}} [u^{\mathcal{T}}]$.

Cette reconstruction préserve, au niveau discret, les propriétés structurelles du problème continu (citées ci-dessus).

Objet de l'exposé

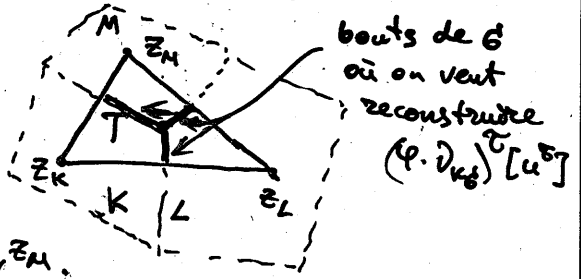
- présenter les schémas VF pour le problème non linéaire qui préservent également ces propriétés structurelles
- analyser leur convergence (en se servant des dites propriétés).

Schémas disponibles

- B.A. & M. Gutnic, P. Wittbold SIAM YNA'04

marche par reconstruction affine (morceaux sur des maillages duaux aux maillages triangulaires).

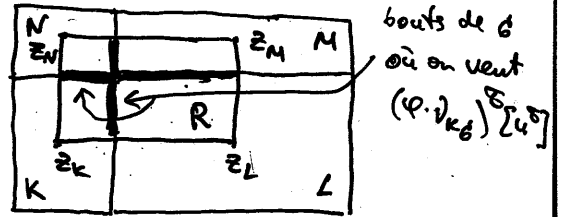
Dans chaque triangle $T = z_k z_l z_m$ on prend sur ∂T , $(\varphi \cdot \nu_k)_T [u^T] = \varphi(\nabla_T u^T) \cdot \nu_k$ où $\nabla_T u^T$ est l'interpolation affine des valeurs de u_k, u_l, u_m aux points z_k, z_l, z_m .



- B.A. & F. Boyer, F. Hubert M2AN'04; Num. Math. '05; IMANNA'06?

marche pour les problèmes variationnels sur maillages cartésiens

Dans chaque rectangle $R = z_k z_l z_m z_n$ on prend pour $(\nu_R u^R)$ une valeur reconstruite à partir de u_k, u_l, u_m, u_n . Il n'y a pas de reconstruction "naturelle", mais la structure rigide (on veut que le schéma discret dérive d'un potentiel) fait qu'il n'y a qu'une famille à un paramètre de reconstructions possibles de $(\nu_R u^R)$ (les formules sont compliquées, sauf pour maillages uniformes). On observe que l'équation d'Euler-Lagrange pour minimisation de la fonctionnelle discrète $J_T[u^T] = \sum_R m_R \Phi((\nu_R u^R))$ est exactement un schéma de volumes finis.



- B.A. & F. Boyer, F. Hubert Num. PDE'06?, proc. FVCA-4'05

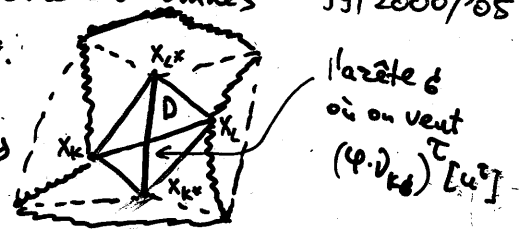
marche sur maillages "en duplex": deux maillages recouvrent, chacun, le domaine, on a deux jeux d'inconnues et d'équations discrètes.

Gain: approximation du gradient discret simple et très efficace, structure de maillages très générale (raffinement local etc.)

Idee: Coudière-Vila-Villedieu / Herminie / Domélevo-Omnès 199/2000/'05

On partitionne Ω en "diamants" $D = z_k z_{k^*} z_l z_{l^*}$.

Les points $(z_k)_k$ et $(z_{k^*})_{k^*}$ engendrent deux maillages de Ω . On prend le double jeu d'inconnues $u^T = \{(u_k)_k, (u_{k^*})_{k^*}\}$ et on écrit le schéma de volumes finis (une équation pour chaque k et une pour chaque k^*). Dans chaque diamant D , on reconstruit un vecteur $\nabla_D u^T$ et on prend



$(\varphi \cdot \nu_k)_D [u^T] = \varphi(\nabla_D u^T) \cdot \nu_k$
et idem pour k^*, l, l^* .

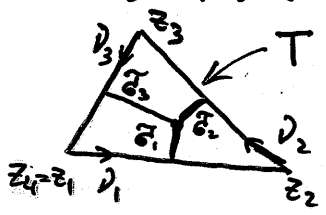
- J. Droniou '06? (cf. le cas linéaire Eymard, Droniou '05?)

marche sur maillages presque arbitraires, grâce à la relaxation du lien flux-gradient. On garde les flux φ_i via chaque arête comme inconnues et on les relie aux inconnues u^T par une relation bien choisie qui complète le système discret.

Dans tous les cas, la structure est héritée du pb. continu grâce à la dualité des opérateurs -div et grad discrets.

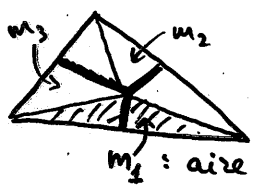
Exemple (des schéma) en continu : $\int_{\Omega} -\text{div } \varphi(\sigma_u) \cdot v = \int_{\Omega} \varphi(\sigma_u) \cdot \nabla v \quad \forall v \in W_0^{1,p}$

en discret :
 v^T nulle sur les mailles de bord



$$\sum_k a_k [u^T] v^T = \sum_k \sum_{\sigma \in \sigma_k} \int_{\sigma} (\varphi \cdot \nu_{\sigma})^T [u^T] v^T = \sum_T \sum_{i=1}^3 m_{\sigma_i} \varphi(\nabla_T u^T) \cdot \nu_i (\sigma_i \cdot \nu_i) v_i = \sum_T \varphi(\nabla_T u^T) \cdot \sum_{i=1}^3 m_{\sigma_i} \nu_i (\sigma_i \cdot \nu_i) \nu_i = \sum_T m_T \varphi(\nabla_T u^T) \cdot \nabla_T v^T$$

(on regroupe les termes par semi-arête σ dans T et on fait la sommation par parties)



où la dernière inégalité vient du calcul suivant :

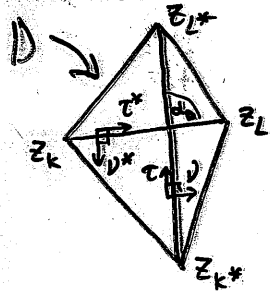
Lemme $\frac{2}{m_T} \sum_{i=1}^3 m_{\sigma_i} (\xi \cdot \nu_i) \nu_i = \xi \quad \forall \xi \in \mathbb{R}^2$

De la dualité, on déduit par ex. la coercivité :

$$\sum_k a_k [v^T] v^T = \sum_T m_T \varphi(\nabla_T v^T) \cdot \nabla_T v^T \geq C \sum_T |\nabla_T v^T|^p$$

mais aussi la monotonie, les dépendances Hölderiennes.

Exemple (schémas duplex)



Dans chaque diamant D, on reconstruit le gradient :

$$\nabla_D u^T = \frac{1}{\sin \nu_D} \left(\frac{u_L - u_k}{|z_L - z_k|} \nu + \frac{u_{L^*} - u_{k^*}}{|z_{L^*} - z_{k^*}|} \nu^* \right)$$

Lemme $\frac{1}{\sin \nu_D} (\xi \cdot \nu) \nu + (\xi \cdot \nu^*) \nu^* = \xi \quad \forall \xi \in \mathbb{R}^2$

Corollaire (div-grad dualité discrète)

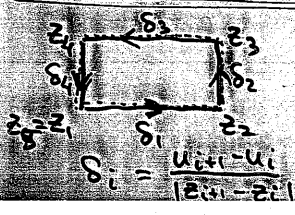
On a $\frac{1}{2} \sum_k a_k [u^T] v_k + \frac{1}{2} \sum_{k^*} a_{k^*} [u^T] v_{k^*} = \sum_D m_D \varphi(\nabla_D u^T) \cdot \nabla_D v^D$,
 pour toutes u^T, v^T fonctions duplex discrètes avec v^T nulle dans les volumes de bord.

Propriétés du schéma duplex

- inégalités de Poincaré, injections, interpolations discrètes;
- existence et unicité de solution discrète, estimation $W^{1,p}$ a priori;
- si φ dérive d'un potentiel Φ , le problème discret garde la structure variationnelle : $a_k [u^T] = \frac{\partial}{\partial u_k} \sum_D m_D \Phi(\nabla_D u^T)$.

Ces propriétés sont quantifiées par une constante $\text{reg}(h)$ (où h = taille du maillage); $\text{reg}(h)$ est censée rester bornée lorsque $h \rightarrow 0$.

Schémas sur maillages cartésiens



Si maillage uniforme : $|\nabla_R u^T|^2 = \frac{1}{2} \sum_{i=1}^4 \delta_i^2 + \frac{1}{2} (\delta_1 + \delta_3)^2$,
 où $\frac{1}{2}$ (un paramètre de "torsion") peut varier dans certaines limites. En général, on obtient une formule (avec paramètre θ) via les proportions de R.
 Ces schémas ont aussi la propriété de dualité div-grad.

Conclusion Pour les schémas "duplex" (les mêmes résultats sont/peuvent être obtenus sur maillages cartésiens et sur maillages duaux aux triangulaires)

- convergence, pour tout terme source $f \in L^p$ (et $W^{1,p}$ si on fait entrer la partie singulière dans le flux)

• estimations d'erreur en normes $W^{1,p,T}$ et L^p en $\begin{cases} h^{\frac{1}{p-1}}, & p \geq 2 \\ h^{p-1}, & 1 < p < 2 \end{cases}$ à condition que f donne lieu à une solution $u_e \in W^{1,p}$.

- Numériquement, ces taux sont pessimistes (en particulier, les ordres de convergence obtenus tendent vers zéro lorsque $p \rightarrow +\infty$ ou $p \rightarrow 1$).
- Sauf pour $1 < p \leq 2$, on ne dispose pas de résultats de régularité déjà pour le p -laplacien.

On va alors chercher à améliorer ces résultats en deux directions

① → Tant qu'à faire, on suppose la régularité aussi élevée que nécessaire pour avoir un meilleur taux de convergence: B.A, F. Boyer, F. Hubert INMATH'06

② → On considère les solutions génériques mais on utilise que leur régularité naturelle: $u_e \in \begin{cases} B_{\infty}^{1+\frac{1}{p-1}, p}, & p \geq 2 \\ W^{2,p}, & 1 < p < 2 \end{cases}$ (espace de Besov)
B.A, F. Boyer, F. Hubert Num. Math'05

Idee de la preuve d'est. d'erreur de base

Le schéma: $a_k[u^T] = m_k f_k$
L'éq. continu: $\sum_{G \in \mathcal{G}_k} \int_G \varphi(\nabla u_e) \cdot \nu_{kG} = m_k f_k$

On retranche et multiplie par $u^T - u_e^T$:

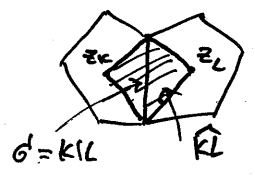
$$\begin{aligned} \sum_k (a_k[u^T] - a_k[u_e^T]) (u_k - u_{ek}) &= \\ &= \sum_k \sum_{G \in \mathcal{G}_k} m_G (\Delta f)_G (u_k - u_{ek}) = \\ &= \sum_{G=K|L} m_G |z_k - z_l| (\Delta f)_G \frac{(u_k - u_{ek}) - (u_l - u_{el})}{|z_k - z_l|} \end{aligned}$$

on le note $(\Delta f)_G^T$

$$a_k[u_e^T] + \sum_{G \in \mathcal{G}_k} m_G \int_G \varphi(\nabla u_e) \cdot \nu_{kG} - (\varphi(\nabla u_e) \cdot \nu_{kG}) [u_e^T]$$

sa joue le rôle de perturbation du terme source.

ici on retrouve la mesure de la partie de Ω entourant l'arête $G = K|L$ qui sépare K et L ; on la note $m_{K|L}$



Ainsi, tout se passe comme si $\Delta f := \sum_{K|L} (\Delta f)_{K|L} \mathbb{1}_{K|L}$ était une perturbation du terme source.

- Or, on a :
- $|\nabla^T u_e^T - \nabla u_e| \sim h$ lorsque u_e a deux dérivées
 - la continuité Hölder de φ fait que $\|\Delta f\|_{2,p'} \sim \begin{cases} h, & p \geq 2 \\ h^{p-1}, & 1 < p \leq 2 \end{cases}$
 - la continuité Hölder de φ' fait que $\|\nabla^T u_e^T\|_{2,p} \sim \begin{cases} \|\Delta f\|_{\frac{1}{p-1}}, & p \geq 2 \\ \|\Delta f\|, & 1 < p \leq 2 \end{cases} \sim \begin{cases} (h^{\frac{1}{p-1}})^{\frac{1}{p-1}}, & p \geq 2 \\ (h^{p-1})^{\frac{1}{p-1}}, & 1 < p \leq 2 \end{cases}$

Améliorations

(16)

① Sur maillages cartésiens uniformes, on a $(\Delta f)_\delta \sim h$, pour $p \geq 2$,
 mais $\sum_{\delta \in \mathcal{T}_h} m_\delta (\Delta f)_\delta \sim h^4$ pour $p \geq 4$ (ou $p=2$) et des solutions $u_e \in W^{4, q, qch}$

Alors $\sum_K \sum_{\delta \in \mathcal{T}_K} m_\delta (\Delta f)_\delta (u_K - u_{eK}) = \sum_K m_K (\tilde{\Delta f})_K (u_K - u_{eK})$ avec $(\tilde{\Delta f})_K = \frac{1}{m(K)} \sum_{\delta \in \mathcal{T}_K} m_\delta (\Delta f)_\delta$.

On peut considérer $\tilde{\Delta f} = \sum_K (\tilde{\Delta f})_K \mathbb{1}_K$ comme une perturbation de terme source; cette perturbation est maintenant de l'ordre h^2 et on double le taux de convergence:

$$\| \nabla^T e^T \|_{L^p} \sim \begin{cases} h^{\frac{2}{p-1}}, & p \geq 4 \text{ ou } p=2 \\ h^{\frac{p-2}{p-1}}, & 3 < p < 4 \end{cases}$$

On peut faire encore mieux (idée de Barrett-Liu '93 + techniques d'interpolation) si on remarque que la continuité Hölder de φ' s'améliore lorsque l'on contrôle les petites valeurs de $|\nabla u_e|$: loin de zéro, φ' est Lipschitz!

Alors en fonction de ν tq $\int |\nabla u_e|^{-\nu} < +\infty$, on peut améliorer le taux.

Cas extrêmes: $\nu = 0$, alors les taux sont donnés ci-dessus ($h^{\frac{2}{p-1}}$, $p \geq 4$)
 $\nu = +\infty$ (i.e., il s'agit de solutions $W^{4, q, qch}$ sans points critiques)

→ ordres $\left. \begin{array}{l} h^2 \text{ en norme } L^\infty \text{ et } W^{1/2, \infty} \\ h^{\frac{4}{p}} \text{ en norme } W^{1/p, \infty} \end{array} \right\} \text{ pour } p \geq 2$
 numériquement, cet ordre se confirme comme optimal → ordre $\left. \begin{array}{l} h^2 \text{ en norme } W^{1/p, \infty} \\ h^{\frac{3p-2}{p}} \text{ en norme } L^\infty \end{array} \right\} \text{ pour } 1 < p < 2.$

Pour $0 < \nu < +\infty$ (on a notamment $\nu = \delta(p-2)$ si $|f|^{-\delta} \in L^1$), on montre des ordres intermédiaires.

② Sur maillages cartésiens uniformes (ou gentiment raffinés - work in progress)

On adapte au cadre EF l'idée de Tyukhtin '82, utilisée par Chow '89 et on fait l'usage de la régularité naturelle Besov $B^{1+\frac{1}{p}, p}$ de u_e .

On se place dans le cadre variationnel: e.g. pour le p -laplacien,

$$a_K = \frac{\partial}{\partial u_K} \mathcal{J}_\tau [u^T] \text{ où } \mathcal{J}_\tau [u^T] = \frac{1}{p} \sum_R m_R |\nabla_R u^T|^p - \sum_K m_K f_K u_K,$$

la fonctionnelle \mathcal{J}_τ approche la fonctionnelle continue $\mathcal{J}[u] = \frac{1}{p} \int_\Omega |\nabla u|^p - \int_\Omega f u$.

Rq Numériquement, on peut alors résoudre le schéma par des méthodes simples du type descente; la nonlinéarité du schéma est alors facile à gérer!

L'estimation d'erreur vient de la convexité de \mathcal{J}_τ , qui peut s'exprimer, de manière équivalente, comme

$$\| \nabla^T u^T - \nabla^T v^T \|_{L^p}^p \leq (\mathcal{J}_\tau [u^T] - \mathcal{J}_\tau [v^T]) \cdot (u^T - v^T) \\ \text{(cas } p \geq 2) \leq [\mathcal{J}_\tau [v^T] - \mathcal{J}_\tau [u^T] - \nabla \mathcal{J}_\tau [u^T] \cdot (v^T - u^T)]$$

- La première forme utilise l'optimalité de u^c pour J_c et donne la preuve de l'estimation de base de tout à l'heure.
- La deuxième permet d'utiliser également l'optimalité de u_c pour J :

$$\begin{aligned} \|\nabla^T u^c - \nabla^T u_c^c\|_{L^p}^p &\leq J_c[u_c^c] - J_c[u^c] - \underbrace{\nabla J_c[u^c]}_{=0 \text{ car } u^c \text{ est optimale pour } J_c} (u_c^c - u^c) \leq \\ &\leq J_c[u_c^c] - J[\pi_T u_c] + J[\pi_T u_c] - J[u_c] + \underbrace{J[u_c] - J[\pi_T u^c]}_{=0 \text{ car } u \text{ est optimale pour } J} + J[\pi_T u^c] - J_c[u^c] \end{aligned}$$

où π_T est une interpolation affine/morceaux sur maillage \mathcal{T} telle que $\pi_T u_c = \pi_T u_c^c$.

Il faut alors contrôler trois termes, pour lesquelles on a:

- $|J[u_c] - J[\pi_T u_c]| \leq Ch^{2d}$ lorsque $u_c \in B_{\infty}^{1+d, p}$
- $|J_c[u_c^c] - J[\pi_T u_c^c]| \leq Ch^{2k}$ lorsque $u_c^c, u^c \in B_{\infty}^{1+d, p}$ - discret
- $|J_c[u^c] - J[\pi_T u^c]|$

On en déduit $\|\nabla^T u^c\|_{L^p}^p \leq Ch^{2d} \Rightarrow$ l'ordre $W^{1, p, T}$ en $h^{\frac{2d}{p}}$.

Oz, on connaît

Th (Y. Simon '78) Pour $p \geq 2$, Ω rectangle + cond. Dirichlet homogène, pour tout $f \in L^p(\Omega)$ la solution u de $-\Delta_p u = f$ est dans $B_{\infty}^{1+\frac{1}{p-1}, p}$.

(et cette régularité est optimale dans les Sobolev fractionnaires).

On montre alors

Proposition • si $u_c \in B_{\infty}^{1+d, p}$, alors $u_c^c \in B_{\infty}^{1+d, p}$ - discret

• toute solution u^c du schéma VF sur maillage cartésien uniforme pour $-\Delta_p u = f$, $f \in L^p$, est dans $B_{\infty}^{1+\frac{1}{p-1}, p}$ - discret.

(autrement dit, le schéma VF hérite de la régularité naturelle du problème continu!)

Corollaire Pour tout $f \in L^p$, sans hypothèse a priori sur u_c , on trouve en normes $W^{1, p, T}$ et L^p l'ordre de convergence en $\boxed{h^{\frac{2}{p(p-1)}}$ pour $p \geq 2$.

En numériquement, cet ordre se confirme comme optimal.

// //

Conclusion : on sait refaire la théorie "variationnelle" (fonction test u) du cas continu dans le cadre VF. La clé: la propriété de dualité div-grad.

? Problème ouvert: la théorie L^1 (pour le pb. d'évolution): fonct. test sig