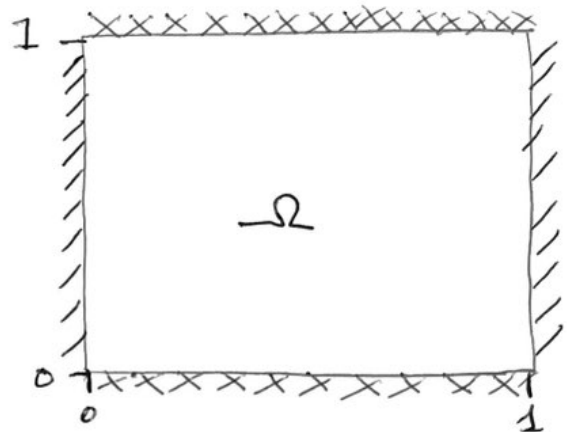# Numerics Exercise 3

i)

/// = Dirichlet boundary conditions on walls $\partial\Omega_D$

XXX = Neumann boundary conditions on walls $\partial\Omega_N$

Dirichlet BC: $u = 0$ on $\partial\Omega_D$

Neumann BC: $\frac{\partial u}{\partial y} = 0$ on $\partial\Omega_N$

Domain: $\Omega \in [0,1]^2$

Poisson equation: $-\nabla^2 u = f(x,y)$

Imposed force: $f(x,y) = 2\pi^2 \sin(\pi x)\cos(\pi y)$

Exact solution: $u_e(x,y) = \sin(\pi x)\cos(\pi y)$

To verify this:

$$\nabla^2 u_e = \partial_x^2 (\sin(\pi x)\cos(\pi y)) + \partial_y^2 (\sin(\pi x)\cos(\pi y))$$

$$= -\pi^2 u_e - \pi^2 u_e$$

$$= -2\pi^2 u_e$$

$$= -2\pi^2 \sin(\pi x)\cos(\pi y) \quad \text{as expected}$$

$$= -f$$

Ritz Gallerkin principle:

for the energy functional: $J[u]$
~~I[u]~~ $J[u] = \iint_\Omega \frac{1}{2}(\nabla u)^2 - f(x,y)u \, dx\, dy$

we find our solution when $\delta J[u] = J[u + \varepsilon\delta u] - J[u] = 0$

ie when the ~~potential~~ energy of the system is minimised

( continued )

$$\delta J[u] = J[u + \cancel{\delta} \varepsilon \delta u] - J[u]$$

$$= \iint_\Omega \frac{1}{2} \left( \nabla(u + \delta u) \right)^2 - f(x,y) \, (u + \delta u) \, dx \, dy$$

$$- \left( \iint_\Omega \frac{1}{2} (\nabla u)^2 - \cancel{f(x,y) u} \; dx \, dy \right)$$

$$= \iint_\Omega \frac{1}{2} \left( \cancel{(\nabla u)^2} + 2\nabla(\varepsilon \delta u) \cdot \nabla u + \varepsilon^2 (\nabla \delta u)^2 - \cancel{(\nabla u)^2} \right) - \varepsilon f \, \delta u \, dx \, dy$$

$$= \iint_\Omega \frac{\cancel{2} \varepsilon \nabla \delta u \cdot \nabla u}{2} + \frac{\varepsilon^2 (\nabla(\delta u))^2}{2} - \varepsilon f \delta u \, dx \, dy$$

$$\frac{\delta J[u]}{d\varepsilon} = \iint_\Omega \cancel{\frac{\nabla(\delta u) \nabla u}{2}} \nabla(\delta u) \cdot \nabla u + \varepsilon (\nabla \delta u)^2 - \cancel{\varepsilon} f \delta u \, dx \, dy$$

in $\displaystyle \lim_{\varepsilon \to 0}$

$$= \iint_\Omega \nabla(\delta u) \cdot \nabla u \; \cancel{\ast} - f \delta u \, dx \, dy$$

∴ we have obtained :

$$\iint_\Omega \nabla(\delta u) \cdot \nabla u - f(x,y) \, \delta u \, dx \, dy = 0$$

Continuous Ritz Gallerkin Formulation

Conditions on variable $\delta u$ : $\delta u(x,y)$ must be a smooth function within the domain $\Omega$ so that $\nabla(\delta u)$ is well defined

$\delta u$ must also be square integrable ie $\iint_\Omega \delta u(x,y) \, dx \, dy$ must give a finite solution

1 continued )

## Weak Formulation

$$-\nabla^2 u = f \quad \text{in} \quad \Omega \in [0,1]^2 \qquad\qquad u = 0 \text{ on } \partial\Omega_D$$

$$\partial\Omega = \partial\Omega_N \cup \partial\Omega_D \qquad\qquad \nabla u \cdot n = 0 \text{ on } \partial\Omega_N$$

$$-\langle \nabla^2 u, \nu \rangle = \langle f, $$ ~~(crossed out)~~

$$-\langle \nabla^2 u, \nu \rangle = \langle f, \nu \rangle \quad \text{where } \nu \text{ is chosen as an arbitrary test}$$
function that satisfies the boundary
conditions

$$\langle u, \nu \rangle = \iint_\Omega u(x,y)\, \nu(x,y)\, dx\, dy \qquad \text{Definition of the inner product}$$

$\therefore$ we get:

$$-\iint_\Omega (u_{xx} + u_{yy}) \,\nu(x,y)\, dx\, dy = \iint_\Omega f(x,y)\nu(x,y)\, dx\, dy$$

LHS: $-\iint_\Omega (u_{xx} + u_{yy}) \nu(x,y)\, dx\, dy$

we know: $\partial_x (u_x \nu) = u_{xx}\nu + u_x \nu_x$

$$\partial_x \partial_y (u_y \nu) = u_{yy}\nu + u_y \nu_y$$

rearrange $\hookrightarrow$

$$u_{yy}\nu = \partial_y (u_y \nu) - u_y \nu_y$$

$$u_{xx}\nu = \partial_x (u_x \nu) - u_x \nu_x$$

$\Big\}$ sub into integral

$$-\iint (u_{xx} + u_{yy})\nu\, dx\, dy = -\iint_\Omega \partial_y (u_y \nu) + \partial_x (u_x \nu)\, dx\, dy \quad ]-\text{Ⓐ}$$

$$+ \iint_\Omega (u_y \nu_y + u_x \nu_x)\, dx\, dy \quad ]-\text{Ⓑ}$$

1 continued)

to solve (A): apply green's theorem

$$\int_{\partial\Omega} P\,dx + Q\,dy = \iint_{\Omega} (Q_x - P_y)\,dy\,dx$$

$$-\iint_{\Omega}\left[\overbrace{\partial_y(u_yv)}^{-P_y} + \overbrace{\partial_x(u_xv)}^{+Q_x}\right]dx\,dy = -\int_{\partial\Omega} u_xv\,dy - u_yv\,dx$$

$$= \int_{\partial\Omega_N} u_yv\,dx - \int_{\partial\Omega_0} u_xv\,dy$$

we can choose $v$ (as it is arbitrary) such that

$$\int_{\partial\Omega_0} u_xv\,dy = 0$$

in addition: we know that $u_y = 0$ on $\partial\Omega_N$

∴ we get: (A) $= 0$

∴ our weak formulation is:

$$\iint_{\Omega}[u_yv_y + u_xv_x]\,dx\,dy = \iint_{\Omega} f(x,y)\,v(x,y)\,dx\,dy$$

in 2D: $u_yv_y + u_xv_x = \nabla u \cdot \nabla v$

∴ we can write this as:

$$\iint_{\Omega} \nabla u \cdot \nabla v - fv\,dx\,dy = 0$$

Continuous
Weak formulation

where $\Omega$: $0 \leq y \leq 1$, $f(x,y) = 2\pi^2\sin(\pi x)\cos(\pi y)$
$\quad\quad\quad 0 \leq x \leq 1$

1 continued )

We can see that the two must be equivalent by looking
Comparing the two integrals continuous Ritz Gallerkin
formulation and the weak continuous weak formulation.

$$\iint_{\Omega} \nabla(\delta u) \cdot \nabla u - f(x,y) \delta u \; dx \, dy = 0 \qquad \text{RG form}$$

$$\iint_{\Omega} \nabla v \cdot \nabla u - f(x,y) v \; dx \, dy = 0 \qquad \text{Weak form}$$

both must be :  · smooth so that $\nabla v / \nabla(\delta u)$ is well defined
  · square integrable so $\iint_{\Omega} u \, dx \, dy$ is finite

as $v$ is an arbitrary test function (as is $\delta u$) and the two forms share
the same integral formulation, it is clear that $v$ as and $\delta u$ are
equivalent. They are both test functions that perturb the system and
belong to a specific Hilbert space. specified by the boundary conditions.

2) $\underline{\text{Finding discrete Ritz Gallerkin Formulation}}$

from before :  $J[u] = \iint_{\Omega} \sum_i \sum_j$ 

NB: $u_i$ and $v_j$ are
interchangeable

$$J[u] = \iint_{\Omega} \tfrac{1}{2} |\nabla u|^2 - f(x,y) u \; dx \, dy$$

let :  $u \simeq u_n \approx \sum_{j=1}^{N_n} u_j \varphi_j(\bar{x})$

$u_j$ : coefficients

$\iint_{\Omega} \tfrac{1}{2} \left( \nabla(\sum_i u_i \varphi_i) \cdot \nabla(\sum_j u_j \varphi_j) \right) - f u_j$

$\varphi_j$ : basis function

$- \sum_j f u_j \varphi_j \; dx \, dy = 0$

$N_n$ : number of nodes where
solution is unknown

$\nabla(u_i \varphi_i) = \underbrace{\varphi_i \nabla u_i}_{=0} + u_i \nabla \varphi_i$

$N_T = N_n + N_k$

total # of nodes ← $N_T$ → # of known solution nodes (nodes on boundary $\geq \Omega_0$)

$= u_i \nabla \varphi_i$

2 continued)

$$\therefore \iint_{\Omega} \frac{1}{2} \left( \sum_i \sum_j u_i \varphi_i + u_i \nabla \varphi_i \cdot u_j \nabla \varphi_j \right) - \sum_j u_j f \varphi_j \; dxdy = 0$$

$$\therefore I[u] = \iint_{\Omega} \sum_i \sum_j \frac{u_i \nabla \varphi_i \cdot u_j \nabla \varphi_j}{2} - \sum_j u_j f \varphi_j \; dxdy \quad \text{is the discretised Ritz galerkin functional}$$

$$= \sum_i \sum_j \frac{A_{ij} u_i u_j}{2} - \sum_j b_j u_j$$

where we have defined: $\quad A_{ij} = \iint_{\Omega} \nabla \varphi_i \cdot \nabla \varphi_j \; dxdy$

it is convenient not to absorb factor of $1/2$ into $A_{ij}$

$$b_j = \iint_{\Omega} f \varphi_j \; dx \, dy$$

$$\delta I[u] = I[u + \varepsilon \delta u] - I[u]$$

Performing a variation on $\delta I[u]$

$$\text{have:} \quad \begin{matrix} u_i \longrightarrow u_i + \varepsilon \delta u_i \\ u_j \longrightarrow u_j + \varepsilon \delta u_j \end{matrix} \Big\} \text{ for } I[u + \varepsilon \delta u]$$

$$\therefore \delta I[u] = \sum_i \sum_j \frac{A_{ij}}{2} (u_i + \varepsilon \delta u_i)(u_j + \varepsilon \delta u_j) - \sum_j b_j (u_j + \varepsilon \delta u_j)$$

$$- \sum_i \sum_j \frac{A_{ij}}{2} u_i u_j + \sum_j b_j u_j$$

$$= \sum_i \sum_j \frac{A_{ij}}{2} \left( u_i u_j + u_i \varepsilon \delta u_j + \varepsilon \delta u_i u_j + \varepsilon^2 \delta u_j \delta u_j - u_i u_j \right) - \sum_j b_j \varepsilon \delta u_j$$

$$= \sum_i \sum_j \frac{A_{ij}}{2} \left( \varepsilon (u_i \delta u_j + \delta u_i u_j) + \varepsilon^2 \delta u_j \delta u_i \right) - \sum_j b_j \varepsilon \delta u_j$$

$$\frac{\delta I[u]}{d\varepsilon} = \sum_i \sum_j \frac{A_{ij}}{2} \left( u_i \delta u_i + \delta u_i u_j + 2\varepsilon \delta u_j \delta u_i \right) - \sum_j b_j \delta u_j$$

let $\varepsilon \to 0$

$$= \sum_j \sum_i \frac{A_{ij}}{2} (u_i \delta u_j + u_j \delta u_i) - b_j \delta u_j$$

2)

2 continued)

Since we defined: $A_{ij} = \iint_\Omega \nabla \varphi_i \cdot \nabla \varphi_j \, dx \, dy$

it is clear to see that $A_{ij} = A_{ji}$ (symmetric matrix)
∴ we can swap $i \longleftrightarrow j$ indices

∴ $\dfrac{\delta I[u]}{d\varepsilon} = \sum_j \sum_i \dfrac{A_{ij}}{2} \left( u_j \delta u_i + u_j \delta u_i \right) - b_j \delta u_j$

$= \sum_j \sum_i A_{ji} u_j \delta u_i - b_j \delta u_j = 0$

$\sum_j \sum_i \cancel{A_{ji} u_j \delta u_j} \, A_{ij} u_i \delta u_j - b_j \delta u_j = 0$

$\sum_j \sum_i \left( A_{ij} u_i - b_j \right) \delta u_j = 0$

writing in einstein notation:

$A_{ji}$
$\cancel{A_{ij}} u_i - b_j = 0$
$A_{ij} u_j - b_i = 0$

result when variational principle applied to discretised ritz galerkin functional

## Discretised Weak Formulation

Continuous weak form: $\iint_\Omega \nabla u \cdot \nabla v \, dx \, dy = \iint_\Omega f v \, dx \, dy$

let: $u(\underline{x}) \simeq u_n(\underline{x}) = \sum_{j=1}^{N_n} u_j \varphi_j$    where $\varphi_j$ is square integrable basis function ← a

$\sum_i \sum_j \iint_\Omega \nabla (u_j \varphi_j) \cdot \nabla v \cancel{\,dx} - f v \, dx \, dy = 0$ let $v = \varphi_i$ ie we choose our test function to be the same as our basis function

$\sum_i \sum_j \iint_\Omega \nabla (u_j \varphi_j) \cdot \nabla \varphi_i - f \varphi_i \, dx \, dy = 0$    this works as long as the basis functions are smooth and, square integrable and satisfy the boundary conditions

2 continued)

$$\sum_i \sum_j \iint_\Omega u_j \nabla \varphi_j \cdot \nabla \varphi_i - f \varphi_i \, dx \, dy = 0$$

let us define:     $A_{ij} = \iint_\Omega \nabla \varphi_i \cdot \nabla \varphi_j \, d\Omega$          $d\Omega = dx \, dy$

$$b_i = \iint_\Omega \cancel{f a} f \varphi_i \, d\Omega$$

$$\sum_i \sum_j A_{ij} u_j - b_i = 0$$

$\downarrow$

Einstein summation     $A_{ij} u_j - b_i = 0$
notation

This is the same expression that we obtained
from applying the variational principle to the
discretised Ritz Gallerkin Formulation,

NB: $A_{ij}$ is <u>not</u> a square matrix as $u$ is known on $\partial \Omega_D$ but not on $\partial \Omega_N$
we can rewrite it in the following way ~~to~~ to turn it into a square
matrix by seperating out the Neumann entries (ie those on $\partial \Omega_N$)

$$A_{ik} u_k - b_i = - \sum_{m - N_n + 1}^{N_T} A_{im} u_m$$

in the case of the poisson equation, I think the RHS disappears?

**3)** We can use a barycentric coordinate system that expresses the position of a point P within ~~a~~ a rectangular element ~~by the following~~ in the following way



$$(\vec{x}, \vec{y}) = \vec{v} = \frac{1}{4}(1 \pm \xi_1)(1 \pm \xi_2)\vec{x}_i$$

$$x_0\,\xi_1 + v_1\,\xi_2 + \dot{x}_2'\,\Big($$

you map an irregular triangle / squad to a unit square/ triangle in $\xi, \eta$ space

general triangle

$(1 - \xi_1)$

standard triangle

$x_1 \rightarrow (0,0) \qquad x_2 \rightarrow (1,0) \qquad x_3 \rightarrow 0,1$

basis functions :

$$\phi_1(\xi, \eta) = 1 - \xi - \eta$$

$$\phi_2(\xi, \eta) = \xi$$

$$\phi_3(\xi, \eta) = \eta$$

$$x = x_1 + (x_2 - x_1)\xi + (x_3 - x_1)\eta$$

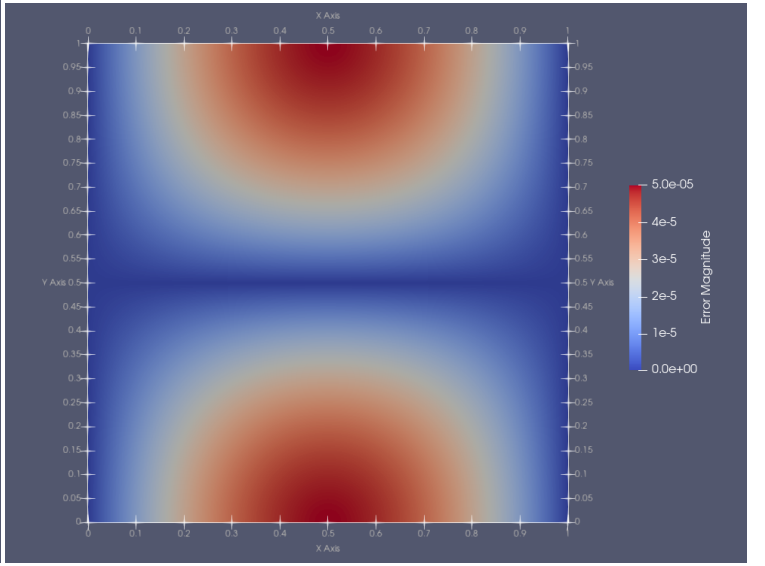$$y = y_1 + (y_2 - y_1)\xi + (y_3 - y_1)\eta$$

4)

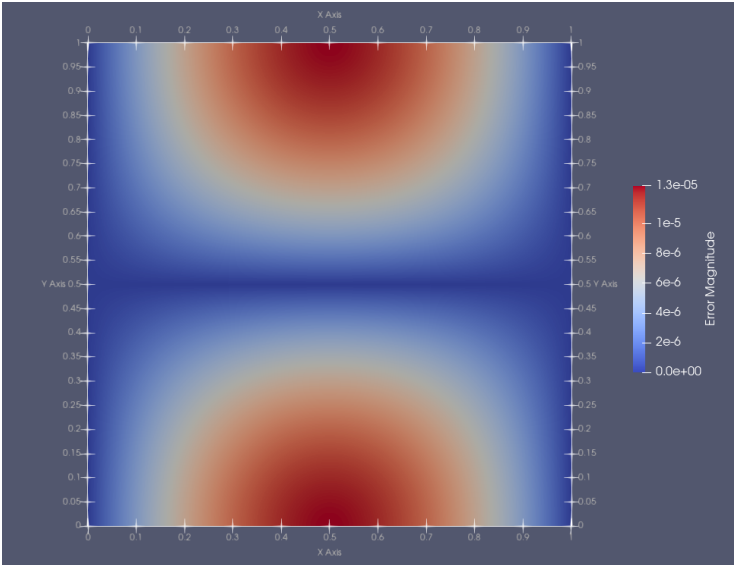

(a) $\{n, p\} = \{16, 1\}$

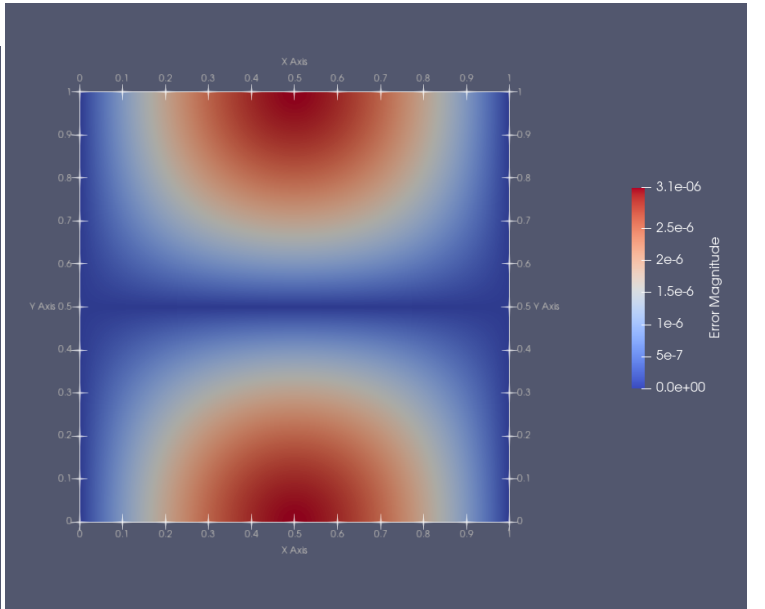(b) $\{n, p\} = \{32, 1\}$

(c) $\{n, p\} = \{64, 1\}$

(d) $\{n, p\} = \{128, 1\}$

(e) $\{n, p\} = \{256, 1\}$

(f) $\{n, p\} = \{512, 1\}$

Figure 1: A figure showing how $\sigma = |u_h - u_e|$ changes as $n$ increases, where $u_h$ is the numerical approximation, $u_e$ is the exact solution and $n$ is the number of nodes within each dimension. The scale on each colour map remains the same, highlighting that the error associated with discretising the Poisson equation using a 1st order FEM is minimised when the resolution of the mesh increases.

(a) $\{n, p\} = \{16, 1\}$

(b) $\{n, p\} = \{32, 1\}$

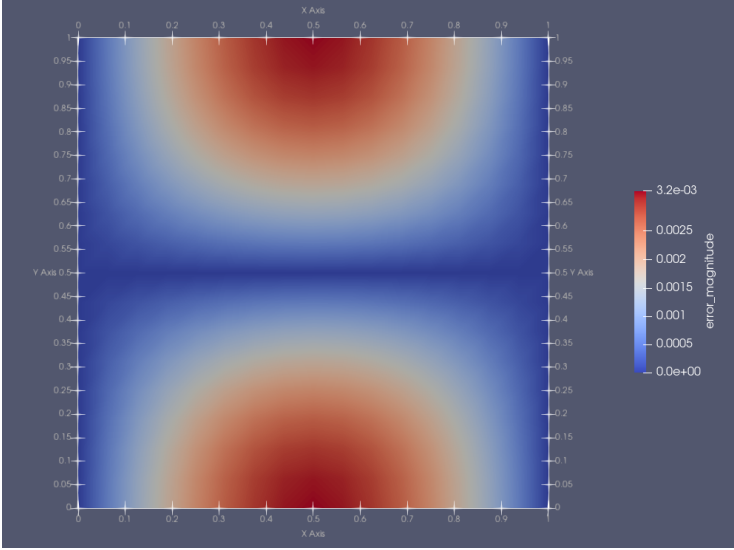(c) $\{n, p\} = \{64, 1\}$

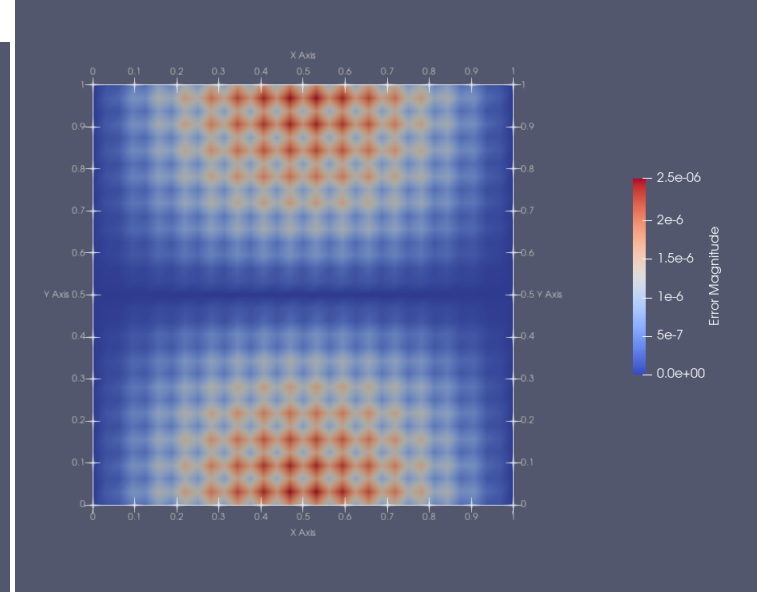(d) $\{n, p\} = \{128, 1\}$

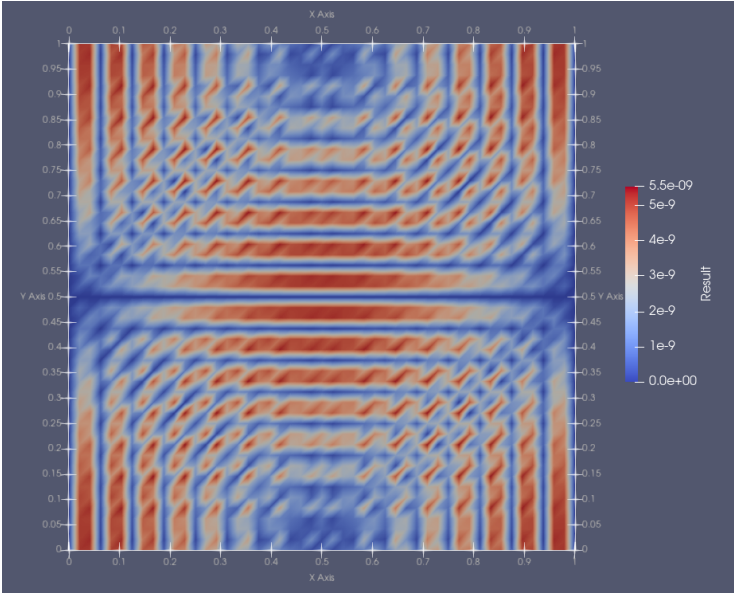(e) $\{n, p\} = \{256, 1\}$

(f) $\{n, p\} = \{512, 1\}$

Figure 2: The same results as in Figure 1, but now each colourbar has been rescaled to show the extremities of $|u_h - u_e|$ for the $n$ refinement. By looking at the maxima of each scale, we observe that doubling $n$ decreases the error by a factor of 4. Therefore, we can deduce that $\sigma \propto n^{-2}$ in 2D or (more generally) $|u_h - u_e| \propto n^{-d}$ where $d$ is the number of dimensions.
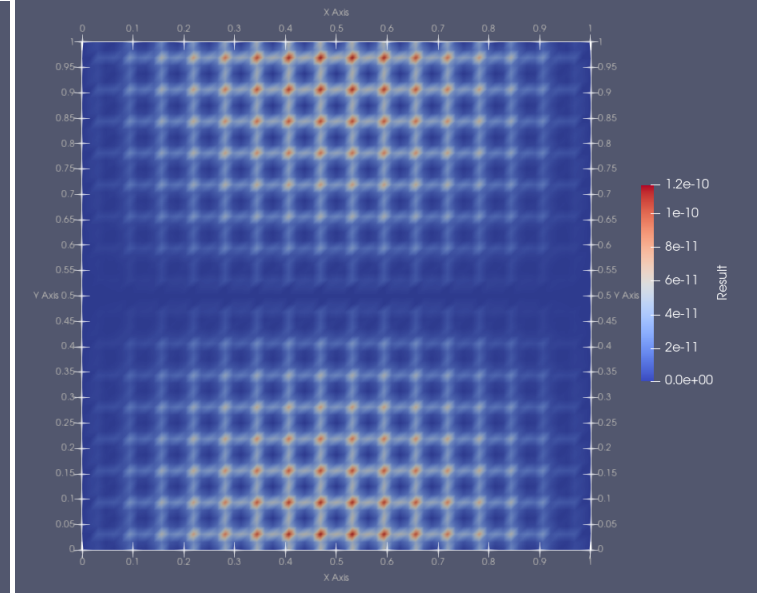
(a) $\{n, p\} = \{16, 1\}$
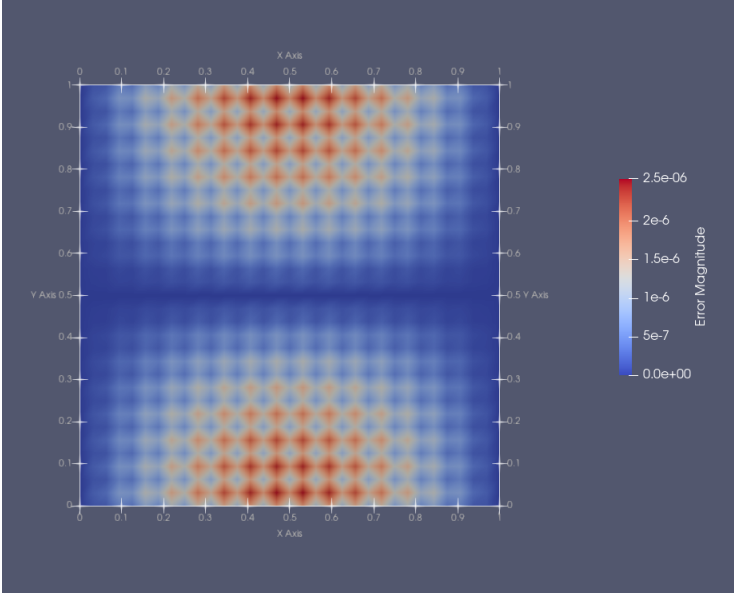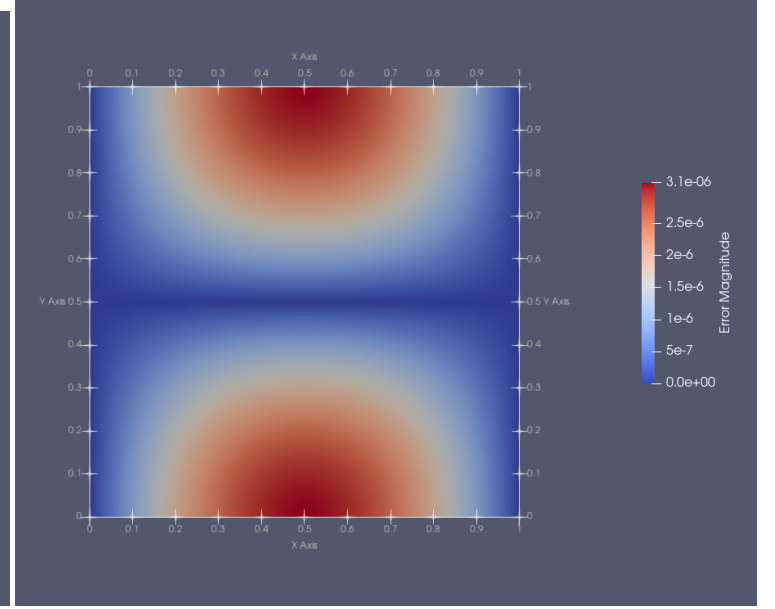
(b) $\{n, p\} = \{16, 2\}$

(c) $\{n, p\} = \{16, 3\}$
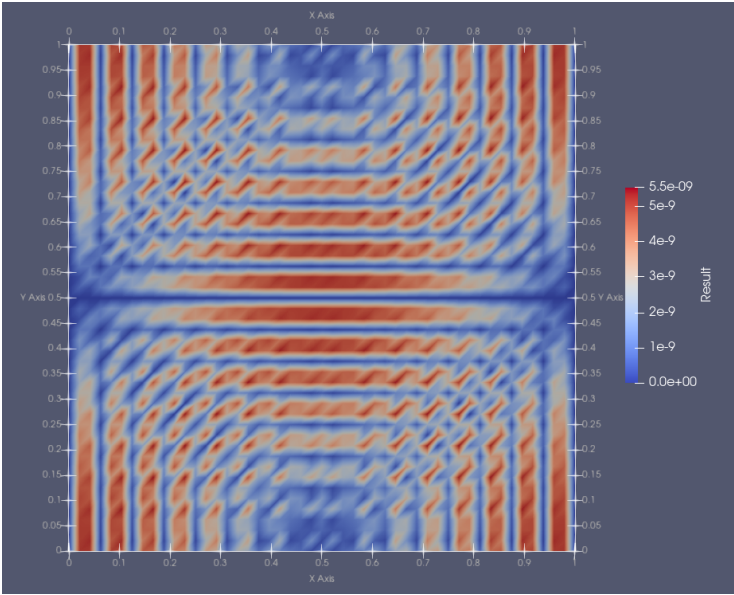
(d) $\{n, p\} = \{16, 1\}$

Figure 3: A $p$ refinement showing the effect of increasing the order $p$ of the Finite Element Approximation when $n = 16$. Notably, increasing $p$ greatly reduces the numerical error $\sigma$ but does not change the spatial resolution.
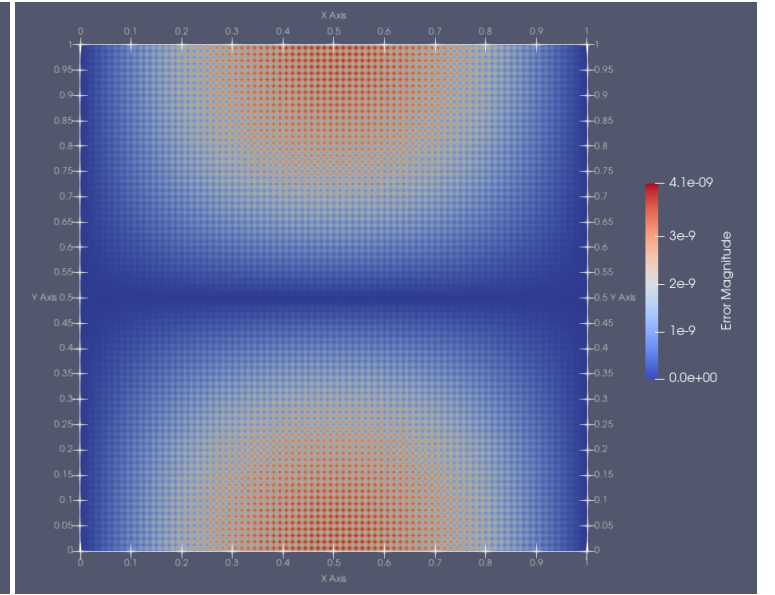
(a) $\{n, p\} = \{16, 2\}$
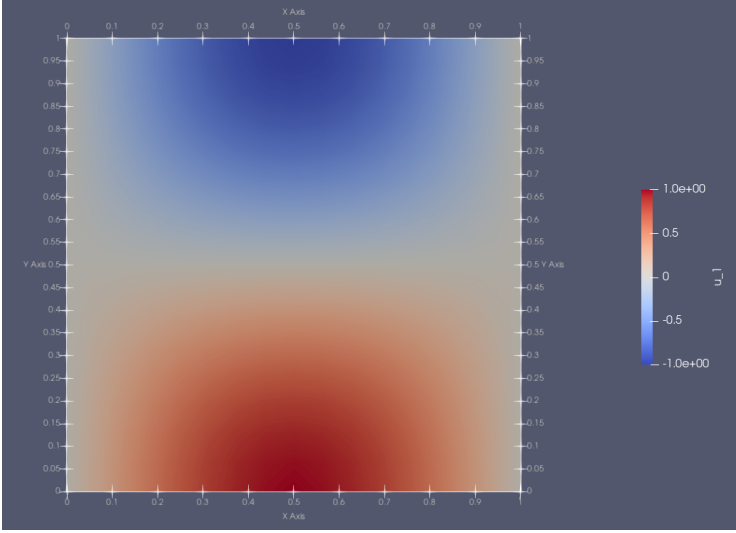


(b) $\{n, p\} = \{256, 1\}$
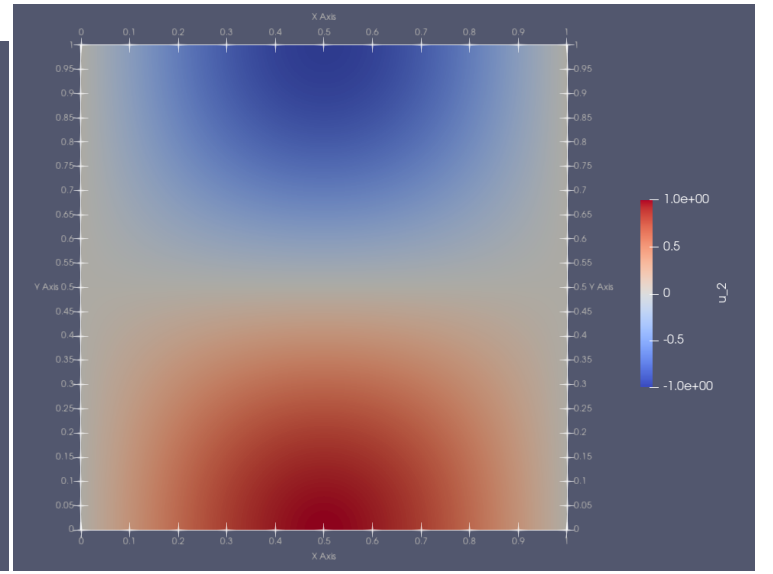


(c) $\{n, p\} = \{16, 3\}$



(d) $\{n, p\} = \{80, 2\}$

Figure 4: Two different $\{n, p\}$ pairs: $[(a), (b)]$ and $[(c), (d)]$. This highlights that increasing $p$ is a very effective way of decreasing the error but does not increase the resolution of the mesh.

(a) $\{n, p\} = \{16, 1\}$

(b) $\{n, p\} = \{256, 1\}$

Figure 5: The numerical approximation to the Poisson equation $u_h(x, y)$ plotted over the domain for two different values of $n$. Since this solution is quite simple, the low spatial resolution of $(a)$ does not significantly affect the results that are obtained within our simulation.

Looking at these results, we can make the following deductions:

- Increasing $n$ (the number of nodes in each dimension) increases the resolution of the graph (which is most easily observed in Figure 4) and the error $\sigma = |u_h - u_e| \propto n^{-2}$ in 2D. We know that $n = \frac{1}{h}$ where $h$ is the side length of each cell within a 2D square mesh. Therefore, an $n$ refinement is the inverse of a $h$ refinement in this case. From this we note that the error scales as $O(h^2)$ when $p = 1$ and the mesh is 2 dimensional. More generally, we can deduce that $\sigma \propto h^{-d}$ for the domain $\{0,1\}^d$ spanned by cells with side length $h$ and volume $h^d$.

- $p$ refinement is far more effective than $h$ refinement for minimising the error but it does not increase the spatial resolution of the mesh. Therefore, it is desirable to use a higher value of $p$ to minimise the error but there is a trade-off as you are only minimising the error at a point. Therefore, if $h$ is too low, you do not have sufficient spatial resolution to accurately resolve the flow and you miss important flow features.

- To expand upon this, the convergence scaling which relates $p$ and $h$ is $O(h^p)$ which tells you how the error decreases as $h$ increases. This suggests that there are different scaling laws for different values of $p$, which makes sense intuitively as a higher order polynomial approximation samples each cell more times (it has a higher degree of freedom). This scaling causes certain $\{h, p\}$ pairs to have a similar amount of numerical error (as observed in Figure 4).

- By comparing Figure 5 with the other figures, we observe that they share the same general profile and symmetry. In particular, the error is maximised where $|u_h(x, y)|$ is large, suggesting that an irregular mesh that places additional cells at these regions could help to resolve the flow more effectively (as an alternative to using a regular mesh and performing refinements of $h$ or $p$).

In conclusion, $h$ refinement is primarily used as a tool to increase mesh resolution, while $p$ refinement is used to reduce numerical error. Using a mixture of the two is usually the best approach for achieving accurate and highly resolved simulations. Since the two different methods of calculating $u_h$ are mathematically equivalent, we can use either method within this question. I personally opted to use the second method (Ritz-Gallerkin principle) as I think it is more elegant, but you could also use the weak formulation and the results should be equivalent (with any error being attributed to machine precision).

5) I have added some comments at the start of the code that I uploaded to Github in an attempt to translate what Firedrake is doing at each step.

6) I may go back and update my FV solutions over the next week or so ( I am still currently deciding whether I want to do this or not).