

# Rapport du projet de la Méthode de base réduite pour la résolution d'EDPs dépendantes de paramètre

*Jincheng KE*

Une présentation du projet sur les résultats obtenus.

Mars 2024

# Table des matières

<b>1</b>	<b>Mise en forme du problème paramétrique</b>	<b>2</b>
1.1	Solution numérique des différents paramètres . . . . .	2
1.2	La décomposition affine . . . . .	2
1.3	Validation numérique . . . . .	3
1.4	Solveur base réduite . . . . .	6
1.4.1	La décomposition affine . . . . .	6
<b>2</b>	<b>Réduction par l'approche Greedy</b>	<b>8</b>
2.1	Calcul du résidu base réduite . . . . .	8
2.2	L'algorithme Greedy . . . . .	10
2.3	Certification . . . . .	11
<b>3</b>	<b>Application</b>	<b>13</b>
3.1	Offline et Online . . . . .	13
3.2	Estimation à posteriori . . . . .	13
3.3	Résultats . . . . .	13

# Chapitre 1

## Mise en forme du problème paramétrique

### 1.1 Solution numérique des différents paramètres

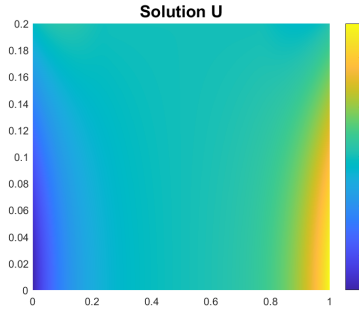


FIGURE 1.1 –  $\kappa_1 = 1, \kappa_2 = 1$

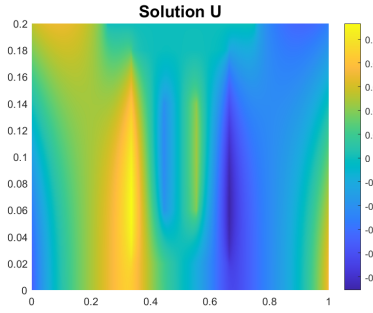


FIGURE 1.2 –  $\kappa_1 = 0.1, \kappa_2 = 1$

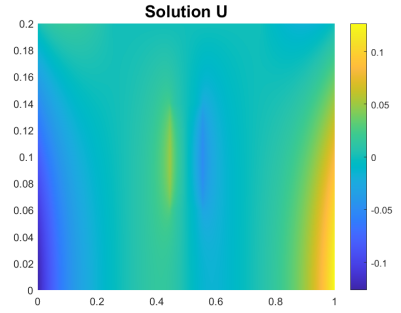


FIGURE 1.3 –  $\kappa_1 = 1, \kappa_2 = 0.1$

### 1.2 La décomposition affine

On a la matrice  $\mathbf{A}$  peut être décomposée par

$$\mathbf{A} = \mathbf{A}_0 + \kappa_1 \mathbf{A}_1 + \kappa_2 \mathbf{A}_2$$

avec

$$\mathbf{A}_q = \mathbf{K} \cdot \mathbb{1}_{\{\text{RefTriangles}=q\}}, \forall q = 0, 1, 2$$

Et comme on a

$$\mathbf{L} = \mathbf{M} \cdot \mathbf{F} + \mathbf{S}_1 \cdot \mathbf{G}_1 + \mathbf{S}_2 \cdot \mathbf{G}_2 + \mathbf{S}_4 \cdot \mathbf{G}_4 + \mathbf{S}_5 \cdot \mathbf{G}_5 - \mathbf{A} \cdot \Phi_r$$

Il y a que la matrice  $\mathbf{A}$  qui dépend du paramètre  $\kappa$ , on peut donc décomposer le vecteur  $\mathbf{L}$  en décomposant la matrice  $\mathbf{A}$ , on obtient la formule suivante :

$$\mathbf{L} = \mathbf{L}_0 + \kappa_1 \mathbf{L}_1 + \kappa_2 \mathbf{L}_2$$

avec

$$\mathbf{L}_0 = \mathbf{M} \cdot \mathbf{F} + \mathbf{S}_1 \cdot \mathbf{G}_1 + \mathbf{S}_2 \cdot \mathbf{G}_2 + \mathbf{S}_4 \cdot \mathbf{G}_4 + \mathbf{S}_5 \cdot \mathbf{G}_5 - \mathbf{A}_0 \cdot \Phi_r$$

et

$$\mathbf{L}_q = -\mathbf{A}_q \cdot \Phi_r, \forall q = 1, 2$$

cf. code `FE_assemblages.m`

### 1.3 Validation numérique

Pour vérifier la décomposition affine pour la matrice  $\mathbf{A}$  et pour le vecteur  $\mathbf{L}$ , on compare l'erreur absolue et l'erreur relative entre la solution de référence  $\mathbf{U}_{\text{ref}}$  de l'équation  $\mathbf{A}_{\text{ref}} \cdot \mathbf{U}_{\text{ref}} = \mathbf{L}_{\text{ref}}$  et la solution  $\mathbf{U}$  de l'équation avec la composition affine  $(\mathbf{A}_0 + \kappa_1 \mathbf{A}_1 + \kappa_2 \mathbf{A}_2) \cdot \mathbf{U} = \mathbf{L}_0 + \kappa_1 \mathbf{L}_1 + \kappa_2 \mathbf{L}_2$ . Pour  $\kappa = (0.1, 1)$ , l'erreur absolue est environ  $4.658e - 14$ , l'erreur relative est environ  $1.123e - 14$ . Pour  $\kappa = (1, 0.1)$ , l'erreur absolue est environ  $1.408e - 14$ , l'erreur relative est environ  $5.171e - 15$ .

L'erreur est assez petite pour pouvoir valider cet algorithme de la décomposition affine.

# Réduction par l'approche POD

## Construction d'une base réduite par POD

L'implémentation cf. MAIN\_RBPOD.m

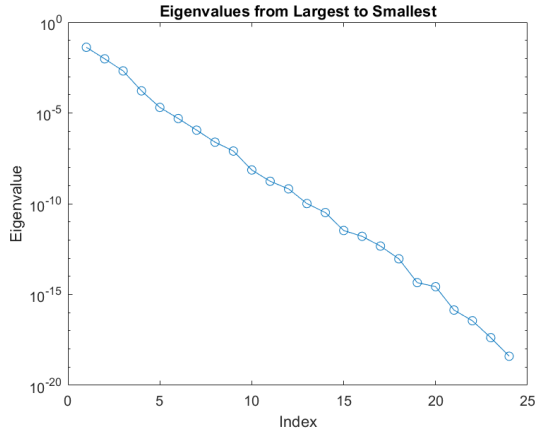


FIGURE 1.4 – type cartésien avec  $n_1 = n_2 = 5$

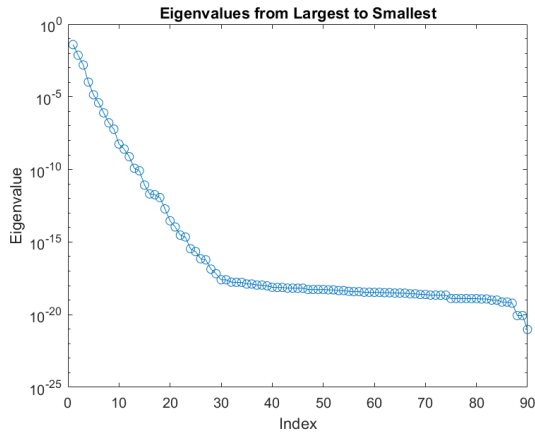


FIGURE 1.5 – type cartésien avec  $n_1 = n_2 = 12$

Pour le plan d'expérience cartésien avec  $n_1 = n_2 = 5$ , on a observé que la vitesse de la décroissance des valeurs propres est uniforme et rapide, il y a toujours une tendance décroissante avec la même vitesse. De plus, la plus petite valeur propre est environ  $1e-19$ . Comme il y a une tendance décroissante avec la même vitesse, et la plus petite valeur propre ici est déjà assez petit, il semble que le problème est réductible. Mais si on veut une précision plus précise, il semble que c'est mieux de prendre plus d'échantillons.

Pour le plan d'expérience cartésien avec  $n_1 = n_2 = 12$ , on a observé que la tendance de la vitesse de la décroissance uniforme et rapide se poursuit jusqu'à  $n = 30$ . Quand  $n \geq 30$ , la vitesse devient vraiment petite, les valeurs propres se stabilisent, et au final, la plus petite valeur propre est environ  $1e-20$ . Comme les valeurs propres se stabilisent quand  $n \geq 30$ , il semble que le nombre de RB de POD ne permettra plus d'améliorer la précision quand il est grand.

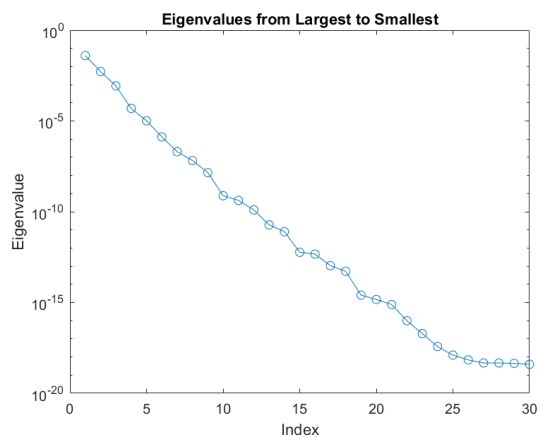


FIGURE 1.6 – type 'random' avec  $n = 35$

Pour le plan d'expérience aléatoire avec  $n = 35$ , on a observé que globalement, il y a la même tendance de la décroissance rapide des valeurs propres avec le cas de l'expérience cartésien en prenant  $n_1 = n_2 = 5$ . La plus petit valeur propre est aussi environ  $1e-19$ . Mais localement, on peut voir que la vitesse de décroissance se varie. Cette variation de la vitesse locale est produite par la sensibilité à l'échantillonnage.

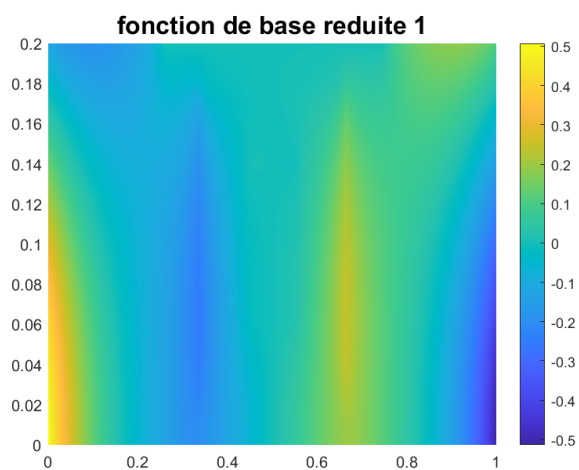


FIGURE 1.7 –  $n_1 = n_2 = 12$  base 1

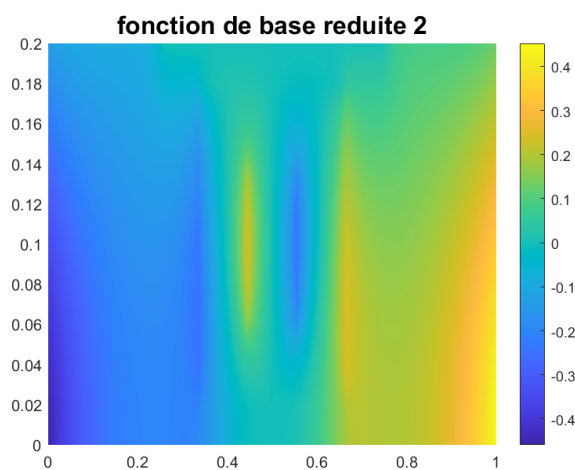


FIGURE 1.8 –  $n_1 = n_2 = 12$  base 2

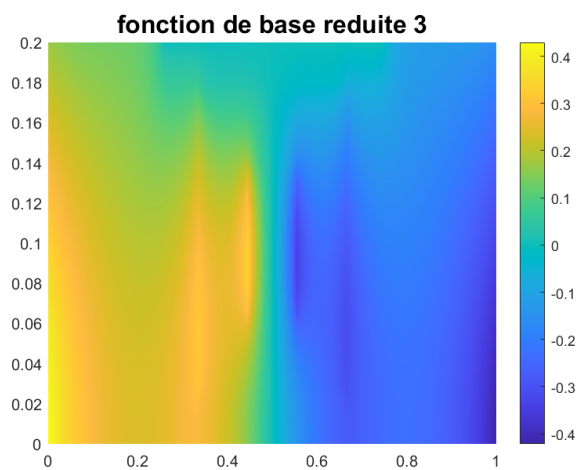


FIGURE 1.9 –  $n_1 = n_2 = 12$  base 3

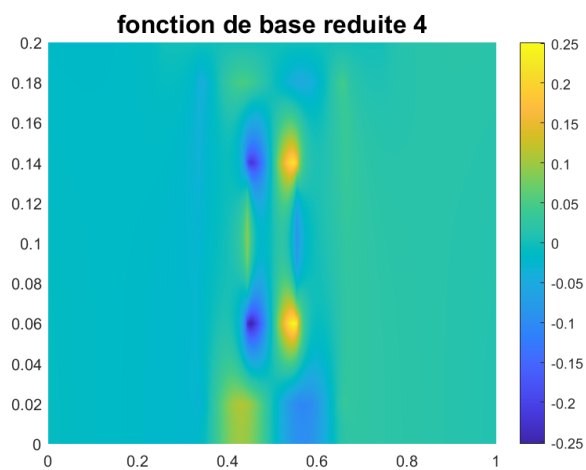
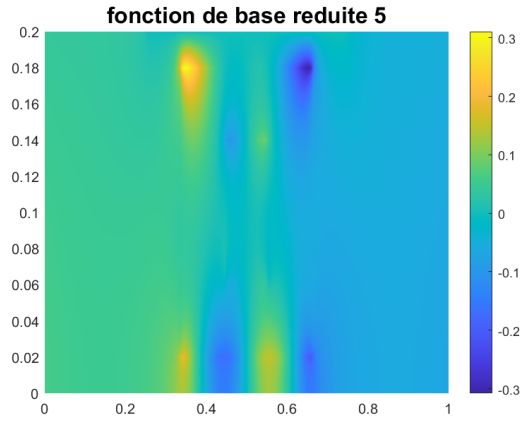


FIGURE 1.10 –  $n_1 = n_2 = 12$  base 4



Les images FIGURE 1.7, 1.8, 1.9, 1.10 et 1.11 sont les 5 premières RB de POD pour le plan d'expérience cartésien avec  $n_1 = n_2 = 12$ .  
L'implémentation cf. `MAIN_RBPOD.m`

FIGURE 1.11 –  $n_1 = n_2 = 12$  base 5

## 1.4 Solveur base reduite

### 1.4.1 La décomposition affine

On a

$$\mathbf{A}^{\text{rb}} = \mathbf{P}^T \cdot \mathbf{A} \cdot \mathbf{P}$$

et

$$\mathbf{L}^{\text{rb}} = \mathbf{P}^T \cdot \mathbf{L}$$

En utilisant la décomposition affine de  $\mathbf{A}$  et  $\mathbf{L}$ , on aura la décomposition suivante :

$$\mathbf{A}^{\text{rb}} = \mathbf{A}_0^{\text{rb}} + \kappa_1 \mathbf{A}_1^{\text{rb}} + \kappa_2 \mathbf{A}_2^{\text{rb}}$$

$$\mathbf{L}^{\text{rb}} = \mathbf{L}_0^{\text{rb}} + \kappa_1 \mathbf{L}_1^{\text{rb}} + \kappa_2 \mathbf{L}_2^{\text{rb}}$$

avec

$$\mathbf{A}_q^{\text{rb}} = \mathbf{P}^T \cdot \mathbf{A}_q \cdot \mathbf{P}$$

$$\mathbf{L}_q^{\text{rb}} = \mathbf{P}^T \cdot \mathbf{L}_q, \forall q = 0, 1, 2$$

L'implémentation cf. `RB_reduced_decomp.m`

### Complexité

Pour `RB_solve.m`, comme on a déjà décomposé affinement, alors toutes les opérations sont sur la matrice de taille  $N \times N$  ( $N \approx 30$ ) et le vecteur de taille  $N$ , donc la complexité est  $\mathcal{O}(N^3)$ .

### Validation numérique

L'implémentation cf. `MAIN_RBSOLVER.m`

On va tracer les erreurs d'approximation base réduite maximales (max error) et moyennes (mean error) en fonction de  $N$  avec  $N = 4, 8, 12, 16, 20, 24, 28, 32$  et on va tracer également erreurs de l'approximation RB optimale (error2).

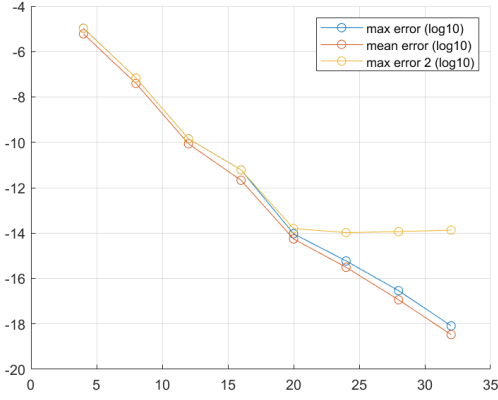


FIGURE 1.12 – type cartésien avec  $n_1 = n_2 = 12$

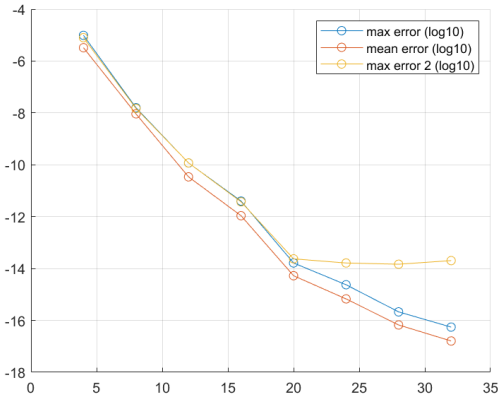


FIGURE 1.13 – type 'random' avec  $n = 35$

Pour le plan d'expérience cartésien avec  $n_1 = n_2 = 12$ , on a observé que avec le nombre de RB qui devient de plus en plus grand, l'erreur maximale et l'erreur moyenne devient aussi de plus en plus petit, et comme la courbe est une ligne à l'échelle de semi-log, on a la vitesse de convergence est exponentielle. Pourtant, pour l'erreur de l'approximation RB optimale, elle décroît aussi vite que l'erreur maximale / moyenne au début, mais jusqu'à  $n = 20$ . Quand  $n \geq 20$ , elle se stabilise à valeur de environ  $10e-14$ . Cela implique on a atteint à la limite de cette échantillonnage. Cette valeur est la norme de projection de  $\mathbf{U}$  dans l'espace  $\mathbf{V}_{\text{NTrain}}^\perp$  (NTrain = 144 ici).

Pour le plan d'expérience aléatoire avec  $n = 35$ , on a observé que avec le nombre de RB qui devient de plus en plus grand, l'erreur maximale et l'erreur moyenne devient aussi de plus en plus petit, et comme la courbe est une ligne à l'échelle de semi-log, on a la vitesse de convergence est exponentielle. Pourtant, pour l'erreur de l'approximation RB optimale, elle décroît aussi vite que l'erreur maximale / moyenne au début, mais jusqu'à  $n = 20$ . Quand  $n \geq 20$ , elle se stabilise à valeur de environ  $10e-14$ . On peut expliquer ce phénomène avec la même raison pour le plan d'expérience cartésien.



## Chapitre 2

# Réduction par l'approche Greedy

### 2.1 Calcul du résidu base réduite

Le coût mémoire de **Respart\_1l** est  $Q_l^2$ .

Le coût mémoire de **Respart\_1a** est  $Q_l \cdot Q_a \cdot N$ .

Le coût mémoire de **Respart\_aa** est  $Q_a^2 \cdot N^2$ .

Pour **Respart\_1l**, en utilisant la symétrie de la matrice  $B^{-1}$ , on peut optimiser le coût mémoire à  $0.5 \cdot Q_l^2$ .

De même, pour **Respart\_aa**, en utilisant la symétrie de  $B^{-1}$  et de  $P$ , on peut optimiser le coût mémoire à  $0.5 \cdot Q_a^2 \cdot N^2$ .

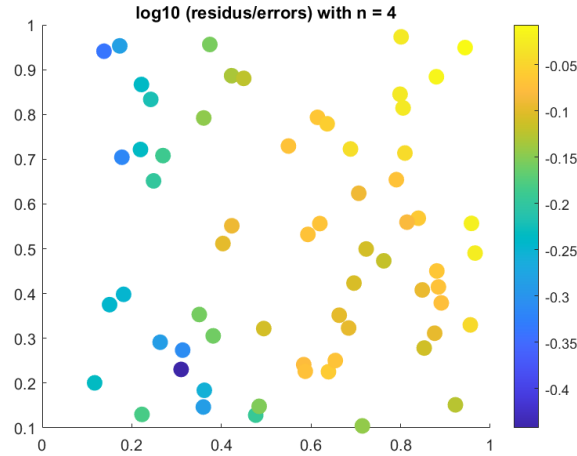


FIGURE 2.1 –  $\log_{10} \frac{\text{résidus}}{\text{errors}}$  avec  $n = 4$

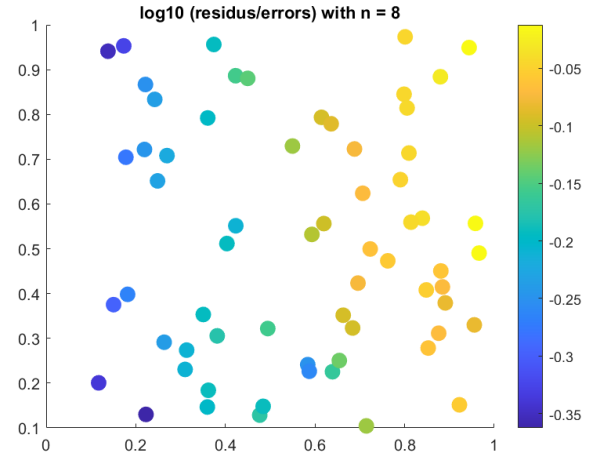


FIGURE 2.2 –  $\log_{10} \frac{\text{résidus}}{\text{errors}}$  avec  $n = 8$

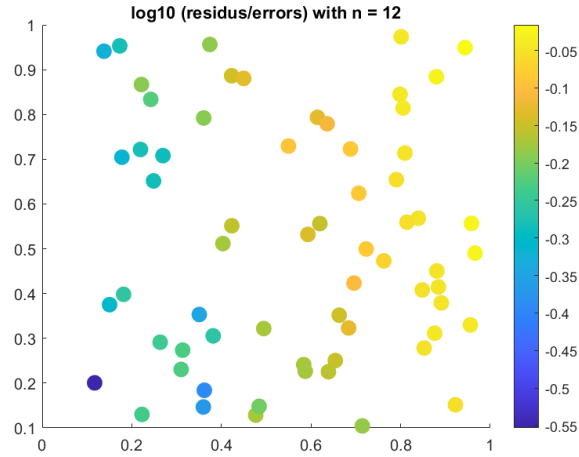


FIGURE 2.3 –  $\log_{10} \frac{\text{résidus}}{\text{errors}}$  avec  $n = 12$

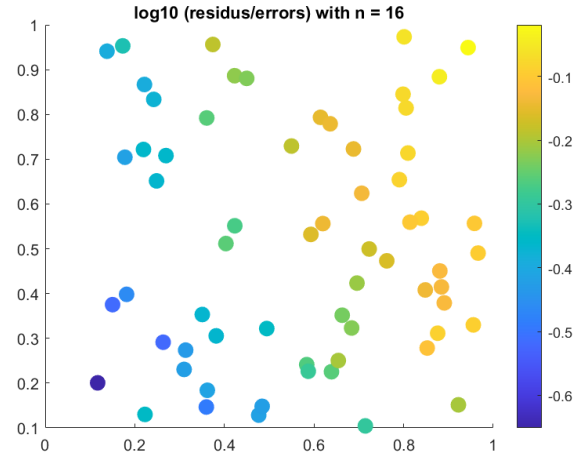


FIGURE 2.4 –  $\log_{10} \frac{\text{résidus}}{\text{errors}}$  avec  $n = 16$

Les images FIGURE 2.1, 2.2, 2.3 et 2.4 sont le rapport  $\text{résidus/errors}$  lors de nombre de POD  $N = 4, 8, 12, 16$ .

On souhaite que ce rapport soit proche de 1 (sa valeur logarithmique proche de 0). On a observé que la plus part de points sont entre 0 et -0.15 pour  $n = 4, 8$ . Pour  $n = 12$ , environ la moitié de points sont entre 0 et -0.15. Pour  $n = 16$ , il reste environ moins d'un tiers de points qui sont entre 0 et -0.15. Cela implique que la qualité de l'algorithme de résidu devient de moins en moins bien. Ce problème est dû à la précision numérique de cet algorithme qui n'est pas suffisamment précise.

## 2.2 L'algorithme Greedy

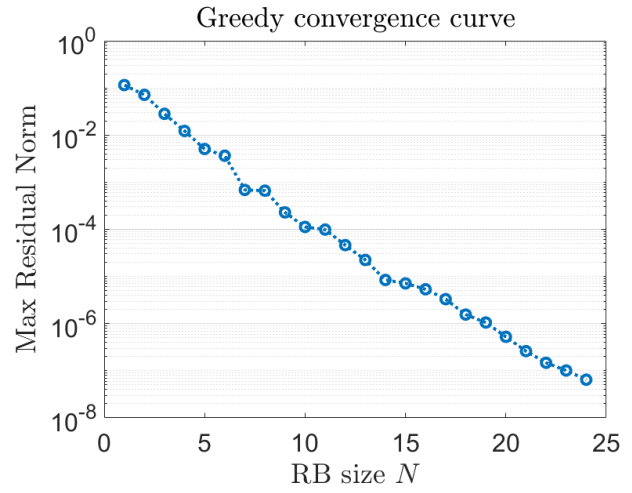
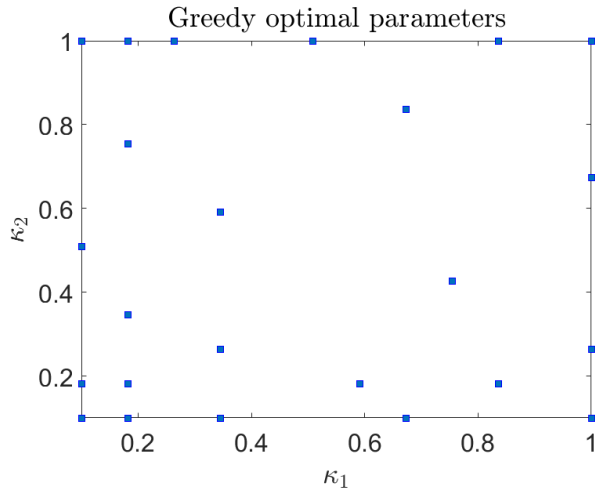


FIGURE 2.5 – algorithme classique avec  $\text{tol} = 10e - 7$

FIGURE 2.6 – algorithme classique avec  $\text{tol} = 10e - 7$

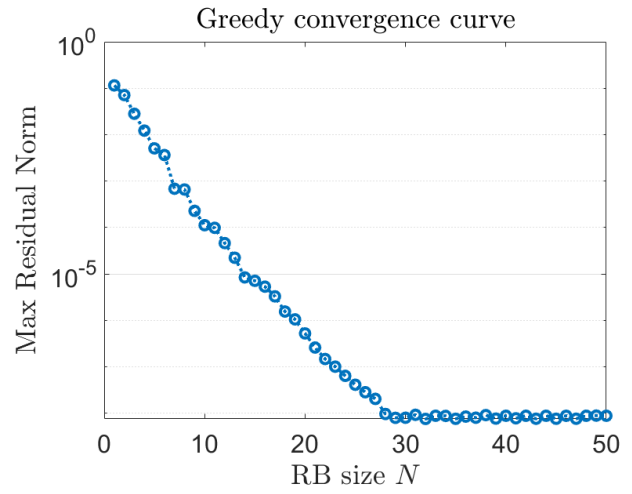
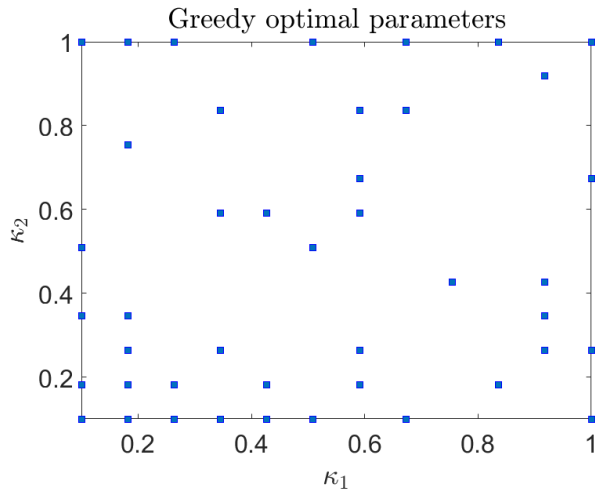


FIGURE 2.7 – algorithme classique avec  $\text{tol} = 10e - 10$

FIGURE 2.8 – algorithme classique avec  $\text{tol} = 10e - 10$

Les images FIGURE 2.5 et 2.6 sont le résultats pour la méthode classique avec une tolérance à  $10e - 7$ . On peut observer que la méthode fonctionne bien, le résidu décroît très vite.

Les images FIGURE 2.7 et 2.8 sont le résultats pour la méthode classique avec une tolérance à  $10e - 10$ . On peut observer que la méthode ne fonctionne plus, le résidu décroît très vite jusqu'à  $10e - 8$  et qui se stabilise. Donc avec la méthode classique, on ne peut pas avoir une précision de niveau  $10e - 10$ .

De plus, le nombre de résidu négatif peut également expliquer le problème de précision numérique. Pour tolérance à  $10e - 7$ , tous les résidus sont positifs. Mais pour tolérance à  $10e - 10$ , par exemple pour itération 50, il y a 20 résidus qui sont strictement négatifs. Or le résidu devrait être positif, on a trouvé alors le problème de précision numérique.

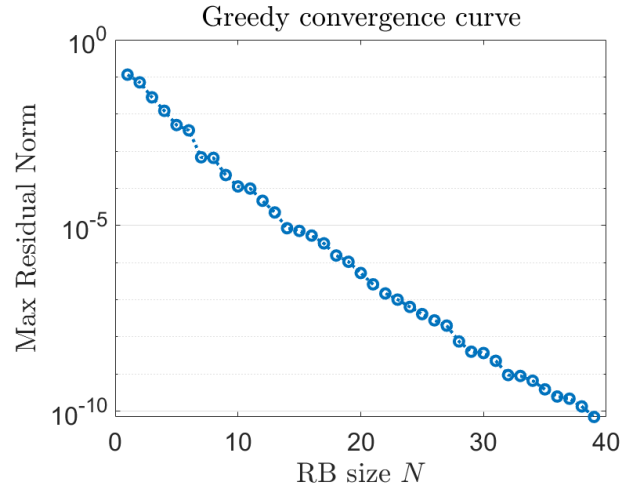
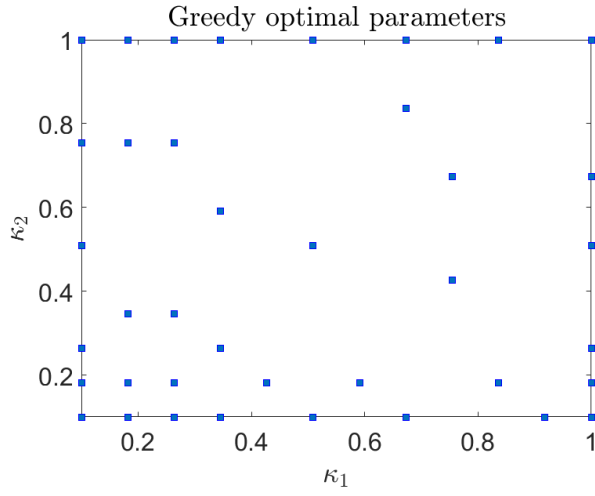


FIGURE 2.9 – algorithme stabilisé avec  $\text{tol} = 10e - 10$  FIGURE 2.10 – algorithme stabilisé avec  $\text{tol} = 10e - 10$   
 Les images FIGURE 2.9 et 2.10 sont le résultats pour la méthode stabilisé avec une tolérance à  $10e - 10$ .  
 On peut observer que la méthode fonctionne bien, le résidu décroît très vite.

De plus, on a observé aussi que tous les résidus sont positifs. Cela implique que la méthode stabilisée améliore le problème de précision numérique.

## 2.3 Certification

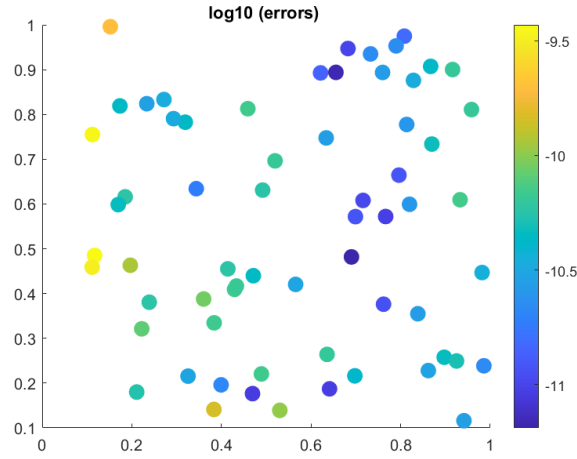


FIGURE 2.11 – errors avec RB Greedy

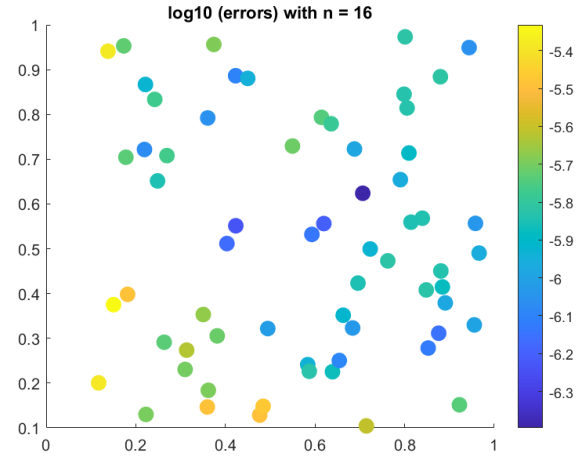
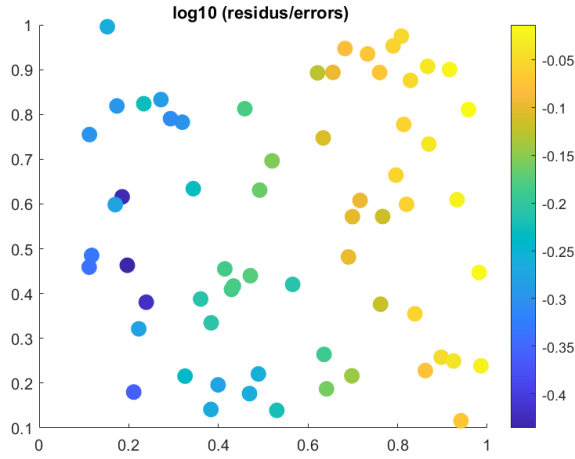


FIGURE 2.12 – errors avec RB POD  $N = 16$



Les images FIGURE 2.11 et 2.13 sont le résultats pour la base réduite de Greedy avec la méthode stabilisé avec une tolérance à  $10e - 10$ . En comparant avec FIGURE 2.12, on peut observer que la précision pour la base réduite de Greedy est beaucoup plus précise que celle pour la base POD lors de  $N = 16$ .

De plus, on a observé aussi que le rapport **residus/errors** pour la base réduite de Greedy avec la méthode stabilisé a une bonne qualité. Il y a environ plus de moitié de points qui sont entre 0 et -0.15.

FIGURE 2.13 –  $\log_{10} \frac{\text{résidus}}{\text{errors}}$  avec RB Greedy

# Chapitre 3

## Application

### 3.1 Offline et Online

Pour la phase **Offline**, On a  $j^{\text{rb}} = P^T \cdot j$ .

Pour la phase **Online**, on a  $Q^{\text{rb}}(\mu) = j^{\text{rb}T} \cdot x(\mu)$ .

La complexité pour obtenir  $x(\mu)$  est  $\mathcal{O}(N^3)$ .

La complexité pour le produit  $j^{\text{rb}T} \cdot x(\mu)$  est  $\mathcal{O}(N)$ .

### 3.2 Estimation à posteriori

Pour obtenir estimation  $|Q^{\text{rb}}(\mu) - Q^{\mathcal{N}}(\mu)| \leq \delta_N^Q(\mu)$ , on a  $|Q^{\text{rb}}(\mu) - Q^{\mathcal{N}}(\mu)| = |j^T \cdot (P \cdot x(\mu) - u^{\mathcal{N}}(\mu))|$ .

Comme on a  $\|j\| = \sup_{\|v\| \neq 0} \frac{|j^T \cdot v|}{\|v\|}$ , on a alors  $|j^T \cdot (P \cdot x(\mu) - u^{\mathcal{N}}(\mu))| \leq \|j\| \cdot \|(P \cdot x(\mu) - u^{\mathcal{N}}(\mu))\|$ .

De plus, on a déjà  $\|(P \cdot x(\mu) - u^{\mathcal{N}}(\mu))\| \leq \Delta_N(\mu)$ , où  $\Delta_N(\mu)$  est l'estimation à posteriori par le résidu avec un algorithme stabilisé.

Finalement, on peut poser  $\delta_N^Q(\mu) = \|j\| \cdot \Delta_N(\mu)$ .

### 3.3 Résultats

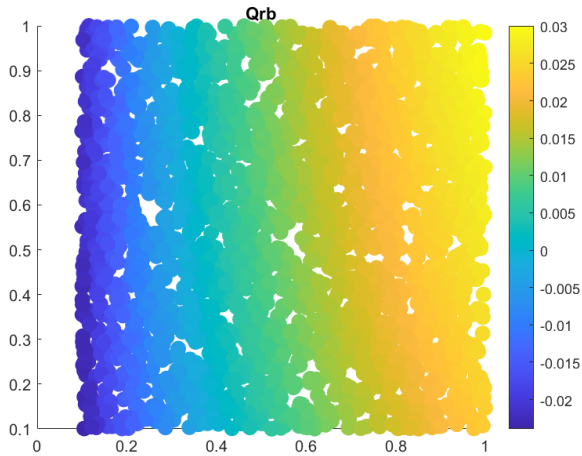


FIGURE 3.1 –  $Q^{\text{rb}}(\mu)$  pour 3000 points

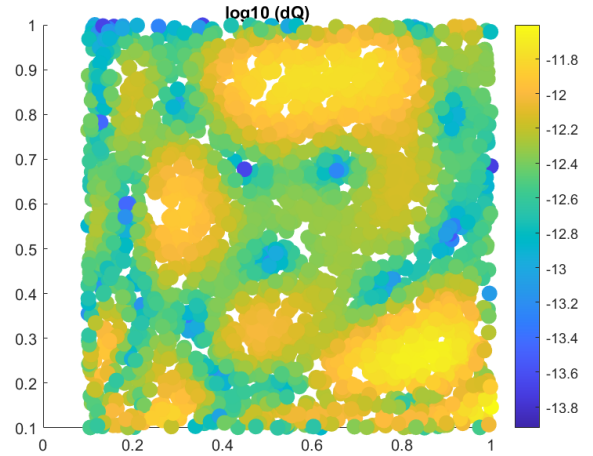


FIGURE 3.2 –  $\delta_N^Q(\mu)$  pour 3000 points

Les images FIGURE 3.1 et 3.2 sont le résultats pour  $Q^{\text{rb}}(\mu)$  et  $\delta_N^Q(\mu)$ .

Selon la valeur de  $\delta_N^Q(\mu)$ , on peut dire qu'on a une très bonne précision pour  $Q^{\text{rb}}(\mu)$ .

De plus, pour le temps de calcul, le temps utilisé total est 0.5251s. Le temps pour calculer les 3000 différents  $Q^{\text{rb}}(\mu)$  est 0.2280s. Le temps pour calculer l'estimation à posteriori est 0.2971s. Grâce à la méthode de base réduite, on a vraiment gagné beaucoup de temps de calculs et obtenir une très bonne précision (environ  $10e - 11$ ).