



**HAL**  
open science

# Some problems related to statistical error in stochastic homogenization

William Minvielle

► **To cite this version:**

William Minvielle. Some problems related to statistical error in stochastic homogenization. Numerical Analysis [math.NA]. Université Paris-Est Marne la Vallée, 2015. English. NNT: . tel-01429837

**HAL Id: tel-01429837**

**<https://hal.science/tel-01429837>**

Submitted on 9 Jan 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



L'École Doctorale Mathématiques et Sciences et Technologies de  
l'Information et de la Communication (MSTIC)

THÈSE DE DOCTORAT

présentée pour obtenir le titre de

Docteur de l'Université Paris-Est  
Spécialité : Mathématiques

présentée par

**William MINVIELLE**

**Sujet :** *Quelques problèmes liés à l'erreur statistique  
en homogénéisation stochastique*

*(Some problems related to statistical error in stochastic homogenization)*

Soutenue le 25 septembre 2015 devant le jury composé de :

*Rapporteurs :* M. Yves ACHDOU  
M. James NOLEN

*Examineurs :* M. Grégoire ALLAIRE  
M. Antoine LEJAY

*Directeurs de Thèse :* Claude LE BRIS  
Frédéric LEGOLL



**Titre :** Quelques problèmes liés à l'erreur statistique en homogénéisation stochastique.

**Résumé :** Le travail de cette thèse a porté sur le développement de techniques numériques pour l'homogénéisation d'équations dont les coefficients présentent des hétérogénéités aléatoires à petite échelle. Les difficultés liées à la résolution de telles équations aux dérivées partielles peuvent être résolues grâce à la théorie de l'homogénéisation stochastique. On substitue alors la résolution d'une équation dont les coefficients sont aléatoires et oscillants à l'échelle la plus fine du problème par la résolution d'une équation à coefficients constants. Cependant, une difficulté subsiste : le calcul de ces coefficients dits homogénéisés sont définis par une moyenne ergodique, que l'on ne peut atteindre en pratique. Seuls des approximations aléatoires de ces quantités déterministes sont calculables, et l'erreur commise lors de l'approximation est importante. Ces questions sont développées en détail dans le **Chapitre 1** qui tient lieu d'introduction. L'objet du **Chapitre 2** de cette thèse est de réduire l'erreur de cette approximation dans un cas nonlinéaire, en réduisant la variance de l'estimateur par la méthode des variables antithétiques. Dans le **Chapitre 3**, on montre comment obtenir une meilleure réduction de variance par la méthode des variables de contrôle. Cette approche repose sur un modèle approché, disponible dans le cas étudié. Elle est plus invasive et moins générique, on l'étudie dans un cas linéaire. Dans le **Chapitre 4**, à nouveau dans un cas linéaire, on introduit une méthode de sélection pour réduire l'erreur commise. Enfin, le **Chapitre 5** porte sur l'analyse d'un problème inverse, où l'on recherche des paramètres à l'échelle la plus fine, ne connaissant que quelques quantités macroscopiques, par exemple les coefficients homogénéisés du modèle.

**Mots clés :** Equations aux dérivées partielles, Homogénéisation, Matériaux aléatoires, Méthodes de Monte Carlo, Réduction de variance, Problème inverse.

**Title:** Some problems related to statistical error in stochastic homogenization.

**Abstract:** In this thesis, we design numerical techniques to address the homogenization of equations the coefficients of which exhibit small scale random heterogeneities. Solving such elliptic partial differential equations is prohibitively expensive. One may use stochastic homogenization theory to reduce the complexity of this task. We then substitute the random, fine scale oscillating coefficients of the equation with constant homogenized coefficients. These coefficients are defined through an ergodic average inaccessible to practical computation. Only random approximations thereof are available. The error committed in this approximation is significant. These issues are detailed in the introductory **Chapter 1**. In **Chapter 2**, we show how to reduce the error in this approximation, in a nonlinear case, by using an antithetic variable estimator that has a smaller variance than the standard Monte Carlo estimator. In **Chapter 3**, in a linear case, we show how to obtain an even better variance reduction with the control variate method. Such a method is based on a surrogate model. In **Chapter 4**, we use a selection method to reduce the global error. **Chapter 5** is devoted to the analysis of an inverse problem, wherein we seek parameters at the fine scale whilst only being provided with a handful of macroscopic quantities, among which the homogenized coefficients.

**Keywords:** Partial differential equations, Homogenization, Random materials, Monte Carlo methods, Variance reduction, Inverse problem.

# Contents

<b>1</b>	<b>Introduction</b>	<b>11</b>
1.1	Introduction générale	11
1.2	Introduction mathématique	13
1.2.1	Modélisation	13
1.2.2	Stationnarité	14
1.2.3	Homogénéisation	15
1.2.4	Décomposition de l'erreur	16
1.2.5	Quantification de la convergence	16
1.3	Résumé des travaux	17
1.3.1	Réduction de variance par variables antithétiques pour un problème d'homogénéisation stochastique nonlinéaire convexe.	17
1.3.2	Une approche par variables de contrôle basée sur une théorie perturbative pour la réduction de variance en homogénéisation stochastique.	20
1.3.3	Structures quasi-aléatoires spéciales : une approche de sélection pour l'homogénéisation stochastique.	22
1.3.4	Un problème d'identification de paramètres en homogénéisation stochastique.	25
1.4	Perspectives	28
1.4.1	Comparaison des méthodes de réduction de variance	28
1.4.2	Extension à la stationnarité continue des variables de contrôle	28
1.4.3	Utilisation d'autres variables de contrôle	29
1.4.4	Sélection pour des problèmes nonlinéaires	29
1.4.5	Convergence du problème inverse en dimension supérieure	29
1.4.6	Perspectives générales dans le champ de recherche	30
<b>2</b>	<b>Variance reduction using antithetic variables for a nonlinear convex stochastic homogenization problem</b>	<b>31</b>
2.1	Introduction	32
2.1.1	Homogenization theoretical setting	33
2.1.2	The questions we consider	35
2.2	Description of the proposed approach and main results	36
2.2.1	Statement of our main results	36
2.2.2	Classical results on antithetic variables	40
2.2.3	Derivatives of the corrector and of the homogenized energy density	40
2.2.4	Monotonicity properties	41
2.2.5	Proof of Propositions 2.1, 2.2 and 2.5	45
2.2.6	Examples satisfying our structure assumptions	46
2.3	Numerical results	48
2.3.1	Newton algorithm to solve the truncated corrector problem	48

2.3.2	Overview of numerical results . . . . .	49
2.3.3	Test Case 1 . . . . .	50
2.3.4	Test Case 2 . . . . .	53
2.3.5	Test Case 3 . . . . .	55
2.4	Appendix: Proof of (2.8) . . . . .	55
<b>3</b>	<b>A control variate approach based on a defect-type theory for variance reduction in stochastic homogenization</b>	<b>59</b>
3.1	Introduction . . . . .	60
3.1.1	Homogenization theoretical setting . . . . .	62
3.1.2	Practical approximation of the homogenized matrix . . . . .	63
3.1.3	Control variate approach . . . . .	64
3.2	A weakly random setting: rare defects in a periodic structure . . . . .	66
3.2.1	Presentation of the model . . . . .	66
3.2.2	Weakly-random homogenization result . . . . .	67
3.3	Control variate approaches for stochastic homogenization . . . . .	69
3.3.1	A first-order model . . . . .	69
3.3.2	A second-order model . . . . .	72
3.4	Elements of theoretical analysis . . . . .	74
3.4.1	One-dimensional case . . . . .	74
3.4.2	Multi-dimensional case . . . . .	79
3.5	Numerical results . . . . .	81
3.5.1	Low contrast test-case . . . . .	82
3.5.2	High contrast test-case . . . . .	84
3.5.3	Using a Reduced Basis (RB) approach . . . . .	84
<b>4</b>	<b>Special Quasirandom Structures: a selection approach for stochastic homogenization</b>	<b>91</b>
4.1	Introduction . . . . .	92
4.1.1	Overview . . . . .	92
4.1.2	Theoretical setting . . . . .	94
4.1.3	Numerical approximation of the homogenized matrix . . . . .	95
4.2	Variance reduction approach . . . . .	96
4.2.1	Original formulation of the SQS approach . . . . .	96
4.2.2	Formal derivation of the SQS conditions using a perturbative setting . . . . .	97
4.2.3	Practical evaluation of the SQS conditions . . . . .	101
4.2.4	Selection Monte Carlo sampling . . . . .	105
4.3	Elements of theoretical analysis . . . . .	107
4.3.1	Proof of convergence of the approach . . . . .	107
4.3.2	Complete analysis in some simple cases . . . . .	114
4.4	Numerical experiments . . . . .	117
4.4.1	Robustness of the approach . . . . .	117
4.4.2	Efficiency of the approach . . . . .	118
4.5	Proof of Lemma 4.4 . . . . .	125
<b>5</b>	<b>A parameter identification problem in stochastic homogenization</b>	<b>129</b>
5.1	Introduction . . . . .	131
5.2	Discrete homogenization theory . . . . .	133
5.2.1	Homogenization result . . . . .	133

5.2.2	Approximation on finite boxes . . . . .	135
5.2.3	Physical problem . . . . .	137
5.2.4	The one dimensional case . . . . .	141
5.3	A parameter fitting problem . . . . .	145
5.3.1	General case . . . . .	145
5.3.2	The one-dimensional case . . . . .	146
5.4	Numerical results . . . . .	150
5.4.1	Optimization algorithm . . . . .	151
5.4.2	Numerical results . . . . .	151
5.5	Appendix: Computation of the derivatives of (5.50) . . . . .	154
<b>A Neumann, Dirichlet and periodic approximations of the corrector problem</b>		<b>155</b>
A.1	Introduction . . . . .	155
A.2	Approximations of the corrector . . . . .	155
A.2.1	The Dirichlet approximation . . . . .	156
A.2.2	The periodic approximation . . . . .	156
A.2.3	Neumann approximation . . . . .	157
A.3	Flux formulation . . . . .	158
A.3.1	The $H^{\text{div}}$ space . . . . .	158
A.3.2	The periodic $H_{\text{per}}^{\text{div}}$ space . . . . .	159
A.3.3	Flux formulation of the periodic problem (A.4) . . . . .	159
A.3.4	Flux formulation of the Neumann problem (A.9) . . . . .	160
A.4	Comparison of the approximations . . . . .	161
A.4.1	Comparison of the Dirichlet and periodic approximations . . . . .	161
A.4.2	Comparison of the Neumann and periodic approximations . . . . .	162
<b>Bibliography</b>		<b>165</b>





# Publications

- [LM15b] F. Legoll and W. Minvielle.  
Variance reduction using antithetic variables for a nonlinear convex stochastic homogenization problem.  
*Discrete and Continuous Dynamical Systems - Series S*, 8(1):1–27, 2015.
- [LM15a] F. Legoll and W. Minvielle.  
A control variate approach based on a defect-type theory for variance reduction in stochastic homogenization.  
*Multiscale Modeling & Simulation*, 13(2):519–550, 2015.
- [LMOS15] F. Legoll, W. Minvielle, A. Obliger and M. Simon.  
A parameter identification problem in stochastic homogenization.  
*ESAIM: Proc.*, 48:190–214, 2015.
- [LBLM] C. Le Bris, F. Legoll and W. Minvielle.  
Special Quasirandom Structures: a selection approach for stochastic homogenization.  
To be submitted.



# Chapter 1

## Introduction

### 1.1 Introduction générale

Les matériaux multiéchelles sont des matériaux présentant des caractéristiques qui varient selon l'échelle à laquelle on les examine. Ainsi, un béton contient des pores inclus dans une matrice solide, la taille de ces pores pouvant être micrométrique. Le comportement (par exemple élastique) macroscopique de l'ensemble est une combinaison du comportement des différentes phases, mais ne saurait se contenter d'une simple moyenne algébrique des rigidités de chaque phase. En d'autres termes, bien que le calcul intéressant l'ingénieur (par exemple, certifier la résistance du matériau) soit fait à une échelle de l'ordre du mètre au moins, la description mécanique du matériau est effectuée à une échelle bien plus fine. Ces matériaux multiéchelles sont présents dans nombre de produits industriels (automobile, aéronautique, génie civil). Par exemple, dans le secteur de l'énergie nucléaire, des matériaux multiéchelles sont utilisés pour construire des réacteurs nucléaires (bétons, aciers) et envisagés pour le stockage des déchets radioactifs (argiles).

Comment garantir la sécurité des ouvrages et anticiper leur comportement ? Les techniques usuelles d'expérimentation physique peuvent être inadaptées si le matériau coûte trop cher à produire, ou si l'on souhaite prédire ou concevoir un matériau nouveau. L'expérimentation *numérique* prend alors toute sa place. Cependant, la résolution des équations de la mécanique sur des échelles macroscopiques pour des matériaux multiéchelles n'est pas *a priori* aisée, car reproduire fidèlement la microstructure requiert déjà de mailler à une échelle très fine, sans compter la résolution subséquente des équations. Pour autant, il existe des techniques spécialisées pour de tels problèmes, et l'homogénéisation - dont il est question dans cette thèse - en est une.

L'homogénéisation est une théorie mathématique qui permet de simplifier la description du comportement de matériaux multiéchelles. A l'échelle macroscopique, les propriétés du milieu sont décrites par des propriétés moyennes, constantes ou variant lentement selon la variable d'espace. C'est un modèle adapté pour mener des calculs à l'échelle la plus grande, car il n'est plus nécessaire d'avoir une résolution de calcul très fine. Le contrôle de l'erreur commise lors de cette substitution est exprimée par un paramètre de *séparation d'échelle*  $\varepsilon \ll 1$ , qui quantifie le rapport de tailles caractéristiques entre la structure microscopique et la taille macroscopique de l'échantillon sur lequel porte le calcul. Il existe des situations où un continuum d'échelles doivent être prises en compte, et la séparation d'échelle n'a pas lieu. Les techniques (pratiques ou théoriques) pour les traiter sont encore largement exploratoires, et nous ne couvrons pas ce cas dans cette thèse. Lorsqu'il y a séparation d'échelle, plus  $\varepsilon$  est petit, plus le modèle homogénéisé est précis.

Cette théorie propose ainsi de remplacer le calcul dans un milieu hétérogène par un calcul dans un milieu homogène. Le calcul des coefficients du milieu homogène est l'étape cruciale en homogénéisation, et aussi la plus coûteuse.

Une hypothèse de *structure* est nécessaire pour obtenir des expressions explicites. Dans le livre fondateur [BLP78], les auteurs supposent que la structure microscopique est la répétition *périodique* d'un même motif. Le comportement moyen est déductible de la résolution d'un problème "typique" sur la maille périodique (le problème du *correcteur*, parfois appelé problème de cellule dans ce cas).

Cependant, l'hypothèse de périodicité est parfois restrictive. Si des matériaux peuvent assez bien la satisfaire, d'autres s'en éloignent nettement. Un exemple d'hypothèse de structure peu contraignante est la stationnarité (voir les articles fondateurs [Koz79, PV81]). On suppose que le matériau observé correspond à *une réalisation* d'un matériau aléatoire, qui possède une propriété d'invariance spatiale de sa *loi*. La périodicité correspond alors à un cas particulier. Certaines situations, par exemple un réseau cristallin perturbé par un défaut local ayant une influence importante (dopage de semi-conducteurs par exemple), ne peuvent pas être modélisé comme un problème aléatoire stationnaire, toutefois cette classe est très vaste.

Cette généralité vient avec certains inconvénients. D'un point de vue pratique, il faut connaître de la loi du milieu de référence, ce qui peut être difficile s'agissant de quantités microscopiques. De plus, l'équation du correcteur permettant de calculer le comportement moyen change de niveau de complexité par rapport au cas périodique : il s'agit d'une équation à coefficients aléatoires, posée dans un domaine non borné. Intuitivement, on ne peut pas se contenter d'un domaine fini pour décrire la variété infinie de la loi de cette variable aléatoire. Ainsi, pour simplifier la description du comportement macroscopique, il faut d'abord en passer par un calcul presque aussi complexe que le problème original, avec toutefois une simplification importante : le chargement correspond à un chargement de référence et non au chargement originel pour lequel on cherche à calculer la solution. Par conséquent, d'un point de vue pratique, l'homogénéisation aléatoire ne permet pas de réduire significativement la complexité du problème de départ si l'objectif n'est d'effectuer qu'un seul calcul, pour un seul chargement et une seule occurrence des paramètres de la loi de la microstructure. Toutefois, il existe des situations où il faut effectuer des calculs *répétés*, pour de nombreux chargements différents. Ainsi, pour résoudre un problème d'évolution temporelle du matériau multiéchelle, il faudra résoudre de manière répétée, pour différents chargements, l'équation statique de ce matériau. De même, la résolution d'un problème *inverse* comprend nécessairement l'évaluation répétée, pour différentes valeurs de paramètres, du calcul *direct*. Citons enfin la décomposition de domaine, où la solution du problème est obtenue par itérations successives de la résolution du problème par sous-domaines, les itérations ayant pour but de calculer la solution aux interfaces. L'homogénéisation aléatoire réduit alors globalement la difficulté : au lieu de devoir effectuer de très nombreux calculs coûteux, un seul sera nécessaire - afin de déterminer milieu homogène associé - après quoi les suivants seront réalisés au moyen du milieu homogène, donc rapides.

D'autres méthodes multiéchelles, alternatives à l'homogénéisation, ont été développées. Citons par exemple l'approche MsFEM (méthode des éléments finis multiéchelles). L'homogénéisation constitue alors un guide pour justifier ces méthodes récentes et encore exploratoires.

L'intérêt de l'homogénéisation étant établi, il reste que le calcul du coefficient homogénéisé est coûteux. Les approches considérées dans cette thèse ont pour objectif de

rendre ce calcul plus efficace.

Puisque le problème du correcteur est posé sur un domaine non borné, il faut pour le résoudre numériquement le tronquer sur un domaine borné. Du fait de cette troncature, on obtient alors une approximation aléatoire du milieu homogène, de laquelle on approche l'espérance grâce à une méthode de Monte Carlo : on considère un grand nombre de réalisations de l'environnement pour lesquelles on résout le problème du correcteur, et l'on considère la moyenne empirique de ces coefficients homogénéisés *apparents* pour approcher le coefficient homogénéisé exact. L'erreur totale est en grande partie liée à la dispersion statistique des coefficients apparents.

Dans les **Chapitre 2** et **Chapitre 3**, nous construisons des approximations ayant la même espérance que l'estimateur de Monte Carlo usuel, mais avec une dispersion statistique plus faible. L'approximation du coefficient homogénéisé construit selon une telle approche sera moins entaché d'erreur. Dans le **Chapitre 2**, on a recours à une méthode de réduction de variance classique et peu dépendante du problème, la technique des *variables antithétiques*. Dans le **Chapitre 3**, on utilise la technique des *variables de contrôle*, nécessitant une connaissance du milieu de référence plus précise et un travail préparatoire. Nous poursuivons vers une méthode davantage liée au problème, donc obtenant une meilleure réduction d'erreur, dans le **Chapitre 4** où nous développons une méthode de *sélection*. L'objectif du **Chapitre 5** est différent. On s'intéresse à un problème *inverse*, et l'on cherche les paramètres de la loi de l'environnement microscopique connaissant des quantités macroscopiques (par exemple les coefficients homogénéisés).

## 1.2 Introduction mathématique

### 1.2.1 Modélisation

Commençons par un exemple représentatif des problèmes qui seront traités dans cette thèse. Soit  $d$  la dimension ambiante. On considère un domaine borné et régulier  $\mathcal{D} \subset \mathbb{R}^d$  qui représente le solide que l'on cherche à modéliser. En élasticité linéaire, la grandeur caractéristique de ce solide est le tenseur d'élasticité  $A_\varepsilon$  qui est un tenseur d'ordre 4, symétrique coercif et borné, reliant le déplacement à la contrainte selon une relation constitutive linéaire. L'indice  $\varepsilon$  souligne que cette rigidité est une caractéristique *locale* du solide, qui varie à l'échelle  $\varepsilon \ll 1$ . Un problème typique consiste à chercher le déplacement  $u_\varepsilon \in H^1(\mathcal{D})^d$  d'un solide soumis à des forces volumiques  $f \in H^{-1}(\mathcal{D})^d$  et attaché au bord  $\partial\mathcal{D}$ . L'équation aux dérivées partielles résultante s'écrit alors

$$\begin{cases} -\operatorname{div} A_\varepsilon : \nabla^s u_\varepsilon = f & \text{dans } \mathcal{D}, \\ u_\varepsilon = 0 & \text{sur } \partial\mathcal{D}. \end{cases} \quad (1.1)$$

On a noté  $\nabla^s : H^1(\mathcal{D})^d \mapsto L^2(\mathcal{D})^{d \times d}$  la partie symétrique du gradient. Il est possible d'adapter les techniques, classiques ou nouvelles, ainsi que leurs justifications mathématiques, à cette équation. Cependant, pour plus de simplicité dans l'expérimentation numérique, nous avons choisi dans cette thèse de travailler sur une version scalaire de cette équation, qui conserve néanmoins toute la difficulté du caractère multiéchelle :

$$\begin{cases} -\operatorname{div} A_\varepsilon \nabla u_\varepsilon = f & \text{dans } \mathcal{D}, \\ u_\varepsilon = 0 & \text{sur } \partial\mathcal{D}. \end{cases} \quad (1.2)$$

Dans (1.2), le second membre  $f \in H^{-1}(\mathcal{D})$  est une fonction scalaire, et non vectorielle. Le champ  $A_\varepsilon \in L^\infty(\mathcal{D})^{d \times d}$  est matriciel et l'inconnue  $u_\varepsilon \in H^1(\mathcal{D})$  est scalaire. De très

nombreuses situations physiques sont modélisées par (1.2). La conduction thermique en est une : l'équation (1.2) gouverne la température  $u_\varepsilon$  ;  $A_\varepsilon$  est le coefficient de diffusion thermique et  $f$  le terme source. L'électrostatique, où  $u_\varepsilon$  est le potentiel scalaire,  $f$  la distribution de charge et  $A_\varepsilon$  la perméabilité diélectrique du milieu en est une autre.

L'équation (1.2) est typique en sciences de l'ingénieur, mais d'autres modèles, parfois plus complexes, existent. En particulier, dans le **Chapitre 2** sera explorée une version plus générale et nonlinéaire de ce problème, et dans le **Chapitre 5** nous verrons une version discrète de l'équation aux dérivées partielles (1.2).

### 1.2.2 Stationnarité

Nous supposons que  $A_\varepsilon$  s'écrit sous la forme  $A_\varepsilon(x, \omega) = A\left(\frac{x}{\varepsilon}, \omega\right)$  où  $\omega \in \Omega$  représente la variable d'aléa, et  $x \in \mathbb{R}^d$  la variable d'espace. C'est une hypothèse importante, selon laquelle le champ  $A_\varepsilon$  est identique en chaque point macroscopique. En outre, nous supposons que  $A$  est *stationnaire*, c'est à dire possède une propriété d'invariance spatiale de sa loi : celle-ci est la même en chaque point microscopique. Le paragraphe suivant précise cette propriété.

Soit  $(\Omega, \mathcal{F}, \mathbb{P})$  un espace de probabilité. On suppose que le groupe  $(\mathbb{Z}^d, +)$  agit sur  $\Omega$ . On note  $(\tau_k)_{k \in \mathbb{Z}^d}$  cette action, et on suppose qu'elle préserve la mesure  $\mathbb{P}$ , c'est-à-dire que pour tout  $k \in \mathbb{Z}^d$  et  $B \in \mathcal{F}$ ,  $\mathbb{P}(\tau_k B) = \mathbb{P}(B)$ . Nous supposons que cette action est *ergodique* : si  $B \in \mathcal{F}$  est invariant selon tous les  $\tau_k$ , alors  $\mathbb{P}(B) = 0$  ou 1.

**Définition 1.1.** Une fonction  $F \in L^1_{\text{loc}}(\mathbb{R}^d, L^1(\Omega))$  est stationnaire si

$$\forall k \in \mathbb{Z}^d, \quad F(x + k, \omega) = F(x, \tau_k \omega) \quad \text{presque partout et presque sûrement.} \quad (1.3)$$

Les fonctions stationnaires (au sens de la stationnarité *discrète* comme ci-dessus, ou *continue* comme ci-après dans la Remarque 1.5, pour ces dernières  $\mathbb{R}^d$  remplace  $\mathbb{Z}^d$  dans l'équation (1.3)) fournissent le socle de la théorie de l'homogénéisation aléatoire et vérifient le théorème ergodique (voir [Kre85, Shi84, Tem72]) suivant :

**Théorème 1.2.** Soit une fonction stationnaire  $F \in L^\infty(\mathbb{R}^d, L^1(\Omega))$ . On note, pour  $k \in \mathbb{Z}^d$ ,  $|k|_\infty = \sup_{1 \leq i \leq d} |k_i|$ . Alors

$$\frac{1}{(2N+1)^d} \sum_{|k|_\infty \leq N} F(x, \tau_k \omega) \xrightarrow{N \rightarrow \infty} \mathbb{E}(F(x, \cdot)) \quad \text{dans } L^\infty(\mathbb{R}^d), \text{ presque sûrement.}$$

En conséquence, en appelant  $Q$  le cube unité de  $\mathbb{R}^d$ ,

$$F\left(\frac{x}{\varepsilon}, \omega\right) \xrightarrow[\varepsilon \rightarrow 0]{*} \mathbb{E}\left(\int_Q F(x, \cdot) dx\right) \quad \text{dans } L^\infty(\mathbb{R}^d), \text{ presque sûrement.} \quad (1.4)$$

Autrement dit, les fonctions stationnaires admettent des moyennes sur les grands volumes.

**Remarque 1.3.** Dans le Théorème 1.2, si la fonction  $F$  est  $(0, 1)^d$ -périodique, on retrouve alors une version du Lemme de Riemann-Lebesgue :

$$\forall \varphi \in L^1(\mathbb{R}^d), \quad \int_{\mathbb{R}^d} \varphi(x) F\left(\frac{x}{\varepsilon}\right) dx \xrightarrow{\varepsilon \rightarrow 0} \int_{(0,1)^d} F(y) dy \int_{\mathbb{R}^d} \varphi(x) dx.$$

**Remarque 1.4.** Une autre hypothèse du même type, mais différente, est la stationnarité continue. On suppose que  $\mathbb{R}^d$  agit sur  $\Omega$  selon une action  $(\tau_y)_{y \in \mathbb{R}^d}$ , et la stationnarité de  $F$  s'écrit alors

$$\forall y \in \mathbb{R}^d, \quad F(x + y, \omega) = F(x, \tau_y \omega) \quad \text{presque partout et presque sûrement.} \quad (1.5)$$

La relation (1.4) demeure en remplaçant l'opérateur de moyenne  $\mathbb{E} \int_Q$  par  $\mathbb{E}$ .

### 1.2.3 Homogénéisation

Sous l'hypothèse de stationnarité de  $A$ , on peut énoncer le résultat d'homogénéisation suivant (voir [Koz79, PV81]) :

**Théorème 1.5.** Soit  $A \in L^\infty(\mathbb{R}^d \times \Omega)^{d \times d}$ . On suppose que

- (ellipticité) : il existe une constante  $\alpha > 0$  telle que, pour presque tout  $x \in \mathbb{R}^d$ , pour tout  $p \in \mathbb{R}^d$  et presque sûrement en  $\omega \in \Omega$ ,

$$\alpha |p|^2 \leq p \cdot A(x, \omega) p.$$

- (stationnarité) :  $A$  est stationnaire au sens de (1.3).

Soit  $\mathcal{D}$  un domaine borné de  $\mathbb{R}^d$ ,  $f \in H^{-1}(\mathcal{D})$  et  $A_\varepsilon(x, \omega) := A(x/\varepsilon, \omega)$ . Alors, presque sûrement, la solution du problème (1.2) converge faiblement dans  $H_0^1(\mathcal{D})$  vers  $u^* \in H_0^1(\mathcal{D})$  solution de l'équation

$$\begin{cases} -\operatorname{div} A^* \nabla u^* = f & \text{dans } \mathcal{D}, \\ u^* = 0 & \text{sur } \partial \mathcal{D}. \end{cases} \quad (1.6)$$

La matrice constante, déterministe et  $\alpha$ -coercive  $A^* \in \mathbb{R}^{d \times d}$  est donnée par

$$\forall p \in \mathbb{R}^d, \quad A^* p := \mathbb{E} \int_Q A(p + \nabla w_p), \quad (1.7)$$

où le correcteur  $w_p \in L^2(\Omega; L_{\text{loc}}^2(\mathbb{R}^d))$  est l'unique solution (à une constante additive près) dont le gradient est dans l'espace  $\nabla w_p \in L^2(\Omega; L_{\text{unif}}^2(\mathbb{R}^d))^d$  de l'équation

$$\begin{cases} -\operatorname{div} A(p + \nabla w_p) = 0 & \text{dans } \mathbb{R}^d, \\ \nabla w_p & \text{stationnaire et } \mathbb{E} \int_Q \nabla w_p = 0. \end{cases} \quad (1.8)$$

La convergence donnée par le Théorème 1.5 justifie l'approximation de la solution de (1.2) par celle de (1.6) lorsque  $\varepsilon \ll 1$ . Au contraire de (1.2), le problème (1.6) est facile à résoudre numériquement puisque  $A^*$  est constant et que l'équation est posée sur un domaine borné.

Pour autant, il demeure une difficulté dans ce programme: calculer  $A^*$  est en général difficile. Plus précisément, pour calculer  $A^*$ , il faut résoudre  $d$  problèmes du correcteur (1.8) (un dans chaque direction  $p = e_i$ ,  $1 \leq i \leq d$ ), lesquels sont posés dans  $\mathbb{R}^d$ . La solution recherchée n'est pas décroissante à l'infinie, seulement strictement sous-linéaire. Contrairement au cas périodique, dans le cas aléatoire, on ne peut pas transformer ce problème en un problème de cellule sur un domaine compact<sup>1</sup>.

<sup>1</sup>Plus exactement, le domaine permettant de recouvrer assez de compacité est l'espace de probabilité. Malheureusement, cet espace est trop abstrait pour aider à la résolution numérique du problème du correcteur.



Une solution pour résoudre (1.8) consiste à tronquer le domaine de résolution. Nous résolvons l'équation sur un sous-domaine  $Q_N := (0, N)^d$  de taille finie et aussi grande que possible en la munissant de conditions aux limites adéquates. Par exemple, on considère dans la suite le problème

$$\begin{cases} -\operatorname{div} A(p + \nabla w_p^N) = 0 & \text{dans } \mathbb{R}^d, \\ w_p^N \text{ est } Q_N \text{ périodique.} \end{cases} \quad (1.9)$$

Les conditions aux limites dans (1.9) ont été choisies périodiques. D'autres conditions peuvent être utilisées ; elles seront discutées dans la Section A.2. En conséquence de ces conditions de bords périodiques, le problème (1.9) se ramène à un problème sur le domaine compact  $Q_N$  vu comme un tore. De cette approximation du correcteur nous déduisons une approximation  $A^{*,N}(\omega)$  de  $A^*$  définie par

$$\forall p \in \mathbb{R}^d, \quad A^{*,N}(\omega)p := \frac{1}{|Q_N|} \int_{Q_N} A(x, \omega)(p + \nabla w_p^N(x, \omega)) dx. \quad (1.10)$$

Cette approximation converge presque sûrement vers  $A^*$  d'après un résultat de [BP04, Theorem 1]. Voir aussi la Section A.2.2.

#### 1.2.4 Décomposition de l'erreur

Pour autant, cette approximation souffre de plusieurs faiblesses. Premièrement, en quantifier la convergence est un problème délicat. Nous y revenons dans la Section 1.2.5. Deuxième obstacle, il s'agit d'une approximation aléatoire d'une quantité déterministe. Pire, même l'espérance de cette approximation est distincte de  $A^*$ . Commençons par décomposer l'erreur selon deux contributions orthogonales dans  $L^2(\Omega)$  :

$$A^* - A^{*,N}(\omega) = A^* - \mathbb{E}[A^{*,N}] + \mathbb{E}[A^{*,N}] - A^{*,N}(\omega). \quad (1.11)$$

La première partie de l'erreur dans (1.11), appelée erreur *systematique*, est complètement déterministe. Nous verrons dans la Section 1.2.5 qu'elle n'est pas dominante. La seconde erreur est appelée erreur *statistique*, elle dérive de la dispersion statistique de  $A^{*,N}$  autour de sa moyenne. Une première manière de réduire cette erreur est de considérer  $M$  copies indépendantes et identiquement distribuées de  $A^{*,N}$ , notées  $(A_j^{*,N})_{1 \leq j \leq M}$ , en résolvant (1.9) pour différents environnements aléatoires, puis d'en prendre la moyenne empirique

$$\mathcal{I}_{MC}^M(\omega) := \frac{1}{M} \sum_{j=1}^M A_j^{*,N}(\omega), \quad (1.12)$$

afin d'obtenir un *estimateur* (approximation aléatoire) de  $\mathbb{E}[A^{*,N}]$ . Par le théorème de la limite centrale, l'erreur commise  $|\mathcal{I}_{MC}^M - \mathbb{E}[A^{*,N}]|$  décroît en loi comme  $\sqrt{\frac{\operatorname{Var} A^{*,N}}{M}}$  lorsque  $M$  tend vers l'infini. Notons que  $\operatorname{Var} A^{*,N}$  tend vers 0 lorsque  $N$  tend vers l'infini. Dans les **Chapitre 2** et **Chapitre 3**, nous introduisons des méthodes pour réduire cette erreur statistique, en construisant une approximation de même espérance  $\mathbb{E}[A^{*,N}]$  que  $\mathcal{I}_{MC}^M$  mais de variance plus faible. L'erreur statistique est donc réduite, et l'erreur systématique reste la même. Dans le **Chapitre 4**, nous réduisons globalement l'erreur totale.

#### 1.2.5 Quantification de la convergence

La motivation principale d'une partie des travaux de cette thèse est l'amélioration de la convergence d'estimateurs particuliers vers  $A^*$ , plus performants que  $\mathcal{I}_{MC}^M$  et que nous introduirons ci-dessous. Étudier théoriquement l'avantage de ces nouveaux estimateurs n'est

pas simple. En effet, la simple *quantification* théorique de l'erreur commise dans (1.11), par exemple sous la forme d'une estimation d'erreur  $|A^* - \mathbb{E}[A^{*,N}]| \leq CN^{-\alpha}$ , est complexe. Les premiers travaux en ce sens proviennent de [Yur86]. Ces résultats non optimaux n'ont été que récemment complétés, d'abord en homogénéisation discrète (voir le **Chapitre 5** pour une présentation de cette modélisation). Nous renvoyons à l'article [GNO14] pour une revue récente de ces avancées. Dans un second temps, ils ont été pour partie étendus à des cas continus, certaines études étant toujours en cours.

L'intérêt pour notre travail, est multiple. Nous l'utilisons premièrement comme justification de nos motivations : l'analyse d'erreur confirme qu'il faut réduire l'erreur statistique, principale source de l'erreur. Deuxièmement, il s'agit d'une source d'inspiration pour établir des résultats nouveaux, et pour comparer les hypothèses requises.

## 1.3 Résumé des travaux

### 1.3.1 Réduction de variance par variables antithétiques pour un problème d'homogénéisation stochastique nonlinéaire convexe.

Dans le **Chapitre 2**, nous nous intéressons au problème d'optimisation multiéchelle et nonlinéaire suivant :

$$\mathcal{E}_\varepsilon(f, \mathcal{D}) := \inf \left\{ \int_{\mathcal{D}} W\left(\frac{x}{\varepsilon}, \omega, \nabla u(x)\right) dx - \int_{\mathcal{D}} fu, \quad u \in W_0^{1,p}(\mathcal{D}) \right\}, \quad (1.13)$$

pour un certain  $1 < p < \infty$ , et un chargement  $f \in W_0^{1,p}(\mathcal{D})'$ .

#### Caractère bien posé de (1.13)

Considérons tout d'abord la version mono échelle et déterministe de ce problème :

$$\inf \left\{ \int_{\mathcal{D}} W(x, \nabla u(x)) dx - \int_{\mathcal{D}} fu, \quad u \in W_0^{1,p}(\mathcal{D}) \right\}. \quad (1.14)$$

Sous les conditions de croissance et convexité suivantes de  $W$ ,

$$\exists 2 \leq p < \infty, \exists c, C > 0, \forall x, \xi \in \mathbb{R}^d, \quad c|\xi|^p \leq W(x, \xi) \leq C(1 + |\xi|^p), \quad (1.15)$$

$$L'application partielle \xi \mapsto W(x, \xi) \text{ est convexe pour presque tout } x, \quad (1.16)$$

le problème d'optimisation (1.14) est bien posé.

**Remarque 1.6.** *Sous les hypothèses (1.15)–(1.16),  $\xi \mapsto W(x, \xi)$  est strictement convexe.*

Ce problème est une déclinaison très simplifiée d'un problème de mécanique nonlinéaire.

**Remarque 1.7.** *Dans l'article [Bal76], l'auteur considère un problème dont l'inconnue  $u : \mathbb{R}^d \mapsto \mathbb{R}^d$  est vectorielle plutôt que scalaire. Il s'agit du champ de déplacement d'un solide soumis à un chargement par des forces volumiques  $f : \mathcal{D} \mapsto \mathbb{R}^d$ , et l'énergie (1.14) représente l'énergie de déformation, qu'il convient de minimiser étant donné un chargement. L'auteur discute des hypothèses de croissance et convexité.*

Nous avons choisi l'hypothèse de croissance (1.15), alors même qu'elle n'a pas de contrepartie physique : il s'agit d'une hypothèse technique. Nous renvoyons à [JKO94, Chapter 15] pour des hypothèses plus générales.

**Remarque 1.8.** *Sous l'hypothèse (1.16), le problème (1.14) est équivalent à l'équation d'optimalité d'Euler-Lagrange, qui s'écrit*

$$\begin{cases} -\operatorname{div}[\partial_\xi W(x, \nabla u(x))] = f(x) & \text{dans } \mathcal{D}, \\ u(x) = 0 & \text{sur } \partial\mathcal{D}. \end{cases} \quad (1.17)$$

Cette équation est linéaire dans le cas où  $W$  s'écrit sous la forme  $W(x, \xi) = \frac{1}{2}\xi \cdot A(x)\xi$  et vérifie (1.15)–(1.16) pour  $p = 2$ .

### Homogénéisation stochastique

Afin de discuter l'homogénéisation du problème (1.13), nous supposons que

$$\text{Pour presque tout } x \in \mathbb{R}^d, \quad \xi \mapsto W(x, \omega, \xi) \text{ est presque sûrement continue,} \quad (1.18)$$

$$\forall \xi \in \mathbb{R}^d, \quad (x, \omega) \mapsto W(x, \omega, \xi) \text{ est stationnaire au sens de (1.3).} \quad (1.19)$$

L'hypothèse de stationnarité (1.19) est une hypothèse de *structure*, permettant d'obtenir des moyennes sur des grands volumes. Elle a été discutée dans la Section 1.2.2. L'hypothèse (1.18) est une condition technique classique. L'application  $(x, \xi) \mapsto W(x, \omega, \xi)$ , mesurable par rapport à la variable  $\xi$  et continue par rapport à la variable  $x$ , est appelée fonction de Carathéodory. Grâce à cette hypothèse, une fonction de la forme  $x \mapsto W(x, \omega, g(x))$  est mesurable si  $g$  l'est.

**Théorème 1.9** ([DMM86a, DMM86b]). *On suppose (1.15)–(1.16) et (1.18)–(1.19). Alors le problème (1.13) admet un unique minimiseur  $u_\varepsilon$ , lequel converge presque sûrement, faiblement dans  $W_0^{1,p}(\mathcal{D})$ , vers la solution du problème d'optimisation homogénéisé*

$$\mathcal{E}_0(f, \mathcal{D}) := \inf \left\{ \int_{\mathcal{D}} W^*(\nabla u^*) - \int_{\mathcal{D}} f u, \quad u \in W_0^{1,p}(\mathcal{D}) \right\}, \quad (1.20)$$

où la densité d'énergie homogénéisée  $W^*$  est donnée par

$$W^*(\xi) := \lim_{N \rightarrow \infty} \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} W(x, \omega, p + \nabla w(x)), \quad w \in W_{\text{per}}^{1,p}(Q_N) \right\}. \quad (1.21)$$

A l'instar du cas linéaire, on ne peut toutefois espérer atteindre cette limite, aussi nous considérons dans la suite l'approximation standard suivante

$$W_N^*(\omega, \xi) := \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} W(x, \omega, p + \nabla w(x)), \quad w \in W_{\text{per}}^{1,p}(Q_N) \right\}, \quad (1.22)$$

qui est l'analogie nonlinéaire de (1.10). La limite (1.21) a lieu presque sûrement, et par définition,  $W^*(\xi)$  est la limite presque sûre de  $W_N^*(\omega, \xi)$ . Les conditions aux limites dans (1.22) peuvent être modifiées pour construire d'autres approximations convergentes. En particulier, on peut remplacer  $w \in W_{\text{per}}^{1,p}(Q_N)$  par  $w \in W_0^{1,p}(Q_N)$  avec le même résultat. Nous renvoyons au Lemma 2.16 pour plus de détails. Le Théorème 1.9 généralise la situation linéaire *symétrique*, réalisée pour  $p = 2$  et  $W(x, \omega, \xi) := \frac{1}{2}\xi \cdot A(x, \omega)\xi$ , quand  $A$  est symétrique. Dans ce cas, la formulation énergétique du problème (1.2) (respectivement (1.6)) correspond exactement au problème (1.13) (respectivement (1.20)), et le Théorème 1.9 se réduit au Théorème 1.5.

### Réalisations antithétiques

La technique des variables antithétiques est une méthode simple et facile à mettre en oeuvre pour réduire la variance d'une approximation. Notons d'abord qu'une variable aléatoire peut souvent s'écrire comme une certaine fonction de plusieurs variables aléatoires indépendantes et identiquement distribuées, uniformes sur  $(0, 1)$ . Nous appellerons  $X_k \sim \mathcal{U}(0, 1)$  ces variables aléatoires qui représentent la « source » de l'aléa. Un cas particulier de fonctions stationnaires qui reviendra souvent dans la thèse est le suivant :

$$\chi(x, \omega) := \sum_{k \in \mathbb{Z}^d} g(X_k(\omega)) \mathbb{1}_{k+Q}(x), \quad (1.23)$$

où  $g : (0, 1) \mapsto \mathbb{R}$  est une fonction bornée, et  $Q = (0, 1)^d$ . Le champ aléatoire  $\chi$  est constant sur chaque cellule  $k+Q$ , où il prend la valeur  $g(X_k)$ . Sous une hypothèse adéquate sur  $X_k$  (par exemple s'ils sont indépendants), le champ  $\chi$  est stationnaire. Ainsi le cas  $W(x, \omega, \xi) = f(\xi)\chi(x, \omega)$  est un exemple de  $W$  satisfaisant les hypothèses (1.18)–(1.19) si  $f$  est continue.

Notons  $X_k^{\text{ant}} := 1 - X_k$  la réalisation *antithétique* à  $X_k$ . On pose  $\chi^{\text{ant}}(x, \omega) := \sum_{k \in \mathbb{Z}^d} g(X_k^{\text{ant}}(\omega)) \mathbb{1}_{k+Q}(x)$ , et  $W^{\text{ant}}(x, \omega, \xi) := f(\xi)\chi^{\text{ant}}(x, \omega)$ . Le champ  $W^{\text{ant}}$  est distribué suivant la même loi que  $W$ . On définit

$$W_N^{\star, \text{ant}}(\omega, \xi) := \inf \left\{ \frac{1}{N^d} \int_{Q_N} W^{\text{ant}}(x, \omega, p + \nabla w(x)), \quad w \in W_{\text{per}}^{1,p}(Q_N) \right\},$$

puis on introduit l'approximation de  $W^{\star}$  suivante :

$$\widetilde{W}_N^{\star}(\omega, \xi) := \frac{1}{2} \left( W_N^{\star}(\omega, \xi) + W_N^{\star, \text{ant}}(\omega, \xi) \right). \quad (1.24)$$

Puisque  $W_N^{\star, \text{ant}}(\cdot, \xi)$  a la même loi que  $W_N^{\star}(\cdot, \xi)$ , l'espérance de  $\widetilde{W}_N^{\star}(\cdot, \xi)$  est identique à celle de  $W_N^{\star}(\cdot, \xi)$ , donc l'erreur systématique définie dans (1.11) n'est pas modifiée par le changement d'estimateur. L'erreur statistique est, elle, modifiée, car la variance de  $\widetilde{W}_N^{\star}$  vient remplacer celle de  $W_N^{\star}$ .

La réduction de variance est alors démontrée grâce à des propriétés de monotonie de  $W_N^{\star}$  par rapport à ses variables d'aléa. On peut en effet utiliser un résultat du type suivant pour conclure.

**Lemme 1.10** ([Liu08]). *Soit  $f : \mathbb{R} \mapsto \mathbb{R}$  une fonction croissante. Soit  $X : \Omega \mapsto \mathbb{R}$  une variable aléatoire distribuée selon la loi uniforme  $\mathcal{U}[0, 1]$ . Alors*

$$\mathbb{V}\text{ar} \left( \frac{1}{2} (f(X) + f(1 - X)) \right) \leq \frac{1}{2} \mathbb{V}\text{ar} (f(X)).$$

### Réduction de l'erreur statistique

La stratégie de calcul présentée dans la Section 1.2.4 revient ici à résoudre  $2M$  problèmes du type (1.22), et de calculer l'approximation aléatoire

$$\mathcal{I}_{MC}^M(\omega) := \frac{1}{2M} \sum_{j=1}^{2M} W_{N,j}^{\star}(\omega),$$

où les  $W_{N,j}^*$  sont des copies indépendantes et identiquement distribuées de  $W_N^*$ . La stratégie des variables antithétiques consiste à mettre en oeuvre cette méthode pour  $\tilde{W}_N^*$  : on considère l'approximation aléatoire

$$\mathcal{I}_{\text{ant}}^M(\omega) := \frac{1}{M} \sum_{j=1}^M \tilde{W}_{N,j}^*(\omega).$$

Puisque chaque calcul de  $\tilde{W}_N^*$  nécessite deux résolutions de (1.22), le coût calcul de ces deux méthodes est identique.

Le résultat principal du **Chapitre 2** est la Proposition 2.1, où nous démontrons que la variance de l'estimateur antithétique  $\mathcal{I}_{\text{ant}}^M$  est plus petite que celle de l'estimateur standard  $\mathcal{I}_{MC}^M$ . Pour ce faire, nous établissons dans (2.13) que la variance de  $\tilde{W}_N^*$  est plus petite que la moitié de la variance de  $W_N^*$ . L'estimateur  $\mathcal{I}_{\text{ant}}^M$  ayant le même coût que  $\mathcal{I}_{MC}^M$  et une variance plus petite, il est donc plus performant.

Nous montrons également, dans le cas unidimensionnel et sous des hypothèses plus contraignantes, que nous réduisons la variance des approximations de  $\xi \cdot \partial_\xi W_N^*(\cdot, \xi)$  (respectivement  $\xi \cdot \partial_\xi^2 W_N^*(\cdot, \xi)\xi$ ) en considérant  $\xi \cdot \partial_\xi \tilde{W}_N^*(\cdot, \xi)$  (respectivement  $\xi \cdot \partial_\xi^2 \tilde{W}_N^*(\cdot, \xi)\xi$ ). Ces quantités sont identiques dans le cas où  $W$  prend la forme  $W(x, \omega, \xi) = \frac{1}{2}\xi \cdot A(x, \omega)\xi$  ce qui revient au cas linéaire traité dans [BCLBL12a, BCLBL12b, CLBL10].

## Résultats numériques en dimension deux

Nous illustrons ces résultats théoriques par des expériences numériques en dimension  $d = 2$ . En particulier, nous observons numériquement la réduction de variance sur  $\partial_\xi W_N^*$  et  $\partial_{\xi\xi}^2 W_N^*$ , bien que nous n'ayons pas prouvé ceci en dimension  $d \geq 2$ .

### 1.3.2 Une approche par variables de contrôle basée sur une théorie perturbative pour la réduction de variance en homogénéisation stochastique.

Dans le **Chapitre 3**, nous introduisons une méthode de variable de contrôle pour obtenir une approximation de  $A^*$  plus performante que  $\mathcal{I}_{MC}^M$  (défini par (1.12)), car de plus faible variance. Considérons le cas où

$$A := A_\eta(x, \omega) = A_{\text{per}}(x) + b_\eta(x, \omega)(C_{\text{per}} - A_{\text{per}})(x) \quad (1.25)$$

est une perturbation à l'ordre  $\eta$  de  $A_{\text{per}}$ . Nous dérivons une approximation  $Y_\eta^2$  précise à l'ordre 2 en loi de  $A_\eta^{*,N}$ . Cependant, nous n'utilisons pas  $Y_\eta^2(\omega)$  dans le régime  $\eta \ll 1$  pour approcher directement  $A^{*,N}(\omega)$ , nous l'utilisons dans le cas où  $\eta$  n'est pas petit, en tant que variable de contrôle. De plus, nous utilisons la méthode de base réduite de [LBT12] pour calculer efficacement  $Y_\eta^2(\omega)$ .

## Méthode des variables de contrôle

Présentons tout d'abord la méthode des variables de contrôle.

Il s'agit d'une méthode standard pour réduire l'erreur sur l'estimation de l'espérance d'une variable aléatoire. Cette méthode est utile lorsqu'est disponible une variable réduite, d'espérance connue, bien corrélée à la variable d'intérêt. Soit  $X \in L^2(\Omega)$  la variable aléatoire dont on cherche à estimer l'espérance, et  $Y \in L^2(\Omega)$  la variable de contrôle. Soit  $\rho \in \mathbb{R}$ , on pose :

$$X_\rho(\omega) := X(\omega) - \rho(Y(\omega) - \mathbb{E}[Y]). \quad (1.26)$$

La méthode des variables de contrôle consiste à approcher l'espérance de  $X_\rho$  par méthode de Monte Carlo :

$$\mathcal{I}_{CV}^{M,\rho}(\omega) := \frac{1}{M} \sum_{k=1}^M X_\rho^k(\omega), \quad (1.27)$$

où les  $X_\rho^k$  sont des copies indépendantes de  $X_\rho$ . Puisque  $\rho$  n'est pas aléatoire, il est évident que  $\mathbb{E}[X_\rho] = \mathbb{E}[X]$ , ainsi  $\mathcal{I}_{CV}^{M,\rho}$  est un estimateur convergent quand  $M$  tend vers l'infini de  $\mathbb{E}[X]$ . Observons également que

$$\text{Var}X_\rho = \text{Var}X - 2\rho\text{Cov}[X, Y] + \rho^2\text{Var}Y,$$

où  $\text{Cov}[a, b] := \mathbb{E}[(a - \mathbb{E}a)(b - \mathbb{E}b)]$  est la covariance de deux variables aléatoires. Puisque  $\text{Var}X_\rho$  est une forme quadratique par rapport à  $\rho$ , elle admet un minimum  $\rho^* = \frac{\text{Cov}[X, Y]}{\text{Var}Y}$ . Pour cette valeur de  $\rho$ , nous obtenons ainsi  $\text{Var}X_{\rho^*} = \text{Var}X \left(1 - \frac{\text{Cov}[X, Y]^2}{\text{Var}X\text{Var}Y}\right)$ . L'estimateur  $\mathcal{I}_{CV}^{M,\rho^*}$  est d'autant plus performant que sa variance est petite, ce qui est le cas si  $X$  et  $Y$  sont fortement corrélés.

On estime en pratique  $\rho^*$  par des moyennes empiriques. Nous discutons le choix de  $Y$  (qui dépend du problème) ci-après.

### Construction d'une expansion en loi de $A_\eta^{*,N}$

Soit  $A_\eta$  de la forme (1.25). Nous supposons de plus que  $b_\eta$  prend la forme

$$b_\eta(x, \omega) := \sum_{k \in \mathbb{Z}^d} B_k^\eta(\omega) \mathbb{1}_{k+Q}(x), \quad (1.28)$$

où  $B_k^\eta$  sont des variables de Bernoulli indépendantes et identiquement distribuées de paramètre  $\eta$  :  $\mathbb{P}[B_k^\eta = 1] = \eta$  et  $\mathbb{P}[B_k^\eta = 0] = 1 - \eta$ . Ainsi, si  $\eta \ll 1$ , les défauts, représentés par  $B_k = 1$ , sont rares, mais leurs effets sont importants car ils changent à l'ordre 1 la valeur locale de la matrice  $A$ . On définit la matrice  $A_1^k := A_{\text{per}} + (C_{\text{per}} - A_{\text{per}}) \mathbb{1}_{k+Q}$ , correspondant au milieu  $A_{\text{per}}$  perturbé par un défaut localisé en  $k \in \mathbb{Z}^d \cap Q_N$ , et le correcteur  $w_p^{1,k,N}$  associé à un défaut localisé en  $k \in \mathbb{Z}^d \cap Q_N$ , solution  $Q_N$ -périodique de l'équation

$$-\text{div} \left[ A_1^k(p + \nabla w_p^{1,k,N}) \right] = 0 \quad \text{sur } Q_N. \quad (1.29)$$

On montre que la quantité  $\bar{A}_{1\text{def}}^{k,N} := \int_{Q_N} A_1^k(p + \nabla w_p^{1,k,N}) - \int_{Q_N} A_{\text{per}}(p + \nabla w_{\text{per}})$ , qui représente la variation de la réponse de l'environnement de référence perturbée par un défaut localisé en  $k + Q$ , joue un rôle important. Ainsi, dans le Lemme 3.4, on montre que

$$\forall \varphi \in C(\mathbb{R}), \quad \mathbb{E}[\varphi(A_\eta^{*,N})] = \mathbb{E} \left[ \varphi \left( A_{\text{per}}^* + \sum_{k \in \mathbb{Z}^d \cap Q_N} B_k^\eta \bar{A}_{1\text{def}}^{k,N} \right) \right] + O(\eta^2). \quad (1.30)$$

Cela suggère d'utiliser  $Y_\eta^1(\omega) := A_{\text{per}}^* + \sum_{k \in \mathbb{Z}^d \cap Q_N} B_k^\eta(\omega) \bar{A}_{1\text{def}}^{k,N}$  comme variable de contrôle. Pour un certain  $\rho \in \mathbb{R}$ , la variable aléatoire suivante :

$$X_\rho(\omega) := A_\eta^{*,N}(\omega) - \rho(Y_\eta^1(\omega) - \mathbb{E}Y_\eta^1), \quad (1.31)$$

d'espérance  $\mathbb{E}[A_\eta^{*,N}]$ , a une variance significativement plus petite que celle de  $A_\eta^{*,N}$ .

Nous avons amélioré (1.31) en considérant plusieurs contrôles,  $Y_\eta^1(\omega)$ , le contrôle lié à l'expansion poussée à l'ordre deux  $Y_\eta^2(\omega)$  (considérant la perturbation de  $A_{\text{per}}$  par un couple de défauts). Nous utilisons ainsi la variable aléatoire

$$X_{\rho_1, \rho_2}(\omega) = A^{*,N}(\omega) - \rho_1(Y_\eta^1(\omega) - \mathbb{E}[Y_\eta^1]) - \rho_2(Y_\eta^2(\omega) - \mathbb{E}[Y_\eta^2]),$$

où  $\rho_1$  et  $\rho_2$  sont choisis afin de minimiser la variance de  $X_{\rho_1, \rho_2}$ .

### Calcul par bases réduites de la partie déterministe du modèle

Il est crucial pour la méthode que le calcul des  $Y_\eta^2$  soit très peu coûteux. Nous avons recours à la méthode de base réduite développée dans [LBT12], afin d'avoir un coût de calcul essentiellement indépendant de  $N$ .

Ainsi, le calcul de la variable de contrôle n'est pas coûteux en pratique (le surcoût représente 13% de temps de calcul dans nos simulations), tandis que la réduction de variance est très significative.

### Analyse de la méthode en dimension 1

Nous montrons que la variance de  $A^{*,N}$  est d'ordre  $N^{-1}$ , alors que la variance de  $X^1 := A^{*,N} - \rho_N^1(Y_\eta^1 - \mathbb{E}[Y_\eta^1])$  (où  $\rho_N^1$  est le paramètre optimal) est d'ordre  $N^{-2}$  (voir Proposition 3.11), puis nous montrons que la variance de  $X^2 := X_{\rho_1^N, \rho_2^N}$  (où  $\rho_1^N, \rho_2^N$  sont les paramètres optimaux) est d'ordre  $N^{-3}$  (voir Proposition 3.13).

### Résultats numériques en dimension 2

Nous testons notre méthode, toujours dans un cas non perturbatif, où le paramètre de la loi de Bernoulli est  $\eta = 0.5$  avec  $A_{\text{per}} = 3Id$  et  $C_{\text{per}} = 23Id$  et nous observons une réduction de variance de l'ordre de 40 dans ce cas, indépendamment de la taille de la supercellule.

#### 1.3.3 Structures quasi-aléatoires spéciales : une approche de sélection pour l'homogénéisation stochastique.

Dans le Chapitre 4, nous introduisons une troisième méthode de réduction de variance. Il s'agit d'une approche de sélection. Nous ne calculons la solution de (1.9) que pour certaines réalisations de  $A(x, \omega)$  sur  $Q_N$ . Nous sélectionnons celles qui vérifient au mieux des propriétés normalement satisfaites seulement asymptotiquement quand  $N$  tend vers l'infini. L'approche de sélection que nous avons développée s'inspire d'une méthode venant de la physique du solide. L'approche consiste à sélectionner des réalisations qui vérifient au mieux certaines propriétés. Par exemple, pour un matériau composé de deux phases  $A$  et  $B$ , on sélectionne les réalisations de sorte que la fraction volumique de chaque phase dans la supercellule soit exactement la proportion de chaque matériau (propriété qui n'est généralement réalisée que dans la limite des grandes supercellules  $Q_N$ ). Dans un alliage biphasique, on sélectionne ainsi les réalisations selon la fraction volumique, et selon d'autres quantités (proportion d'atomes  $B$  voisins des atomes  $A$ , proportions de chaîne d'atomes  $BB$  voisines de  $A$ , etc.).

### Dérivation des conditions

Pour dériver les conditions utilisées pour la sélection des réalisations, nous considérons ici un modèle faiblement aléatoire. Toutefois, ce régime perturbatif n'est utile que pour



construire les critères de sélection, et la méthode est testée ensuite dans un cas non perturbatif.

Nous étudions un cas où la matrice  $A$  prend la forme

$$A := A_\eta(x, \omega) = C_0 + \eta\chi(x, \omega)C_1(x), \quad (1.32)$$

où  $C_0$  est constante,  $C_1$  est périodique, et  $-1 \leq \chi \leq 1$  est stationnaire et prend la forme

$$\chi(x, \omega) := \sum_{k \in \mathbb{Z}^d} X_k(\omega) \mathbf{1}_{k+Q}(x). \quad (1.33)$$

Ainsi, la source d'aléa du problème est entièrement contenue dans le champ aléatoire  $\chi$ , lequel est aléatoire à travers les variables aléatoires scalaires  $(X_k)_{k \in \mathbb{Z}^d}$ . Dans l'article [BCLBL12b], les auteurs établissent que la matrice homogénéisée  $A_\eta^*$  vérifie le développement limité suivant :

$$A_\eta^* = C_0 + \eta A_1^* + \eta^2 A_2^* + o(\eta^2), \quad (1.34)$$

où le coefficient  $A_1^*$  est simplement  $\mathbb{E} \int_Q \chi C_1$ , et où  $A_2^*$  peut être calculé en résolvant dans  $\{u \in L^2_{\text{loc}}(\mathbb{R}^d), \nabla u \in L^2(\mathbb{R}^d)^d\}$  le problème suivant

$$-\operatorname{div} C_0 \nabla \phi_1 = \operatorname{div} \mathbf{1}_Q C_1 p \quad \text{dans } \mathbb{R}^d. \quad (1.35)$$

De même, on peut montrer que la matrice homogénéisée apparente  $A_\eta^{*,N}(\omega)$  vérifie un développement limité similaire

$$A_\eta^{*,N}(\omega) = C_0 + \eta A_1^{*,N}(\omega) + \eta^2 A_2^{*,N}(\omega) + o(\eta^2), \quad (1.36)$$

où le coefficient  $A_1^{*,N}(\omega)$  vaut  $\frac{1}{|Q_N|} \int_{Q_N} \chi(x, \omega) C_1(x) dx$ , et où  $A_2^{*,N}(\omega)$  peut être calculé en résolvant la version tronquée de (1.35) :

$$-\operatorname{div} C_0 \nabla \phi_1^N = \operatorname{div} \mathbf{1}_Q C_1 p \quad \text{dans } Q_N, \text{ et } \phi_1^N \text{ est } Q_N\text{-périodique.} \quad (1.37)$$

On dit qu'un environnement  $\omega \in \Omega$  satisfait la condition d'ordre  $k$  ( $k = 1, 2$  que l'on peut généraliser) si l'égalité  $A_k^{*,N}(\omega) = A_k^*$  est vérifiée. Si un certain  $\omega$  satisfait les conditions jusqu'à l'ordre  $k$ , nous avons alors

$$A_\eta^*(\omega) = A_\eta^{*,N}(\omega) + o(\eta^k),$$

où le reste est uniforme par rapport à  $N$  et  $\omega$ . Jusqu'à l'ordre 2, dans cet exemple, les conditions que nous imposons sont faciles à évaluer.

En toute généralité, la condition à l'ordre 1 s'écrit

$$\frac{1}{|Q_N|} \int_{Q_N} \chi(x, \omega) C_1(x, \omega) p dx = \mathbb{E} \int_Q \chi C_1 p. \quad (1.38)$$

Lorsque de plus  $C_1$  est constante, et si  $\chi$  prend la forme (1.33) alors cette condition se réécrit

$$\frac{1}{N^d} \sum_{|k|_\infty \leq N} X_k(\omega) = \mathbb{E}[X_0]. \quad (1.39)$$

En utilisant seulement cette sélection au premier ordre, nous obtenons l'algorithme suivant :



**Algorithme 1 (Sélection-Monte Carlo).**

Cet algorithme requiert une tolérance  $\text{tol} \geq 0$ .

Pour  $m = 1, \dots, M$ ,

1. Générer un environnement aléatoire  $\omega_m$ .
2. Si  $\left| \frac{1}{Nd} \sum_{|k|_\infty \leq N} X_k - \mathbb{E}[X_0] \right| > \text{tol}$ , retourner à l'étape 1.
3. Résoudre le problème du correcteur tronqué (1.9).
4. Calculer  $A^{*,N}(\omega_m)$  par (1.10).

Calculer l'approximation de  $\mathbb{E}[A^*]$  suivante :  $\mathcal{I}_{SQS}^M := \frac{1}{M} \sum_{m=1}^M A_N^*(\omega_m)$ .

Cet algorithme diffère de l'algorithme usuel par l'étape 2. Le problème du correcteur (1.9) n'est résolu que pour certaines réalisations.

En général, le coût des méthodes de sélection est dominé par l'étape de sélection elle-même, car un grand nombre de tirages est nécessaire pour produire *une* réalisation convenable, et le coût d'un tirage est élevé. Cependant, nous ne sommes pas ici dans un tel cas : le coût de résolution de (1.9) domine complètement le coût des autres étapes. Ainsi, il est pertinent de bien choisir les  $\omega \in \Omega$  pour lesquels on résout (1.9), même si cette sélection s'effectue en rejetant beaucoup de réalisations.

**Contrôle de l'erreur**

Seules certaines réalisations sélectionnées participent à l'estimation de  $A^*$ , ce qui se traduit mathématiquement de la façon suivante : au lieu d'échantillonner  $A^{*,N}(\omega)$  selon la mesure de probabilité de référence  $\mathbb{P}$ , on l'échantillonne suivant une mesure *conditionnelle*. L'algorithme ainsi proposé correspond exactement à une méthode de Monte-Carlo standard pour l'estimation de  $\mathbb{E}^{\mu_{\text{cond}}}[A^{*,N}]$ , où  $\mu_{\text{cond}}$  est la mesure de référence conditionnée à satisfaire les critères (par exemple (1.39)). Afin d'analyser l'erreur à l'instar de (1.11), nous introduisons  $\mathbb{E}^{\mu_{\text{cond}}}[A^{*,N}]$  au lieu de  $\mathbb{E}[A^{*,N}]$  et nous obtenons la décomposition

$$A^* - A^{*,N}(\omega) = A^* - \mathbb{E}^{\mu_{\text{cond}}}[A^{*,N}] + \mathbb{E}^{\mu_{\text{cond}}}[A^{*,N}] - A^{*,N}(\omega). \quad (1.40)$$

La première partie de l'erreur est une erreur de biais. Comme mentionné dans la Section 1.2.5, on peut s'attendre à ce qu'elle ne soit pas dominante, mais l'analyse complète du taux de convergence de cette erreur est une question non résolue. Sous des hypothèses assez générales, nous montrons dans le Théorème 4.8 qu'elle tend vers 0. Cela assure la convergence de l'approche. Nous vérifions numériquement que l'erreur de biais n'est pas plus grande que dans le cas Monte Carlo usuel, voire même est diminuée en utilisant notre méthode.

La seconde erreur est bien sûr de nature statistique. Cependant, elle est contrôlée par la variance de  $A^{*,N}$  sous la mesure conditionnelle, dont on s'attend à ce qu'elle soit plus petite. Nous montrons, dans certains cas simplifiés (en particulier en dimension un d'espace) que c'est le cas : la variance est réduite.

### Expérimentations numériques en dimension deux

En pratique, les conditions que doivent satisfaire les réalisations peuvent ne pas être *exactement* satisfaites, et nous autorisons une certaine déviation comme dans l’Algorithme 1 en autorisant un niveau d’erreur pour la satisfaction du critère, ou alors en faisant un grand nombre de tirages, et en sélectionnant les meilleurs des échantillons. Nous vérifions numériquement, dans un cas non perturbatif et en dimension  $d = 2$ , la robustesse de notre méthode : même en présence d’une certaine erreur dans la condition que nous imposons, notre approche permet de significativement réduire l’erreur (1.40).

Toujours dans un cas non perturbatif et en dimension  $d = 2$ , nous observons numériquement une bonne réduction de variance ainsi qu’une réduction de l’erreur totale.

#### 1.3.4 Un problème d’identification de paramètres en homogénéisation stochastique.

Dans le **Chapitre 5**, nous abordons un problème inverse dans un cadre d’homogénéisation aléatoire. Le cadre d’étude est celui d’équations aux dérivées partielles *discrètes*. Toutefois il serait possible d’effectuer la même démarche dans le cadre d’équations aux dérivées partielles continues considéré jusqu’à présent.

#### Modélisation d’un réseau de pores

Le problème de physique motivant notre étude est un problème de transport d’ions dans un réseau de pores au sein d’argiles. La physique complexe décrivant ce milieu rend la modélisation complète de la géométrie du milieu difficile. Une modélisation usuellement adoptée est celle, simpliste, d’un réseau structuré (le réseau  $\mathbb{Z}^d$ ) dont les noeuds sont les pores et les arêtes les canaux. Le problème d’intérêt est le suivant : calibrer les paramètres de la loi de la taille des canaux (à l’échelle microscopique), à l’aide de simples expériences hydrauliques sur une espèce non chargée (qui correspond à un calcul de perméabilité effective donc d’homogénéisation), pour diminuer le nombre de paramètres du modèle de transport ionique.

#### Équations aux dérivées partielles discrètes

Nous étudions l’équation

$$\nabla_\varepsilon^* [A(x/\varepsilon, \omega) \nabla_\varepsilon u_\varepsilon(x, \omega)] = f(x) \quad \text{dans } \mathcal{D} \cap \varepsilon\mathbb{Z}^d, \quad u_\varepsilon(x, \omega) = 0 \quad \text{dans } (\mathbb{R}^d \setminus \mathcal{D}) \cap \varepsilon\mathbb{Z}^d, \quad (1.41)$$

contrepartie discrète de (1.2). L’opérateur de différences finies  $\nabla_\varepsilon : \ell^2(\varepsilon\mathbb{Z}^d) \mapsto \ell^2(\varepsilon\mathbb{Z}^d)^d$  est défini par

$$\nabla_\varepsilon g : x \mapsto \frac{1}{\varepsilon} \begin{pmatrix} g(x + \varepsilon e_1) - g(x) \\ \vdots \\ g(x + \varepsilon e_d) - g(x) \end{pmatrix}.$$

L’opérateur  $\nabla_\varepsilon^* : \ell^2(\varepsilon\mathbb{Z}^d)^d \mapsto \ell^2(\varepsilon\mathbb{Z}^d)$  (dont l’équivalent au niveau continu est  $-\text{div}$ ) est défini comme l’adjoint dans  $\ell^2(\varepsilon\mathbb{Z}^d)$  de  $\nabla_\varepsilon$ . C’est pourquoi il correspond à une différence finie dont le décentrage est inverse :

$$-\nabla_\varepsilon^* G : x \mapsto \sum_{i=1}^d \frac{G_i(x) - G_i(x - \varepsilon e_i)}{\varepsilon}.$$

### Estimation de paramètres

On suppose que la loi de probabilité de  $A$  est paramétrée par un certain  $\theta \in \Theta$ . Ce paramètre, de faible dimension (de dimension 2 dans notre travail) est supposé inconnu. Ainsi, la loi de probabilité de  $A$  est supposée mal connue, et nous cherchons à estimer  $\theta$ , au moyen de quantités qui sont facilement accessibles à l'expérience : des quantités macroscopiques. Nous utilisons pour ce faire deux quantités (et leur nombre est réminiscent de la dimension de l'espace des paramètres à identifier : il est nécessaire d'avoir au moins deux quantités connues pour en estimer deux), qui sont le premier coefficient de la matrice homogénéisée  $[A^*]_{1,1}$  et la variance relative de la matrice homogénéisée apparente

$$\text{VarR} \left[ (A^{*,N})_{1,1} \right] := \frac{\text{Var} \left[ (A^{*,N})_{1,1} \right]}{\left( \mathbb{E} \left[ (A^{*,N})_{1,1} \right] \right)^2}. \quad (1.42)$$

Comme le problème d'intérêt est isotrope, nous ne nous intéressons qu'à un seul coefficient diagonal de  $A^*$ .

Nous supposons que ces deux quantités sont connues, par exemple grâce à une expérience physique, et nous cherchons alors à retrouver le paramètre  $\theta \in \Theta$  associé. Nous supposons en outre une forme particulière de dépendance de  $A$  par rapport à  $\theta$  : nous supposons, pour des raisons physiques, que  $A = aId$ , où  $a$  suit une loi de Weibull à deux paramètres  $\theta := (\lambda, k) \in (\mathbb{R}_+^*)^2$ . Rappelons que ces variables aléatoires sont positives, et leur densité s'écrit

$$\forall r > 0, \quad f(r; k, \lambda) = \frac{k}{\lambda} \left( \frac{r}{\lambda} \right)^{k-1} \exp \left( - (r/\lambda)^k \right),$$

ce qui correspond à la distribution cumulée

$$F(r; k, \lambda) = \int_0^r f(s; k, \lambda) ds = 1 - \exp \left( - (r/\lambda)^k \right).$$

Nous formulons ce problème *inverse* sous la forme d'un problème d'optimisation aux moindres carrés. On introduit

$$F_{N,M}(\theta, \omega) := \left( \frac{\overline{K}_M^{*,N}(\theta, \omega)}{K_{\text{obs}}^{*,N}} - 1 \right)^2 + \left( \frac{S_M^N(\theta, \omega)}{S_{\text{obs}}^N} - 1 \right)^2, \quad (1.43)$$

où  $\overline{K}_M^{*,N}(\theta, \omega) := \frac{1}{M} \sum_{k=1}^M [A_k^{*,N}(\theta, \omega)]_{1,1}$  est une approximation de  $e_1 \cdot \mathbb{E}[A^{*,N}(\theta, \cdot)]e_1$ , et

$S_M^N(\theta, \omega)$  est une approximation de la variance relative (1.42). Les quantités  $K_{\text{obs}}^{*,N}$  et  $S_{\text{obs}}^N$  sont les données de notre problème, et on suppose qu'il existe un paramètre déterministe  $\theta_{\text{obs}}$  tel que  $K_{\text{obs}}^{*,N} = e_1 \cdot \mathbb{E}[A^{*,N}(\theta_{\text{obs}}, \cdot)]e_1$  et de même pour  $S_{\text{obs}}^N$ . Nous cherchons donc à minimiser  $F_{N,M}$  par rapport à  $\theta$ .

Lorsque  $N$  tend vers l'infini, la fonctionnelle  $F_{N,M}(\theta, \omega)$  admet une limite  $F_\infty(\theta)$ . Nous démontrons, en une dimension d'espace, que le problème d'optimisation  $\inf_{\theta \in \Theta} F_\infty(\theta)$  est bien posé, et admet pour unique solution  $\theta_{\text{obs}}$ .

### Résolution numérique

Pour résoudre ce problème inverse, nous devons *évaluer* de façon répétée  $F_{N,M}(\theta, \omega)$ , ainsi que ses dérivées, pour différentes valeurs de  $\theta \in \Theta$ . C'est là la partie principale du temps

de calcul, et c'est pourquoi nous avons recours à un algorithme d'optimisation efficace pour atteindre une faible erreur en peu d'itérations, minimisant ainsi le nombre d'appels de la fonctionnelle  $F_{N,M}$ .

Afin de résoudre le problème d'optimisation (1.43), nous avons recours à un algorithme de Newton. Nous pouvons en effet calculer les dérivées de  $A^{*,N}(\theta, \omega)$  (et donc de  $F_{N,M}$ ) par rapport à  $\theta$  jusqu'à l'ordre 2 sans recourir à des différences finies et ceci à coût comparable à la simple *évaluation* de  $F_{N,M}$  (voir paragraphe suivant).

Nous obtenons ainsi en quelques itérations une bonne approximation d'un minimiseur local du problème (1.43), qu'on note  $\theta^*(\omega)$ . Puisque  $F_{N,M}$  est une variable aléatoire,  $\theta^*$  est aussi une variable aléatoire. Nous vérifions numériquement que la variance de  $\theta^*$  est comparable à la variance de  $\bar{K}_M^{*,N}$  et  $S_M^N$ . Ainsi l'erreur statistique supplémentaire introduite par la résolution du problème inverse n'est pas trop importante. Nous vérifions en outre que  $\mathbb{E}[\theta^*]$  est une bonne approximation de  $\theta_{\text{obs}}$ .

Les expériences numériques menées dans le **Chapitre 5** ont été conduites en une dimension d'espace, puis généralisées en dimension deux.

### Calcul des dérivées de $A^{*,N}(\theta, \omega)$ par rapport au paramètre $\theta = (k, \lambda)$

Pour simplifier, nous supposons que  $A$  est symétrique. De même, nous utilisons les notations des équations aux dérivées partielles *continues* mais les calculs s'étendent au cadre discret.

Au vu de (1.10) et de la formulation variationnelle du problème (1.9), nous pouvons symétriser davantage l'équation définissant  $A^{*,N}(\omega)$  :

$$\forall p, q \in \mathbb{R}^d, \quad q \cdot A^{*,N} p = \frac{1}{|Q_N|} \int_{Q_N} (q + \nabla w_q^N) \cdot A(p + \nabla w_p^N).$$

En utilisant les équations d'Euler-Lagrange, il vient, pour la dérivée par rapport à  $k$ ,

$$\forall p, q \in \mathbb{R}^d, \quad q \cdot \partial_k A^{*,N} p = \frac{1}{|Q_N|} \int_{Q_N} (q + \nabla w_q^N) \cdot \partial_k A(p + \nabla w_p^N)$$

et de même pour la dérivée par rapport à  $\lambda$ .

On considère maintenant les dérivées secondes de  $A^{*,N}$ . Soient  $\tau_1$  et  $\tau_2$  deux directions de dérivation ( $\tau_1 = k$  ou  $\lambda$  et de même pour  $\tau_2$ ). En dérivant l'équation ci-dessus par rapport à  $\tau_2$ , il vient :

$$\begin{aligned} q \cdot \partial_{\tau_1 \tau_2}^2 A^{*,N} p &= \frac{2}{|Q_N|} \int_{Q_N} \nabla \partial_{\tau_2} w_q^N \cdot \partial_{\tau_1} A(p + \nabla w_p^N) \\ &+ \frac{1}{|Q_N|} \int_{Q_N} (q + \nabla w_q^N) \cdot \partial_{\tau_1 \tau_2}^2 A(p + \nabla w_p^N). \end{aligned}$$

La dérivée par rapport au paramètre  $\tau$  (i.e.  $k$  ou  $\lambda$ ) du correcteur  $w_p^N$  est aisée à calculer : en dérivant la formulation variationnelle satisfaite par  $w_p^N$ , on obtient que, pour toute fonction test  $\varphi \in H_{\text{per}}^1(Q_N)$  :

$$\int_{Q_N} \nabla \varphi \cdot \partial_{\tau} A(p + \nabla w_p^N) + \int_{Q_N} \nabla \varphi \cdot A \nabla \partial_{\tau} w_p^N = 0, \quad (1.44)$$

ce qui est un problème bien posé pour  $\partial_{\tau} w_p^N$ .

Pour conclure, le calcul des dérivées premières et secondes de  $A^{*,N}$  par rapport au paramètre  $\theta$  ne nécessite que la résolution d'un nombre d'équations (1.44) égal à  $\dim(\Theta)$ .

## 1.4 Perspectives

L'homogénéisation aléatoire est une théorie couvrant des cas très différents et parfois très difficile. Prise dans toute sa généralité, elle est peu accessible à l'expérience numérique. En effet, la méthode standard d'estimation de  $A^*$  est inefficace car on explore l'intégralité de l'espace d'aléa, dont la diversité reste importante lorsque la taille de la supercellule  $Q_N$  est finie.

Pour réduire la complexité du problème, une approche consiste à développer  $A^*$  par rapport à un petit paramètre  $\eta \ll 1$  (contraste, proportion de défauts, ...) présent dans la définition du champ  $A$ . On aboutit alors à des approches très peu coûteuses, mais dont le domaine de validité est *a priori* limité au régime perturbatif.

Une approche possible, dans laquelle les **Chapitre 3** et **Chapitre 4** s'inscrivent, est d'utiliser ces approches perturbatives comme un *préconditionnement*. Ceci permet d'aboutir à des méthodes efficaces dans un régime de paramètres bien plus grand que le régime perturbatif  $\eta \ll 1$ .

### 1.4.1 Comparaison des méthodes de réduction de variance

Dans cette thèse, nous avons exploré trois méthodes de réduction de variance. Nous présentons à présent une description synthétique de ces résultats.

La méthode des variables antithétique explorée dans le **Chapitre 2** présente l'avantage d'être peu invasive et n'est pas coûteuse. Nous avons observé en contrepartie une réduction de variance limitée, de l'ordre de 10 (20 dans les meilleurs cas, et 4 dans les pires). Nous renvoyons à la Section 2.3 pour des résultats détaillés.

La méthode des variables de contrôle introduite dans le **Chapitre 3** nécessite un bon modèle approché, dont la construction peut être complexe. Cependant, une fois cette étape franchie, la réduction de variance peut être bien meilleure : de l'ordre de 40 dans les cas considérés (voir la Section 3.5 pour des résultats détaillés).

La méthode de sélection introduite dans le **Chapitre 4** nécessite non seulement un bon critère de sélection, mais aussi un algorithme de sélection (qui toutefois peut-être choisi simple). Cette approche a donné les résultats les plus intéressants : une réduction de variance de l'ordre de 200 dans les cas considérés (voir la Section 5.4.2 pour des résultats détaillés).

### 1.4.2 Extension à la stationnarité continue des variables de contrôle

Une extension possible du **Chapitre 3** serait de généraliser la modélisation pour couvrir des cas de stationnarité *continue*. On remplace dans ce cas l'équation (1.3) de la Définition 1.1 par (1.5). Sous cette définition, la loi d'une variable aléatoire est invariante par toute translation et pas seulement celles d'un multiple entier de la cellule  $Q$ . Cette hypothèse modélise en particulier des milieux peu structurés, où l'emplacement d'inclusions ou de défauts, par exemple circulaires et de rayons aléatoires, est choisi au hasard sans tenir compte d'un réseau sous-jacent.

Dans le contexte du **Chapitre 3**, la méthode d'ordre 1, assez générique, peut s'étendre sans complications méthodologiques : la variable de contrôle d'ordre 1 reste la proportion volumique d'inclusions. Pour mettre en oeuvre la méthode d'ordre 2, qui seule prend en compte la géométrie, une difficulté nouvelle survient. La variable de contrôle à l'ordre 2 fait intervenir des configurations à deux défauts, lesquelles sont maintenant présentes en

un nombre *infini* (au lieu de  $N^d - 1$ ), puisque les défauts ne sont plus positionnés sur le réseau  $\mathbb{Z}^d \cap Q_N$ .

Afin de généraliser rigoureusement la méthode, il faut adapter à ce nouveau cadre le développement faiblement aléatoire introduit par [ALB12]. Il faut aussi adapter à ce nouveau cadre la méthode des bases réduites proposée dans [LBT12]. Enfin, si tous ces éléments sont réunis, il faudra adapter le développement en loi introduit dans le **Chapitre 3**, lequel pourrait alors être utilisé comme variable de contrôle.

### 1.4.3 Utilisation d'autres variables de contrôle

Une autre extension possible serait d'utiliser d'autres variables de contrôle. Pour estimer la matrice homogénéisée  $A^*$ , on fait parfois appel à des méthodes de bornes (bornes de Hashin et Shtrikman). On commence par montrer que le calcul de  $A^{*,N}$  est équivalent à la résolution d'un problème variationnel, dit problème de Hashin–Shtrikman :

$$\inf_{\tau \in L^2(Q_N)} \mathcal{E}_N(\tau, \omega).$$

Les bornes de Hashin–Shtrikman sont obtenues en choisissant des fonctions  $\tau$  particulières.

Une possibilité serait d'utiliser ces bornes (qui sont des quantités aléatoires) sur  $A^{*,N}$  en tant que variables de contrôle.

### 1.4.4 Sélection pour des problèmes nonlinéaires

Il serait intéressant d'étendre la méthode présentée dans le **Chapitre 4** à des cas nonlinéaires. Par exemple, on pourrait considérer le cas (1.13) issu du **Chapitre 2**. Les équations d'Euler-Lagrange associées s'écrivent :

$$-\operatorname{div}[\partial_\xi W(x/\varepsilon, \omega, \nabla u_\varepsilon)] = f.$$

L'équation linéarisée s'écrit sous la forme

$$-\operatorname{div}[A(x/\varepsilon, \omega) \nabla u_\varepsilon] = f.$$

Supposons par simplicité que nous connaissons complètement le coefficient homogénéisé  $A^*$  associé au problème linéaire, réputé plus facile à résoudre. Alors, il serait pertinent de sélectionner suivant le critère  $A^{*,N}(\omega) = A^*$  les réalisations pour lesquelles on va effectivement calculer  $W_N^*$  défini par (1.22).

L'efficacité de l'approche sera sans doute liée au rapport du coût entre résoudre le problème du correcteur dans le cas linéaire ou bien dans le cas nonlinéaire. Notons de plus que, dans le cas nonlinéaire, on peut être intéressé par *tout* le champ  $\xi \mapsto W^*(\xi)$ , et pas seulement l'évaluation de  $W^*(\xi)$  pour un vecteur  $\xi$  particulier. Ainsi, le calcul de  $A^*$  serait vu comme un unique calcul préalable, avant le calcul de  $W^*(\xi)$  pour de nombreux  $\xi$  différents.

### 1.4.5 Convergence du problème inverse en dimension supérieure

Dans la Section précédente, nous avons présenté l'algorithme proposé dans le **Chapitre 5** en dimension d'espace quelconque. Toutefois sa convergence n'est traitée que dans le cas unidimensionnel. Il serait donc intéressant d'introduire une méthode pour prouver la convergence de notre approximation obtenue par optimisation aux moindres carrés en plusieurs dimensions d'espace.

Dans une autre direction, il serait aussi intéressant de prouver un résultat d'homogénéisation sous des hypothèses sur  $A$  plus faibles que l'uniforme coercivité et l'existence d'une borne uniforme en  $x$  et  $\omega$ . En effet, cette hypothèse n'est pas satisfaite dans le cas que nous étudions (une variable aléatoire distribuée suivant la loi de Weibull peut prendre des valeurs aussi grandes et aussi proches de 0 que voulu). De plus, nous montrons en dimension un d'espace qu'elle peut être relaxée (cf. le Theorem 5.13). La généralisation de ce résultat en dimension quelconque reste ouverte.

#### 1.4.6 Perspectives générales dans le champ de recherche

Nous concluons cette introduction par un résumé plus général. Comment s'inscrivent les travaux de la thèse dans le champ de recherche ?

Tout d'abord, ces différentes approches montrent qu'il est possible de mettre en oeuvre des méthodes de réduction de variance en homogénéisation stochastique. Ces méthodes sont souvent efficaces, et ce d'autant plus que la méthode est sophistiquée.

Seuls des cas simples ont été abordés dans cette thèse : nous avons abordé des équations *linéaires elliptiques*, un problème variationnel *convexe*, une équation aux dérivées partielles *discrète*. Les méthodes introduites doivent être considérés comme des preuves de concept, car il est possible d'étendre les techniques à des cas plus difficiles. En outre, il est sans doute possible d'améliorer qualitativement ces techniques.

En conclusion, on pourrait imaginer à terme l'homogénéisation stochastique *par la pratique*, rendue possible par des approches auxiliaires dont pourraient faire partie les méthodes présentées. Cela permettrait d'améliorer des calculs actuellement réalisés par homogénéisation *périodique*, et ainsi d'avoir un modèle plus riche et prédictif.

## Chapter 2

# Variance reduction using antithetic variables for a nonlinear convex stochastic homogenization problem

Ce **Chapitre** reprend l'intégralité d'un article écrit en collaboration avec Frédéric Legoll et accepté dans *Discrete and Continuous Dynamical Systems - Series S* [[LM15b](#)].

Dans un cadre d'homogénéisation stochastique, nous utilisons la méthode des variables antithétiques pour obtenir un estimateur de moindre variance, donc plus précis, de quantités d'intérêt. La spécificité de notre travail par rapport aux travaux antérieurs est que le modèle est *nonlinéaire*. Nous démontrons dans certains cas que la technique des variables antithétique peut encore être utilisée, et nos expérimentations numériques en dimension 2 démontrent son efficacité.



## Variance reduction using antithetic variables for a nonlinear convex stochastic homogenization problem

Frédéric Legoll<sup>1,3</sup> and William Minvielle<sup>2,3</sup>

legoll@lami.enpc.fr, william.minvielle@cermics.enpc.fr

<sup>1</sup> Laboratoire Navier, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal,  
77455 Marne-La-Vallée Cedex 2, France,

<sup>2</sup> CERMICS, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal,  
77455 Marne-La-Vallée Cedex 2, France;

<sup>3</sup> INRIA Rocquencourt, MATHERIALS research-team, Domaine de Voluceau, B.P. 105, 78153 Le  
Chesnay Cedex, France.

**Abstract.** *We consider a nonlinear convex stochastic homogenization problem, in a stationary setting. In practice, the deterministic homogenized energy density can only be approximated by a random apparent energy density, obtained by solving the corrector problem on a truncated domain.*

*We show that the technique of antithetic variables can be used to reduce the variance of the computed quantities, and thereby decrease the computational cost at equal accuracy. This leads to an efficient approach for approximating expectations of the apparent homogenized energy density and of related quantities.*

*The efficiency of the approach is numerically illustrated on several test cases. Some elements of analysis are also provided.*

### 2.1 Introduction

In this article, we consider some theoretical and numerical questions related to variance reduction techniques for some nonlinear convex stochastic homogenization problems. In short, we show here that a technique based on antithetic variables can be used in that context, provide some elements of analysis, and demonstrate numerically the efficiency of that approach on several test cases. This work is a follow-up of the articles [BCLBL12b, BCLBL12b, CLBL10] where the same questions are considered for a *linear* elliptic equation in divergence form.

The stochastic homogenization problem we consider here writes as follows. Let  $\mathcal{D}$  be an open bounded domain of  $\mathbb{R}^d$  and  $2 \leq p < \infty$ . We consider the highly oscillatory problem

$$\inf \left\{ \int_{\mathcal{D}} W \left( \frac{x}{\varepsilon}, \omega, \nabla u(x) \right) dx - \int_{\mathcal{D}} f(x)u(x)dx, \quad u \in W_0^{1,p}(\mathcal{D}) \right\} \quad (2.1)$$

for some  $f$  and some random smooth field  $W$ , which is stationary in a sense made precise below, and satisfies some convexity and growth conditions such that, for any  $\varepsilon > 0$ ,

problem (2.1) is well-posed. See Section 2.1.1 below for a precise description of the mathematical setting, which has been introduced in [DMM86a, DMM86b]. A classical example that motivated this framework is when

$$W(y, \omega, \xi) = \frac{1}{p} a(y, \omega) |\xi|^p,$$

where  $a$  is stationary (see e.g. [DMM86a, page 382]).

In (2.1),  $\varepsilon$  denotes a supposedly small, positive constant that models the smallest possible scale present in the problem. For  $\varepsilon$  small, it is extremely expensive, in practice, to directly attack (2.1) with a numerical discretization. A useful practical approach is to *first* approximate (2.1) by its associated homogenized problem, which reads

$$\inf \left\{ \int_{\mathcal{D}} W^*(\nabla u(x)) dx - \int_{\mathcal{D}} f(x)u(x)dx, \quad u \in W_0^{1,p}(\mathcal{D}) \right\}, \quad (2.2)$$

and *next* numerically solve the latter problem. The two-fold advantage of (2.2) as compared to (2.1) is that *it is deterministic* and *it does not involve the small scale  $\varepsilon$* .

This simplification comes at a price. The homogenized energy density  $W^*$  in (2.2) is given by an integral involving a so-called corrector function, solution to a nonlinear problem (see (2.8) below for a precise formula). As most often in stochastic homogenization, this corrector problem is set on the *entire* space  $\mathbb{R}^d$ . In practice, approximations are therefore in order. A standard approach (see e.g. [BP04] in the linear setting) is to generate realizations of the energy density  $W$  over a finite, supposedly large volume at the microscale, that we denote  $Q_N$ , and approach the homogenized energy density by some empirical means using approximate correctors computed on  $Q_N$ . Although the *exact* homogenized density  $W^*$  is deterministic, its practical approximation is random, due to the truncation procedure. It is then natural to generate several realizations. However, efficiently averaging over these realizations requires to understand how variance affects the result. This is the purpose of the present article to investigate some questions in this direction, both from the theoretical and numerical standpoints.

Before proceeding and for the sake of consistency, we now present the framework of nonlinear stochastic homogenization we adopt, and make precise the questions we consider.

### 2.1.1 Homogenization theoretical setting

To begin with, we introduce the basic setting of stochastic homogenization we employ. We refer to [ES08] for a general, numerically oriented presentation, and to [BLP78, CD99, JKO94] for classical textbooks. We also refer to [LB10] and the review article [ACLB<sup>+</sup>12] (and the extensive bibliography contained therein) for a presentation of our particular setting. Throughout this article,  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space and we denote by  $\mathbb{E}(X) = \int_{\Omega} X(\omega) d\mathbb{P}(\omega)$  the expectation value of any random variable  $X \in L^1(\Omega, d\mathbb{P})$ . We next fix  $d \in \mathbb{N}^*$  (the ambient physical dimension), and assume that the group  $(\mathbb{Z}^d, +)$  acts on  $\Omega$ . We denote by  $(\tau_k)_{k \in \mathbb{Z}^d}$  this action, and assume that it preserves the measure  $\mathbb{P}$ , that is, for all  $k \in \mathbb{Z}^d$  and all  $A \in \mathcal{F}$ ,  $\mathbb{P}(\tau_k A) = \mathbb{P}(A)$ . We assume that the action  $\tau$  is *ergodic*, that is, if  $A \in \mathcal{F}$  is such that  $\tau_k A = A$  for any  $k \in \mathbb{Z}^d$ , then  $\mathbb{P}(A) = 0$  or 1. In addition, we define the following notion of stationarity (see [ACLB<sup>+</sup>12, Section 2.2]): a function  $F \in L_{\text{loc}}^1(\mathbb{R}^d, L^1(\Omega))$  is said to be *stationary* if, for all  $k \in \mathbb{Z}^d$ ,

$$F(y + k, \omega) = F(y, \tau_k \omega) \quad \text{almost everywhere and almost surely.} \quad (2.3)$$

In this setting, the ergodic theorem [Kre85, Shi84, Tem72] can be stated as follows: *Let  $F \in L^\infty(\mathbb{R}^d, L^1(\Omega))$  be a stationary random variable in the above sense. For  $k = (k_1, k_2, \dots, k_d) \in \mathbb{Z}^d$ , we set  $|k|_\infty = \sup_{1 \leq i \leq d} |k_i|$ . Then*

$$\frac{1}{(2N+1)^d} \sum_{|k|_\infty \leq N} F(y, \tau_k \omega) \xrightarrow{N \rightarrow \infty} \mathbb{E}(F(y, \cdot)) \quad \text{in } L^\infty(\mathbb{R}^d), \text{ almost surely.}$$

This implies (denoting by  $Q$  the unit cube in  $\mathbb{R}^d$ ) that

$$F\left(\frac{x}{\varepsilon}, \omega\right) \xrightarrow{\varepsilon \rightarrow 0} \mathbb{E}\left(\int_Q F(y, \cdot) dy\right) \quad \text{in } L^\infty(\mathbb{R}^d), \text{ almost surely.}$$

The purpose of the above setting is simply to formalize that, even though realizations may vary, the function  $F$  at point  $y \in \mathbb{R}^d$  and the function  $F$  at point  $y+k$ ,  $k \in \mathbb{Z}^d$ , share the same law. In the homogenization context we now turn to, this means that the local, microscopic environment (encoded in the energy density  $W$ ) is everywhere the same *on average*. From this, homogenized, macroscopic properties will follow.

We now describe more precisely the multiscale random problem (2.1). The domain  $\mathcal{D}$  is a regular (in the sense its boundaries are Lipschitz-continuous) bounded domain of  $\mathbb{R}^d$ . The right-hand side function  $f$  belongs to  $L^{p'}(\mathcal{D})$ , with  $1/p + 1/p' = 1$  (hence  $f$  is indeed in the dual space of  $L^p(\mathcal{D})$ ). For any  $\xi \in \mathbb{R}^d$ , the random field  $y, \omega \mapsto W(y, \omega, \xi)$  is assumed stationary in the sense (2.3). We assume that it is continuous (and even  $C^3$ ) with respect to the  $\xi$  variable, and that it is measurable with respect to the  $y$  argument. We also assume that there exists  $c_2 \geq c_1 > 0$  such that

$$\forall y \in \mathbb{R}^d, \quad \forall \omega \in \Omega, \quad \forall \xi \in \mathbb{R}^d, \quad c_1 |\xi|^p \leq W(y, \omega, \xi) \leq c_2 (1 + |\xi|^p). \quad (2.4)$$

Furthermore, we assume henceforth that  $W$  is *strictly convex* with respect to the argument  $\xi$ , in the sense that

$$\forall \eta \in \mathbb{R}^d, \quad \forall \xi \in \mathbb{R}^d, \quad \eta^T \partial_\xi^2 W(y, \omega, \xi) \eta > 0 \quad \text{a.e. and a.s.}, \quad (2.5)$$

where  $\partial_\xi^2 W \in \mathbb{R}^{d \times d}$  is the Hessian matrix of  $\xi \mapsto W(y, \omega, \xi)$ . A more demanding assumption is that  $W$  is  $\alpha$ -convex with respect to the argument  $\xi$ , in the sense that there exists  $\alpha > 0$  such that

$$\forall \eta \in \mathbb{R}^d, \quad \forall \xi \in \mathbb{R}^d, \quad \eta^T \partial_\xi^2 W(y, \omega, \xi) \eta \geq \alpha \eta^T \eta \quad \text{a.e. and a.s.} \quad (2.6)$$

Unless otherwise stated, we only assume (2.5) in the sequel. When needed, we will explicitly assume (2.6).

Under (2.4) and (2.5), the variational problem (2.1) is well-posed. In addition, the homogenized limit of (2.1) has been identified in [DMM86a, DMM86b] (see also [GN11, Theorem 3.1]): the unique solution  $u^\varepsilon(\cdot, \omega)$  to (2.1) converges (weakly in  $W^{1,p}(\mathcal{D})$  and strongly in  $L^p(\mathcal{D})$ , almost surely) to some deterministic function  $u^* \in W^{1,p}(\mathcal{D})$ , solution to (2.2), where the homogenized energy density  $W^*$  is given, for any  $\xi \in \mathbb{R}^d$ , by

$$W^*(\xi) = \lim_{N \rightarrow \infty} \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} W(y, \omega, \xi + \nabla w(y)) dy, \quad w \in W_0^{1,p}(Q_N) \right\} \quad (2.7)$$

where  $Q_N = (-N, N)^d$ . The convergence in (2.7) holds almost surely.

We show in Appendix 2.4 below that

$$W^*(\xi) = \lim_{N \rightarrow \infty} \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} W(y, \omega, \xi + \nabla w(y)) dy, \quad w \in W_{\#}^{1,p}(Q_N) \right\} \quad (2.8)$$

where  $W_{\#}^{1,p}(Q_N)$  denotes the set of functions that belong to  $W_{\text{loc}}^{1,p}(\mathbb{R}^d)$  and are  $Q_N$ -periodic (the only difference between (2.7) and (2.8) is thus the boundary conditions that we consider). The convergence in (2.8) again holds almost surely.

In the sequel, we work on the basis of (2.8), namely using periodic boundary conditions. We could as well, up to slight modifications, work with homogeneous Dirichlet boundary conditions, on the basis of (2.7). We choose periodic boundary conditions as they have been shown, in practice, to provide more accurate numerical results (see e.g. [KFG<sup>+</sup>03]).

### 2.1.2 The questions we consider

In practice, we cannot compute  $W^*(\xi)$ , and have to restrict ourselves to finite size domains. We therefore introduce

$$W_N^*(\omega, \xi) := \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} W(y, \omega, \xi + \nabla w(y)) dy, \quad w \in W_{\#}^{1,p}(Q_N) \right\} \quad (2.9)$$

and readily see from (2.8) that

$$W^*(\xi) = \lim_{N \rightarrow \infty} W_N^*(\omega, \xi) \text{ a.s.}$$

As briefly explained above, although  $W^*$  itself is a deterministic object, its practical approximation  $W_N^*$  is random. It is only in the limit of infinitely large domains  $Q_N$  that the deterministic value is attained. This is a standard situation in stochastic homogenization.

Many studies have been recently devoted (at least in the linear case) to establishing sharp estimates on the convergence of the random apparent homogenized quantities (computed on  $Q_N$ ) to the exact deterministic homogenized quantities. We refer e.g. to [BP04, GO12] and to the comprehensive discussion of [BCLBL12b, Section 1.2]. We take here the problem from a slightly different perspective. We observe that the error

$$W^*(\xi) - W_N^*(\omega, \xi) = \left( W^*(\xi) - \mathbb{E}[W_N^*(\cdot, \xi)] \right) + \left( \mathbb{E}[W_N^*(\cdot, \xi)] - W_N^*(\omega, \xi) \right)$$

is the sum of a systematic error (the first term in the above right-hand side) and of a statistical error (the second term in the above right-hand side). We *focus here on the statistical error*, and propose approaches to reduce the confidence interval of empirical means approximating  $\mathbb{E}[W_N^*(\cdot, \xi)]$  (or similar quantities), for a given truncated domain  $Q_N$ .

Recall that a standard technique to compute an approximation of  $\mathbb{E}[W_N^*(\cdot, \xi)]$  is to consider several independent and identically distributed realizations of the energy density  $W$ , solve for each of them the corrector problem (2.9) (thereby obtaining several i.i.d. values  $W_N^{*,m}(\omega, \xi)$ ), and proceed following a Monte Carlo approach:

$$\mathbb{E}[W_N^*(\cdot, \xi)] \approx I_{2M} := \frac{1}{2M} \sum_{m=1}^{2M} W_N^{*,m}(\omega, \xi).$$

In view of the Central Limit Theorem, we know that our quantity of interest  $\mathbb{E}[W_N^*(\cdot, \xi)]$  lies in the confidence interval

$$\left[ I_{2M} - 1.96 \frac{\sqrt{\text{Var}[W_N^*(\cdot, \xi)]}}{\sqrt{2M}}, I_{2M} + 1.96 \frac{\sqrt{\text{Var}[W_N^*(\cdot, \xi)]}}{\sqrt{2M}} \right]$$

with a probability equal to 95 %.

In this article, we show that, using a well known variance reduction technique, the technique of *antithetic variables* [Liu08, page 27], we can design a practical approach that, for finite  $N$  and any vector  $\xi$ , allows to compute a better approximation of  $\mathbb{E}[W_N^*(\cdot, \xi)]$  (and likewise for similar homogenized quantities). Otherwise stated, for an equal computational cost, the approach provides a more accurate (i.e. with a smaller confidence interval) approximation. We thereby extend to this nonlinear convex setting the results of [BCLBL12b, BCLBL12b, CLBL10] obtained in the linear case.

Our article is articulated as follows. In Section 2.2.1, we describe the proposed approach, and state our main results. The ingredients to prove these results are collected in Sections 2.2.2, 2.2.3 and 2.2.4. The actual proof of our main results is performed in Section 2.2.5. We make there several structural assumptions on the form of the energy density  $W$  to obtain these variance reduction results. In Section 2.2.6, we describe a general class of energy densities  $W$  for which our assumptions are indeed satisfied. We next turn in Section 2.3 to some illustrative numerical examples, where we demonstrate the efficiency of the approach, even in cases where the theoretical analysis is incomplete.

## 2.2 Description of the proposed approach and main results

### 2.2.1 Statement of our main results

This section is devoted to the presentation and the analysis of our approach. We first focus on estimating the expectation  $\mathbb{E}[W_N^*(\cdot, \xi)]$  of the apparent homogenized energy density (see Section 2.2.1). Our variance reduction result, Proposition 2.1, shows that the technique of antithetic variables is indeed efficient. As often the case, it is difficult to *quantitatively* assess how efficient the approach is, and this will be the purpose of the numerical tests described in Section 2.3 to address this question.

We then turn to the estimation of the first (and next second) derivatives of  $W_N^*(\cdot, \xi)$  with respect to  $\xi$ . These quantities naturally appear when one solves the convex homogenized problem (2.2) (approximating  $W^*$  by  $W_N^*(\omega, \cdot)$ ), e.g. using a Newton algorithm. For these two quantities, our result is restricted to the one-dimensional setting. See Section 2.2.1 and Proposition 2.2 for the first derivative, and Section 2.2.1 and Proposition 2.5 for the second derivative.

Sections 2.2.2, 2.2.3, 2.2.4 and 2.2.5 are devoted to the proof of the results stated here. In Section 2.2.6, we discuss an explicit class of energy densities  $W$  that falls into our framework.

### Variance reduction on the homogenized energy density

In this section, we make the following two *structure assumptions* on the rapidly oscillating field  $W$  of (2.1). First, we assume that, for any  $N$ , there exists an integer  $n$  (possibly

$n = |Q_N|$ , but not necessarily) and a function  $\mathcal{A}$ , defined on  $Q_N \times \mathbb{R}^n \times \mathbb{R}^d$ , such that the field  $W(y, \omega, \xi)$  writes

$$\forall y \in Q_N, \forall \xi \in \mathbb{R}^d, \quad W(y, \omega, \xi) = \mathcal{A}(y, X_1(\omega), \dots, X_n(\omega), \xi) \quad \text{a.s.}, \quad (2.10)$$

where  $\{X_k(\omega)\}_{1 \leq k \leq n}$  are independent scalar random variables, which are all distributed according to the uniform law  $\mathcal{U}[0, 1]$ . In general, the function  $\mathcal{A}$ , as well as the number  $n$  of independent, identically distributed variables involved in (2.10), depend on  $N$ , the size of  $Q_N$ , although this dependency is not made explicit in (2.10).

Second, we assume that the function  $\mathcal{A}$  in (2.10) is such that, for all  $y \in Q_N$  and all  $\xi \in \mathbb{R}^d$ , the map

$$(x_1, \dots, x_n) \in \mathbb{R}^n \mapsto \mathcal{A}(y, x_1, \dots, x_n, \xi) \quad (2.11)$$

is non-decreasing with respect to each of its arguments.

**Proposition 2.1.** *We assume (2.10)–(2.11). Let  $W_N^*(\omega, \xi)$  be the approximated homogenized energy density field defined by (2.9). We define on  $Q_N$  the field*

$$W^{\text{ant}}(y, \omega, \xi) := \mathcal{A}(y, 1 - X_1(\omega), \dots, 1 - X_n(\omega), \xi),$$

*antithetic to  $W$  defined by (2.10). We associate to this field the approximate homogenized energy density field  $W_N^{\text{ant},*}(\omega, \xi)$ , defined by (2.9) (replacing  $W$  by  $W^{\text{ant}}$ ). Set*

$$\widetilde{W}_N^*(\omega, \xi) := \frac{1}{2} \left( W_N^*(\omega, \xi) + W_N^{\text{ant},*}(\omega, \xi) \right). \quad (2.12)$$

*Then, for any  $\xi \in \mathbb{R}^d$ ,*

$$\mathbb{E} \left[ \widetilde{W}_N^*(\cdot, \xi) \right] = \mathbb{E} \left[ W_N^*(\cdot, \xi) \right] \quad \text{and} \quad \text{Var} \left[ \widetilde{W}_N^*(\cdot, \xi) \right] \leq \frac{1}{2} \text{Var} \left[ W_N^*(\cdot, \xi) \right]. \quad (2.13)$$

*Otherwise stated,  $\widetilde{W}_N^*(\omega, \xi)$  is a random variable which has the same expectation as  $W_N^*(\omega, \xi)$ , and its variance is smaller than half of that of  $W_N^*(\omega, \xi)$ .*

As mentioned above, this result generalizes [BCLBL12b, Proposition 2.1] to the non-linear convex variational setting considered here.

Before proceeding, we briefly explain the usefulness of the above result for variance reduction techniques. Assume we want to compute the expectation of  $W_N^*(\omega, \xi)$ , for some fixed vector  $\xi \in \mathbb{R}^d$ . Following the classical Monte-Carlo method recalled in Section 2.1.2, we estimate  $\mathbb{E} [W_N^*(\cdot, \xi)]$  by its empirical mean. To this end, we consider  $2M$  independent, identically distributed copies  $\{W_m(y, \omega, \xi)\}_{1 \leq m \leq 2M}$  of the random field  $W(y, \omega, \xi)$  on  $Q_N$ . To each copy  $W_m$ , we associate the approximate homogenized energy density  $W_N^{*,m}(\omega, \xi)$  defined by (2.9). We next introduce the empirical mean

$$I_{2M} = \frac{1}{2M} \sum_{m=1}^{2M} W_N^{*,m}(\omega, \xi), \quad (2.14)$$

and consider that, in practice, the mean  $\mathbb{E} [W_N^*(\cdot, \xi)]$  is equal to the estimator  $I_{2M}$  within

an approximate margin of error  $1.96 \frac{\sqrt{\text{Var} [W_N^*(\cdot, \xi)]}}{\sqrt{2M}}$ .

Alternate to considering (2.14), we may consider

$$\tilde{I}_{2M} = \frac{1}{M} \sum_{m=1}^M \widetilde{W}_N^{*,m}(\omega, \xi), \quad (2.15)$$

where  $\widetilde{W}_N^{*,m}$  is defined by (2.12). Again, in practice, the mean  $\mathbb{E}[W_N^*(\cdot, \xi)] = \mathbb{E}[\widetilde{W}_N^*(\cdot, \xi)]$  is equal to  $\tilde{I}_{2M}$  within an approximate margin of error  $1.96 \frac{\sqrt{\text{Var}[\widetilde{W}_N^*(\cdot, \xi)]}}{\sqrt{M}}$ . Observe now that both estimators (2.14) and (2.15) are of equal cost, since they require the same number  $2M$  of corrector problems to be solved. The accuracy of the latter is better if and only if  $\text{Var}[\widetilde{W}_N^*(\cdot, \xi)] \leq \frac{1}{2} \text{Var}[W_N^*(\cdot, \xi)]$ , which is exactly the bound (2.13) of Proposition 2.1.

### Variance reduction on the first derivative of the homogenized energy density

Restricting ourselves to the one-dimensional setting, we now state a variance reduction result for the estimation of  $\mathbb{E}[\xi \partial_\xi W_N^*(\cdot, \xi)]$ . Note that, to distinguish derivatives with respect to  $y$  from derivatives with respect to  $\xi$ , we keep the notation  $\partial_\xi W$ , even though we are in the one-dimensional situation.

We again make the structure assumption (2.10), and observe that it implies that

$$\forall y \in (-N, N), \forall \xi \in \mathbb{R}, \quad \xi \partial_\xi W(y, \omega, \xi) = \mathcal{A}_1(y, X_1(\omega), \dots, X_n(\omega), \xi) \quad \text{a.s.},$$

where  $\{X_k(\omega)\}_{1 \leq k \leq n}$  are scalar i.i.d. random variables, which are all distributed according to the uniform law  $\mathcal{U}[0, 1]$ , and where the function  $\mathcal{A}_1$ , defined on  $(-N, N) \times \mathbb{R}^n \times \mathbb{R}$ , is given by

$$\mathcal{A}_1(y, x, \xi) = \xi \partial_\xi \mathcal{A}(y, x, \xi). \quad (2.16)$$

In addition, we assume that, for all  $y \in (-N, N)$  and all  $\xi \in \mathbb{R}$ , the map

$$(x_1, \dots, x_n) \in \mathbb{R}^n \mapsto \mathcal{A}_1(y, x_1, \dots, x_n, \xi) \quad (2.17)$$

is non-decreasing with respect to each of its arguments.

We recall that the function  $\xi \mapsto W(y, \omega, \xi)$  is strictly convex (see assumption (2.5)) and satisfies (2.4). It therefore has a unique minimizer  $\xi_0(y, \omega)$ . In the sequel, we consider energy densities such that this minimizer is independent of  $y$  and  $\omega$  (see Remark 2.3 below). Without loss of generality, we can assume that  $\xi_0 = 0$ . We thus consider energy densities  $W$  such that

$$\xi \mapsto W(y, \omega, \xi) \text{ attains its minimum at } \xi = 0, \text{ a.e. and a.s.} \quad (2.18)$$

**Proposition 2.2.** *Let  $d = 1$ , and assume (2.10), (2.16), (2.17) and (2.18). We introduce*

$$\widetilde{\xi \partial_\xi W}_N^*(\omega, \xi) := \frac{1}{2} \left( \xi \partial_\xi W_N^{\text{ant},*}(\omega, \xi) + \xi \partial_\xi W_N^*(\omega, \xi) \right), \quad (2.19)$$

where  $W_N^{\text{ant},*}(\omega, \xi)$  and  $W_N^*(\omega, \xi)$  are defined as in Proposition 2.1. Then, for any  $\xi \in \mathbb{R}$ ,

$$\mathbb{E} \left[ \widetilde{\xi \partial_\xi W}_N^*(\cdot, \xi) \right] = \mathbb{E}[\xi \partial_\xi W_N^*(\cdot, \xi)] \quad \text{and} \quad \text{Var} \left[ \widetilde{\xi \partial_\xi W}_N^*(\cdot, \xi) \right] \leq \frac{1}{2} \text{Var}[\xi \partial_\xi W_N^*(\cdot, \xi)]. \quad (2.20)$$



**Remark 2.3.** *Our work is motivated by the modeling of hyperelastic materials experiencing large deformations. In this context, Assumption (2.18) amounts to assuming that there exists a natural configuration of the material. This assumption is well-known in the material science community. Many materials, but not all, indeed satisfy such an assumption. Note also that similar assumptions can be found in different but related contexts, such as homogenization of the Hamilton-Jacobi equation (see e.g. [ACS14, Equation (2.9)]).*

**Remark 2.4.** *The fact that we consider the quantity of interest  $\mathbb{E}[\xi \partial_\xi W_N^*]$  is reminiscent of our work in the linear case, that is when  $W(y, \omega, \xi) = \frac{1}{2} \xi^T A(y, \omega) \xi$  where  $A$  is a coercive bounded stationary symmetric matrix. In that case (see [BCLBL12b, BCLBL12b, CLBL10]), we proved variance reduction for the scalar quantities  $\xi^T A_N^*(\omega) \xi$  for any vector  $\xi \in \mathbb{R}^d$  (note that variance reduction was also observed for other quantities). In the non-linear setting studied in this article, we consider quantities of interest that, if the problem turns out to be linear, are equal to  $\xi^T A_N^*(\omega) \xi$  for some vector  $\xi$ . More general quantities of interest are considered in the numerical tests reported on in Section 2.3.*

### Variance reduction on the second derivative of the homogenized energy density

Considering again the one-dimensional setting as in Section 2.2.1, we eventually state a variance reduction result for the estimation of  $\mathbb{E} \left[ \partial_\xi^2 W_N^*(\cdot, \xi) \right]$ .

Recall that, for any  $y$  and  $\omega$ , the map  $\xi \mapsto \partial_\xi W(y, \omega, \xi)$  is increasing. We can therefore introduce its reciprocal function  $\zeta \mapsto \psi(y, \omega, \zeta)$ , which is also increasing.

We again make the structure assumption (2.10), and observe that it implies that, for any  $y \in (-N, N)$  and any  $\zeta \in \mathbb{R}$ ,

$$\partial_\xi^2 W(y, \omega, \psi(y, \omega, \zeta)) = \mathcal{A}_2(y, X_1(\omega), \dots, X_n(\omega), \zeta) \quad \text{a.s.},$$

where  $\{X_k(\omega)\}_{1 \leq k \leq n}$  are scalar i.i.d. random variables, which are all distributed according to the uniform law  $\mathcal{U}[0, 1]$ , and where the function  $\mathcal{A}_2$ , defined on  $(-N, N) \times \mathbb{R}^n \times \mathbb{R}$ , is given by

$$\mathcal{A}_2(y, x, \zeta) = \partial_\xi^2 \mathcal{A} \left( y, x, [\partial_\xi \mathcal{A}(y, x, \cdot)]^{-1}(\zeta) \right), \quad (2.21)$$

where  $\zeta \mapsto [\partial_\xi \mathcal{A}(y, x, \cdot)]^{-1}(\zeta)$  is the reciprocal function of  $\xi \mapsto \partial_\xi \mathcal{A}(y, x, \xi)$ .

In addition, we assume that, for all  $y \in (-N, N)$  and all  $\zeta \in \mathbb{R}$ , the map

$$(x_1, \dots, x_n) \in \mathbb{R}^n \mapsto \mathcal{A}_2(y, x_1, \dots, x_n, \zeta) \quad (2.22)$$

is non-decreasing with respect to each of its arguments.

**Proposition 2.5.** *Let  $d = 1$ , and assume (2.10), (2.16), (2.17), (2.21) and (2.22). We also assume that (2.18) holds, and that*

$$\begin{aligned} \xi \mapsto \partial_\xi^2 W(y, \omega, \xi) \text{ is non decreasing for } \xi \geq 0 \\ \text{and non increasing for } \xi \leq 0, \text{ a.e. and a.s.} \end{aligned} \quad (2.23)$$

We introduce

$$\widetilde{\partial_\xi^2 W_N^*}(\omega, \xi) := \frac{1}{2} \left( \partial_\xi^2 W_N^{\text{ant},*}(\omega, \xi) + \partial_\xi^2 W_N^*(\omega, \xi) \right),$$

where  $W_N^{\text{ant},*}(\omega, \xi)$  and  $W_N^*(\omega, \xi)$  are defined as in Proposition 2.1. Then, for any  $\xi \in \mathbb{R}$ ,

$$\mathbb{E} \left[ \widetilde{\partial_\xi^2 W_N^*}(\cdot, \xi) \right] = \mathbb{E} \left[ \partial_\xi^2 W_N^*(\cdot, \xi) \right] \quad \text{and} \quad \text{Var} \left[ \widetilde{\partial_\xi^2 W_N^*}(\cdot, \xi) \right] \leq \frac{1}{2} \text{Var} \left[ \partial_\xi^2 W_N^*(\cdot, \xi) \right]. \quad (2.24)$$

The density  $W(y, \omega, \xi) = a(y, \omega) |\xi|^p$ , where  $a$  is positive and bounded away from zero and  $p \geq 2$ , typically satisfies the assumption (2.23).



### 2.2.2 Classical results on antithetic variables

We first recall the following lemma, and refer e.g. to [BCLBL12b, Lemma 2.1] for a proof. This result is crucial for our proof of variance reduction using the technique of antithetic variables, performed in Section 2.2.5.

**Lemma 2.6** ([Liu08], page 27). *Let  $f$  and  $g$  be two real-valued functions defined on  $\mathbb{R}^n$ , which are non-decreasing with respect to each of their arguments. Consider  $X = (X_1, \dots, X_n)$  a vector of random variables, which are all independent from one another. Then*

$$\mathbb{Cov}(f(X), g(X)) \geq 0. \quad (2.25)$$

The following result is a simple consequence of the above lemma (see [BCLBL12b, Corollary 2.3] for a proof).

**Corollary 2.7** ([Liu08]). *Let  $f$  be a function defined on  $\mathbb{R}^n$ , which is non-decreasing with respect to each of its arguments. Consider  $X = (X_1, \dots, X_n)$  a vector of random variables, which are all independent from one another, and distributed according to the uniform law  $\mathcal{U}[0, 1]$ . Then*

$$\mathbb{V}\text{ar} \left( \frac{1}{2} (f(X) + f(1 - X)) \right) \leq \frac{1}{2} \mathbb{V}\text{ar} (f(X)),$$

where we denote  $1 - X = (1 - X_1, \dots, 1 - X_n) \in \mathbb{R}^n$ .

*Proof.* Choosing  $g(x_1, \dots, x_n) = -f(1 - x_1, \dots, 1 - x_n)$  in Lemma 2.6, we obtain that

$$\mathbb{Cov}(f(X), f(1 - X)) = \mathbb{Cov}(f(X_1, \dots, X_n), f(1 - X_1, \dots, 1 - X_n)) \leq 0.$$

We next observe that

$$\begin{aligned} \mathbb{V}\text{ar} \left( \frac{1}{2} (f(X) + f(1 - X)) \right) &= \frac{1}{2} \mathbb{V}\text{ar}(f(X)) + \frac{1}{2} \mathbb{Cov}(f(X), f(1 - X)) \\ &\leq \frac{1}{2} \mathbb{V}\text{ar}(f(X)), \end{aligned}$$

where we have used that  $\mathbb{V}\text{ar}(f(X)) = \mathbb{V}\text{ar}(f(1 - X))$ . □

### 2.2.3 Derivatives of the corrector and of the homogenized energy density

We now introduce the correctors as the solutions to (2.9):

$$w^N(\cdot, \omega, \xi) := \operatorname{arginf} \left\{ \int_{Q_N} W(\cdot, \omega, \xi + \nabla v), \quad v \in W_{\#}^{1,p}(Q_N), \quad \int_{Q_N} v = 0 \right\}.$$

In this section, we derive some useful expressions for the derivatives with respect to  $\xi$  of  $w^N$  and of  $W_N^*$ .

The first order optimality condition in (2.9) reads

$$\forall h \in W_{\#}^{1,p}(Q_N), \quad \int_{Q_N} (\nabla h)^T \partial_{\xi} W(\cdot, \omega, \xi + \nabla w^N) = 0. \quad (2.26)$$

We deduce from that condition that

$$\partial_{\xi} W_N^*(\omega, \xi) = \frac{1}{|Q_N|} \int_{Q_N} \partial_{\xi} W(\cdot, \omega, \xi + \nabla w^N), \quad (2.27)$$

and we note that we do not need to know  $\partial_\xi w^N$  to compute  $\partial_\xi W_N^*$ . Computing the derivative of this equality with respect to  $\xi$ , we obtain that

$$\partial_\xi^2 W_N^*(\omega, \xi) = \frac{1}{|Q_N|} \int_{Q_N} (\text{Id} + \partial_\xi \nabla w^N) \partial_\xi^2 W(\cdot, \omega, \xi + \nabla w^N) \quad (2.28)$$

with the convention that  $[\partial_\xi \nabla w^N]_{jk} = \frac{\partial^2 w^N}{\partial \xi_j \partial y_k}$  for  $1 \leq j, k \leq d$ . We can actually obtain a somewhat more symmetric expression. Computing the derivative of (2.26) with respect to  $\xi$ , we indeed see that

$$\forall h \in W_{\#}^{1,p}(Q_N), \quad \int_{Q_N} (\text{Id} + \partial_\xi \nabla w^N) \partial_\xi^2 W(\cdot, \omega, \xi + \nabla w^N) \nabla h = 0. \quad (2.29)$$

We then infer from (2.28) and (2.29) that

$$\partial_\xi^2 W_N^*(\omega, \xi) = \frac{1}{|Q_N|} \int_{Q_N} (\text{Id} + \partial_\xi \nabla w^N) \partial_\xi^2 W(\cdot, \omega, \xi + \nabla w^N) (\text{Id} + \partial_\xi \nabla w^N)^T. \quad (2.30)$$

**Remark 2.8.** *Using the same kind of arguments, we see that the function  $g_j = \frac{\partial w^N}{\partial \xi_j} \in W_{\#}^{1,p}(Q_N)$  is solution to the variational formulation*

$$\begin{aligned} \forall h \in W_{\#}^{1,p}(Q_N), \quad & \int_{Q_N} (\nabla h)^T \partial_\xi^2 W(\cdot, \omega, \xi + \nabla w^N) \nabla g_j \\ & = - \sum_{i=1}^d \int_{Q_N} \frac{\partial h}{\partial y_i} \frac{\partial^2 W}{\partial \xi_j \partial \xi_i}(\cdot, \omega, \xi + \nabla w^N). \end{aligned} \quad (2.31)$$

Suppose that  $W$  is  $\alpha$ -convex (i.e. satisfies (2.6)). Then problem (2.31) is well-posed and allows to uniquely determine (up to an additive constant)  $g_j$ , by solving a linear elliptic partial differential equation.

Combined with (2.30), this remark provides a practical way to compute  $\partial_\xi^2 W_N^*(\omega, \xi)$  without using any finite difference approximation in  $\xi$ .

We finally note that, in view of (2.27), we have

$$\xi \cdot \partial_\xi W_N^*(\omega, \xi) = \frac{1}{|Q_N|} \int_{Q_N} \xi \cdot \partial_\xi W(\cdot, \omega, \xi + \nabla w^N). \quad (2.32)$$

Likewise, in view of (2.30), we see that

$$\begin{aligned} \xi^T \partial_\xi^2 W_N^*(\omega, \xi) \xi = \\ \frac{1}{|Q_N|} \int_{Q_N} [\xi + \nabla(\xi \cdot \partial_\xi w^N)]^T \partial_\xi^2 W(\cdot, \omega, \xi + \nabla w^N) [\xi + \nabla(\xi \cdot \partial_\xi w^N)]. \end{aligned} \quad (2.33)$$

## 2.2.4 Monotonicity properties

Our goal in this section is to establish monotonicity properties for the homogenization process. Such properties are indeed useful to apply Corollary 2.7 and therefore prove variance reduction.

To simplify the notation, we assume in this section that we are in a *periodic* setting. For any  $\xi \in \mathbb{R}^d$ , the function  $y \mapsto W(y, \xi)$  is supposed to be  $Q$ -periodic (with  $Q = (0, 1)^d$ ),

to satisfy the growth condition (2.4) and to be strictly convex with respect to  $\xi$ . The associated homogenized energy density is then given by

$$W^*(\xi) = \inf \left\{ \int_Q W(y, \xi + \nabla w(y)) dy, \quad w \in W_{\#}^{1,p}(Q), \quad \int_Q w = 0 \right\}. \quad (2.34)$$

We first show a monotonicity property on the homogenized energy density in Section 2.2.4. Next, restricting ourselves to the one-dimensional setting, we show monotonicity properties for the first and the second derivative of the homogenized energy density (see respectively Sections 2.2.4 and 2.2.4).

### On the homogenized energy density

The following result is an extension to the nonlinear setting of a well-known result in the linear setting (see [Tar97, page 12]).

**Lemma 2.9.** *Suppose that the fields  $W_1$  and  $W_2$  satisfy*

$$\forall \xi \in \mathbb{R}^d, \quad W_2(y, \xi) \geq W_1(y, \xi) \text{ a.e. on } Q. \quad (2.35)$$

We denote  $W_1^*$  and  $W_2^*$  the corresponding homogenized energy densities, defined by (2.34). We then have

$$\forall \xi \in \mathbb{R}^d, \quad W_2^*(\xi) \geq W_1^*(\xi). \quad (2.36)$$

*Proof.* Fix  $\xi \in \mathbb{R}^d$ . For any  $v \in W_{\#}^{1,p}(Q)$  with  $\int_Q v = 0$ , we have that

$$W_1^*(\xi) \leq \int_Q W_1(y, \xi + \nabla v(y)) dy \leq \int_Q W_2(y, \xi + \nabla v(y)) dy.$$

Taking the infimum over  $v$ , we obtain the claimed result.  $\square$

**Remark 2.10.** *Consider the case of an energy density that is positively homogeneous of degree  $p$  with respect to its variable  $\xi$ , that is such that  $W(y, \lambda\xi) = |\lambda|^p W(y, \xi)$  for any  $y \in \mathbb{R}^d$ ,  $\xi \in \mathbb{R}^d$  and  $\lambda \in \mathbb{R}$ . A typical example is  $W(y, \xi) = \frac{1}{p} a(y) |\xi|^p$ . We then have, for any  $y$  and  $\xi$ , that*

$$\xi \cdot \partial_{\xi} W(y, \xi) = pW(y, \xi) \quad \text{and} \quad \xi^T \partial_{\xi}^2 W(y, \xi) \xi = p(p-1)W(y, \xi). \quad (2.37)$$

Using successively (2.32), (2.26) and (2.37), we obtain that

$$\begin{aligned} \xi \cdot \partial_{\xi} W^*(\xi) &= \int_Q \xi \cdot \partial_{\xi} W(\cdot, \xi + \nabla w) \\ &= \int_Q (\xi + \nabla w) \cdot \partial_{\xi} W(\cdot, \xi + \nabla w) \\ &= p \int_Q W(\cdot, \xi + \nabla w) \\ &= pW^*(\xi), \end{aligned} \quad (2.38)$$

where  $w$  is the corrector, solution to (2.34).

We next observe that, for any  $\lambda \in \mathbb{R}$  and any  $\xi \in \mathbb{R}^d$ , we have  $w(\cdot, \lambda\xi) = \lambda w(\cdot, \xi)$ . Thus, for any  $y$ , the map  $\xi \mapsto w(y, \xi)$  is homogeneous of degree one, and therefore  $\xi \cdot \partial_\xi w = w$ . We thus infer from (2.33), using (2.37), that

$$\begin{aligned} \xi^T \partial_\xi^2 W^*(\xi) \xi &= \int_Q [\xi + \nabla w]^T \partial_\xi^2 W(\cdot, \xi + \nabla w) [\xi + \nabla w] \\ &= p(p-1) \int_Q W(\cdot, \xi + \nabla w) \\ &= p(p-1) W^*(\xi). \end{aligned} \quad (2.39)$$

Consider now two fields  $W_1$  and  $W_2$  that are positively homogeneous of degree  $p$  with respect to the variable  $\xi$  and satisfy (2.35). Then we deduce from (2.36), (2.38) and (2.39) that, for all  $\xi \in \mathbb{R}^d$ ,

$$\xi \cdot \partial_\xi W_2^*(\xi) \geq \xi \cdot \partial_\xi W_1^*(\xi) \quad \text{and} \quad \xi^T \partial_\xi^2 W_2^*(\xi) \xi \geq \xi^T \partial_\xi^2 W_1^*(\xi) \xi.$$

### On the first derivative of the homogenized energy density

We now establish a monotonicity result on the derivative of  $W^*(\xi)$ , in the one-dimensional setting.

As in Section 2.2.1 (see (2.18)), we consider energy densities  $W$  such that

$$\xi \mapsto W(y, \xi) \text{ attains its minimum at } \xi = 0 \text{ for almost all } y \in Q. \quad (2.40)$$

**Lemma 2.11.** *Let  $d = 1$ , and consider two energy densities  $W_1$  and  $W_2$  satisfying (2.40), and such that*

$$\forall \xi \in \mathbb{R}, \quad \xi \partial_\xi W_2(y, \xi) \geq \xi \partial_\xi W_1(y, \xi) \text{ a.e. on } (0, 1). \quad (2.41)$$

We denote  $W_1^*$  and  $W_2^*$  the corresponding homogenized energy densities, defined by (2.34). We then have

$$\forall \xi \in \mathbb{R}, \quad \xi \partial_\xi W_2^*(\xi) \geq \xi \partial_\xi W_1^*(\xi). \quad (2.42)$$

*Proof.* We first claim that

$$\partial_\xi W^*(\xi) \text{ has the same sign as } \xi. \quad (2.43)$$

To prove this, we note that the corrector equation reads (see (2.26))

$$\frac{d}{dy} \left[ \partial_\xi W \left( y, \xi + \frac{dw}{dy}(y, \xi) \right) \right] = 0 \quad \text{on } (0, 1), \quad w(\cdot, \xi) \text{ is 1-periodic.}$$

We therefore see that  $\partial_\xi W \left( y, \xi + \frac{dw}{dy}(y, \xi) \right)$  is independent of  $y$ , and using (2.27), we obtain that

$$\partial_\xi W \left( y, \xi + \frac{dw}{dy}(y, \xi) \right) = \partial_\xi W^*(\xi) \quad \text{on } (0, 1).$$

Let  $\xi \mapsto \psi(y, \xi)$  be the reciprocal function of  $\xi \mapsto \partial_\xi W(y, \xi)$ , which exists and is increasing thanks to the strict convexity of  $\xi \mapsto W(y, \xi)$ . We deduce from the above equation, after integration over  $(0, 1)$ , that

$$\xi = \int_0^1 \psi(y, \partial_\xi W^*(\xi)) dy. \quad (2.44)$$

We are now in position to prove (2.43). Indeed, we first note that (2.40), that reads  $\partial_\xi W(y, \xi = 0) = 0$ , implies that  $\psi(y, 0) = 0$ . If  $\partial_\xi W^*(\xi) \geq 0$ , then  $\psi(y, \partial_\xi W^*(\xi)) \geq$

$\psi(y, 0) = 0$ , hence, integrating over  $(0, 1)$  and using (2.44), we obtain  $\xi \geq 0$ . Likewise,  $\partial_\xi W^*(\xi) \leq 0$  implies that  $\xi \leq 0$ . The claim (2.43) is proved.

To proceed, we see that the assumption (2.41) equivalently reads, using the reciprocal functions,

$$\forall \zeta \in \mathbb{R}, \quad \zeta \psi_2(y, \zeta) \leq \zeta \psi_1(y, \zeta) \quad \text{a.e. on } (0, 1). \quad (2.45)$$

We now prove (2.42) by contradiction. Assume that  $\xi \partial_\xi W_2^*(\xi) < \xi \partial_\xi W_1^*(\xi)$  for some  $\xi \in \mathbb{R}$ . Without loss of generality, we can assume that  $\xi > 0$ , and therefore  $\partial_\xi W_2^*(\xi) < \partial_\xi W_1^*(\xi)$ . Using (2.43), we additionally have  $0 < \partial_\xi W_2^*(\xi)$ . Using that  $\zeta \mapsto \psi_2(y, \zeta)$  is increasing and (2.45) with  $\zeta = \partial_\xi W_1^*(\xi) > 0$ , we have

$$\psi_2(y, \partial_\xi W_2^*(\xi)) < \psi_2(y, \partial_\xi W_1^*(\xi)) \leq \psi_1(y, \partial_\xi W_1^*(\xi)).$$

Integrating over  $(0, 1)$  and using (2.44) yields

$$\xi = \int_0^1 \psi_2(y, \partial_\xi W_2^*(\xi)) dy < \int_0^1 \psi_1(y, \partial_\xi W_1^*(\xi)) dy = \xi,$$

and we reach a contradiction. This concludes the proof.  $\square$

### On the second derivative of the homogenized energy density

We next turn to monotonicity properties of the second derivative of the homogenized energy density. As in Section 2.2.4, we consider energy densities satisfying (2.40). In the spirit of (2.23), we additionally request that, almost everywhere in  $(0, 1)$ ,

$$\begin{aligned} \xi \mapsto \partial_\xi^2 W(y, \xi) \text{ is non decreasing for } \xi \geq 0 \\ \text{and non increasing for } \xi \leq 0. \end{aligned} \quad (2.46)$$

**Lemma 2.12.** *Let  $d = 1$ , and consider two energy densities  $W_1$  and  $W_2$  satisfying (2.40), (2.41), (2.46) and such that*

$$\forall \zeta \in \mathbb{R}, \quad \partial_\xi^2 W_2(y, \psi_2(y, \zeta)) \geq \partial_\xi^2 W_1(y, \psi_1(y, \zeta)) \quad \text{a.e. on } (0, 1). \quad (2.47)$$

We denote  $W_1^*$  and  $W_2^*$  the corresponding homogenized energy densities, defined by (2.34). We then have

$$\forall \xi \in \mathbb{R}, \quad \partial_\xi^2 W_2^*(\xi) \geq \partial_\xi^2 W_1^*(\xi). \quad (2.48)$$

We recall that  $\zeta \mapsto \psi(y, \zeta)$  is the reciprocal function of  $\xi \mapsto \partial_\xi W(y, \xi)$ .

*Proof.* We first compute the derivative of (2.44) and obtain

$$\frac{1}{\partial_\xi^2 W^*(\xi)} = \int_0^1 \frac{dy}{\partial_\xi^2 W[y, \psi(y, \partial_\xi W^*(\xi))]} \quad (2.49)$$

It is sufficient to prove (2.48) for  $\xi > 0$ . Using (2.42) and the fact that  $\psi_1$  and  $\partial_\xi^2 W_1$  are non-decreasing with respect to their second argument, we have

$$\partial_\xi^2 W_1(y, \psi_1(y, \partial_\xi W_1^*(\xi))) \leq \partial_\xi^2 W_1(y, \psi_1(y, \partial_\xi W_2^*(\xi))).$$

Using (2.47) for  $\zeta = \partial_\xi W_2^*(\xi)$ , we deduce that

$$\partial_\xi^2 W_1(y, \psi_1(y, \partial_\xi W_1^*(\xi))) \leq \partial_\xi^2 W_2(y, \psi_2(y, \partial_\xi W_2^*(\xi))).$$

In view of (2.49), this inequality readily implies (2.48) for  $\xi > 0$ . This concludes the proof.  $\square$

### 2.2.5 Proof of Propositions 2.1, 2.2 and 2.5

Now that we have collected all the necessary ingredients, we are in position to prove our main results.

#### Variance reduction on the homogenized energy density

*Proof of Proposition 2.1.* As  $1 - X_k(\omega)$  and  $X_k(\omega)$  share the same law, so do the fields  $W$  and  $W^{\text{ant}}$  on  $Q_N$ . Hence, the homogenized fields  $W_N^*(\omega, \xi)$  and  $W_N^{\text{ant},*}(\omega, \xi)$  share the same law, and we obtain the first assertion of (2.13).

We now choose a vector  $\xi \in \mathbb{R}^d$ , and denote by  $\mathcal{P}_N^\xi$  the operator that associates to a given  $Q_N$ -periodic energy density the homogenized energy density evaluated at  $\xi$ . We see from (2.9) that  $W_N^*(\omega, \xi)$  is the effective energy density (evaluated at  $\xi$ ) obtained by periodic homogenization of  $W|_{y \in Q_N}$ :

$$W_N^*(\omega, \xi) = \mathcal{P}_N^\xi [W(\cdot, \omega, \cdot)|_{y \in Q_N}] \quad \text{a.s.} \quad (2.50)$$

Using the function  $\mathcal{A}$  of (2.10), we introduce the map

$$\begin{aligned} f : \mathbb{R}^n &\rightarrow \mathbb{R} \\ x &\mapsto \mathcal{P}_N^\xi [\mathcal{A}(\cdot, x, \cdot)], \end{aligned}$$

see that  $f(X(\omega)) = W_N^*(\omega, \xi)$  and that, using the definition (2.12) of  $\widetilde{W}_N^*(\omega, \xi)$ , we have

$$\frac{1}{2} (f(X(\omega)) + f(1 - X(\omega))) = \frac{1}{2} (W_N^*(\omega, \xi) + W_N^{\text{ant},*}(\omega, \xi)) = \widetilde{W}_N^*(\omega, \xi). \quad (2.51)$$

We have used above the notation  $1 - X = (1 - X_1, \dots, 1 - X_n) \in \mathbb{R}^n$  introduced in Corollary 2.7.

We know from Assumption (2.11) that, for any  $y \in Q_N$  and any  $\zeta \in \mathbb{R}^d$ , the function  $\mathcal{A}(y, \cdot, \zeta)$  is non-decreasing with respect to each of its arguments. In view of Lemma 2.9, we obtain that  $f$  is non-decreasing.

We are thus in position to use Corollary 2.7, which yields

$$\text{Var} \left( \frac{1}{2} (f(X) + f(1 - X)) \right) \leq \frac{1}{2} \text{Var} (f(X)).$$

Using (2.51), we obtain

$$\text{Var} \left( \widetilde{W}_N^*(\cdot, \xi) \right) = \text{Var} \left[ \frac{1}{2} (f(X) + f(1 - X)) \right] \leq \frac{1}{2} \text{Var} (f(X)) = \frac{1}{2} \text{Var} (W_N^*(\cdot, \xi)),$$

which concludes the proof of the second assertion of (2.13) and of Proposition 2.1.  $\square$

**Remark 2.13.** *Following Remark 2.10, consider a positively homogeneous energy density  $W$ . We have shown there that  $\xi \cdot \partial_\xi W_N^*(\omega, \xi)$  and  $\xi^T \partial_\xi^2 W_N^*(\omega, \xi) \xi$  are equal (up to a deterministic multiplicative constant) to  $W_N^*(\omega, \xi)$ . Thus, under Assumptions (2.10)–(2.11), variance reduction holds for these two outputs as well.*

### Variance reduction on the first derivative of the homogenized energy density

*Proof of Proposition 2.2.* The proof follows the same lines as that of Proposition 2.1.

As  $1 - X_k(\omega)$  and  $X_k(\omega)$  share the same law, so do the fields  $W$  and  $W^{\text{ant}}$  on  $Q_N$ . Hence, the quantities  $\xi \partial_\xi W_N^*(\omega, \xi)$  and  $\xi \partial_\xi W_N^{\text{ant},*}(\omega, \xi)$  share the same law, which implies the first assertion of (2.20).

To prove the second assertion, we again make use, as in the proof of Proposition 2.1, of the operator  $\mathcal{P}_N^\xi$  that associates to a given  $Q_N$ -periodic energy density the homogenized energy density evaluated at  $\xi$  (here,  $Q_N = (-N, N)$ ). Expression (2.50) holds. Choosing a vector  $\xi \in \mathbb{R}$ , we introduce the function

$$\begin{aligned} f : \mathbb{R}^n &\rightarrow \mathbb{R} \\ x &\mapsto \xi \partial_\xi \left[ \mathcal{P}_N^\xi (\mathcal{A}(\cdot, x, \cdot)) \right], \end{aligned}$$

which obviously satisfies  $f(X(\omega)) = \xi \partial_\xi W_N^*(\omega, \xi)$ . Using the definition (2.19) of  $\widetilde{\xi \partial_\xi W_N^*}$ , we have

$$\frac{1}{2} [f(X(\omega)) + f(1 - X(\omega))] = \frac{1}{2} \left[ \xi \partial_\xi W_N^*(\omega, \xi) + \xi \partial_\xi W_N^{\text{ant},*}(\omega, \xi) \right] = \widetilde{\xi \partial_\xi W_N^*}(\omega, \xi). \quad (2.52)$$

Using (2.16) and (2.17), we infer from Lemma 2.11 that  $f$  is non-decreasing.

Using Corollary 2.7, we write that  $\text{Var} \left( \frac{1}{2} (f(X) + f(1 - X)) \right) \leq \frac{1}{2} \text{Var} (f(X))$ . In view of (2.52), we recast this inequality as

$$\text{Var} \left[ \widetilde{\xi \partial_\xi W_N^*}(\cdot, \xi) \right] \leq \frac{1}{2} \text{Var} (\xi \partial_\xi W_N^*(\cdot, \xi)),$$

and therefore obtain the second assertion of (2.20). This concludes the proof of Proposition 2.2.  $\square$

### Variance reduction on the second derivative of the homogenized energy density

*Proof of Proposition 2.5.* The proof follows the same lines as the proof of Proposition 2.2. Using Assumptions (2.17) and (2.22), we see that Assumptions (2.41) and (2.47) of Lemma 2.12 are satisfied. The monotonicity result of Lemma 2.12 next allows to use Corollary 2.7, which implies (2.24).  $\square$

#### 2.2.6 Examples satisfying our structure assumptions

Before proceeding to the numerical tests, we give here some specific examples of fields  $W$  that satisfy the above assumptions. We consider the case

$$W(y, \omega, \xi) = a(y, \omega) \frac{|\xi|^p}{p} + c(y, \omega) \frac{|\xi|^2}{2}, \quad p \geq 2, \quad (2.53)$$

with  $c(y, \omega) \geq 0$  and  $a(y, \omega) \geq a_- > 0$  a.e. and a.s., and provide sufficient conditions on the scalar fields  $a$  and  $c$  for the structure assumptions (2.10), (2.11), (2.17) and (2.22) to be satisfied. Note that (2.18) and (2.23) are already fulfilled.

Consider two families  $(a_k(\omega))_{k \in \mathbb{Z}^d}$  and  $(c_k(\omega))_{k \in \mathbb{Z}^d}$  of independent, identically distributed random variables, and assume that

$$a(y, \omega) = \sum_{k \in \mathbb{Z}^d} \mathbf{1}_{Q+k}(y) a_k(\omega), \quad c(y, \omega) = \sum_{k \in \mathbb{Z}^d} \mathbf{1}_{Q+k}(y) c_k(\omega), \quad (2.54)$$

where  $Q = (0, 1)^d$  and  $Q + k$  is the cube  $Q$  translated by the vector  $k \in \mathbb{Z}^d$ . The scalar field  $a(y, \omega)$  is therefore constant in each cube  $Q + k$  with i.i.d. values  $a_k(\omega)$ , and likewise for  $c(y, \omega)$ .

We assume that there exist  $\alpha > 0$  and  $\beta < \infty$  such that, for all  $k \in \mathbb{Z}^d$ ,  $0 < \alpha \leq a_k(\omega) \leq \beta < +\infty$  and  $0 \leq c_k(\omega) \leq \beta < +\infty$  almost surely. Consequently, (2.4) holds.

Introduce now the cumulative distribution functions  $P_a(x) = \nu_a(-\infty, x)$ , where  $\nu_a$  is the common probability measure of all the  $a_k$ , and next the non-decreasing functions  $f_a(x) = \inf\{z; P_a(x) \geq z\}$ . Then, for any random variable  $X^a(\omega)$  uniformly distributed in  $[0, 1]$ , the random variable  $f_a(X^a(\omega))$  is distributed according to the measure  $\nu_a$ . As a consequence, we can recast (2.54) in the form

$$a(y, \omega) = \sum_{k \in \mathbb{Z}^d} \mathbf{1}_{Q+k}(y) f_a(X_k^a(\omega)),$$

where  $(X_k^a(\omega))_{k \in \mathbb{Z}^d}$  is a family of independent random variables that are all uniformly distributed in  $[0, 1]$ , and  $f_a$  is non-decreasing. We can proceed likewise for the variables  $c_k$ . This yields an example where (2.10), (2.11) and (2.17) hold. In particular, the function  $\mathcal{A}$  of (2.10) reads

$$\mathcal{A}(y, x_a, x_c, \xi) = \frac{|\xi|^p}{p} \sum_{k \in I_N} \mathbf{1}_{Q+k}(y) f_a(x_k^a) + \frac{|\xi|^2}{2} \sum_{k \in I_N} \mathbf{1}_{Q+k}(y) f_c(x_k^c),$$

where  $I_N = \{k \in \mathbb{Z}^d \text{ s.t. } Q + k \subset Q_N\}$  and  $x_a = \{x_k^a\}_{k \in I_N}$ . As shown in [BCLBL12b], more general fields  $a(y, \omega)$  (where random variables may be correlated) also fall into this framework.

In what follows, we prove that, under assumptions (2.53) and (2.54), and if  $p \leq 3$ , the structure assumption (2.22) holds. Without loss of generality, we may assume that  $y \in (0, 1)$ , and write that

$$\forall y \in (0, 1), \quad \mathcal{A}(x_a, x_c, \xi) = \bar{a} \frac{|\xi|^p}{p} + \bar{c} \frac{|\xi|^2}{2},$$

with  $\bar{a} = f_a(x_0^a)$  and  $\bar{c} = f_c(x_0^c)$ . By a slight abuse of notation, we keep implicit the dependency with respect to  $y$ , work with  $\bar{a}$  and  $\bar{c}$  rather than  $x_a$  and  $x_c$ , and write

$$\mathcal{A}(\bar{a}, \bar{c}, \xi) = \bar{a} \frac{|\xi|^p}{p} + \bar{c} \frac{|\xi|^2}{2}.$$

We compute

$$\partial_\xi \mathcal{A}(\bar{a}, \bar{c}, \xi) = \bar{a} |\xi|^{p-2} \xi + \bar{c} \xi$$

and denote  $\zeta \mapsto g(\bar{a}, \bar{c}, \zeta)$  the reciprocal to the function  $\xi \mapsto \partial_\xi \mathcal{A}(\bar{a}, \bar{c}, \xi)$ :

$$\zeta = \bar{a} |g(\bar{a}, \bar{c}, \zeta)|^{p-2} g(\bar{a}, \bar{c}, \zeta) + \bar{c} g(\bar{a}, \bar{c}, \zeta).$$

The function  $\mathcal{A}_2$  of (2.22) then reads

$$\mathcal{A}_2(\bar{a}, \bar{c}, \zeta) = (p-1) \bar{a} |g(\bar{a}, \bar{c}, \zeta)|^{p-2} + \bar{c}.$$

We are left with showing that  $\mathcal{A}_2$  is non-decreasing with respect to  $\bar{a}$  and  $\bar{c}$ .

A first remark is that since  $g(\bar{a}, \bar{c}, \zeta)$  has the same sign as  $\zeta$  (recall that  $\bar{a} > 0$  and  $\bar{c} \geq 0$ ), we may as well restrict ourselves to  $\zeta > 0$  and  $g(\bar{a}, \bar{c}, \zeta) > 0$ . We hence have

$$\begin{aligned} \mathcal{A}_2(\bar{a}, \bar{c}, \zeta) &= (p-1) \bar{a} g(\bar{a}, \bar{c}, \zeta)^{p-2} + \bar{c}, \\ \zeta &= \bar{a} g(\bar{a}, \bar{c}, \zeta)^{p-1} + \bar{c} g(\bar{a}, \bar{c}, \zeta). \end{aligned} \tag{2.55}$$



We first compute the derivative of  $\mathcal{A}_2$  with respect to  $\bar{a}$ :

$$\frac{\partial \mathcal{A}_2}{\partial \bar{a}} = (p-1)g(\bar{a}, \bar{c}, \zeta)^{p-2} + (p-1)(p-2)\bar{a}g(\bar{a}, \bar{c}, \zeta)^{p-3} \frac{\partial g}{\partial \bar{a}}.$$

Using (2.55) to compute  $\frac{\partial g}{\partial \bar{a}}$ , we obtain that

$$(\bar{c} + (p-1)\bar{a}g^{p-2}) \frac{\partial \mathcal{A}_2}{\partial \bar{a}} = (p-1)\bar{c}g^{p-2} + \bar{a}(p-1)g^{2p-4},$$

and since  $p > 1$ ,  $\bar{a} > 0$  and  $g > 0$ , we deduce that  $\frac{\partial \mathcal{A}_2}{\partial \bar{a}} \geq 0$ .

We next compute the derivative of  $\mathcal{A}_2$  with respect to  $\bar{c}$ . Using again (2.55) to compute  $\frac{\partial g}{\partial \bar{c}}$ , we obtain that

$$(\bar{c} + (p-1)\bar{a}g^{p-2}) \frac{\partial \mathcal{A}_2}{\partial \bar{c}} = \bar{c} - (p-1)(p-3)\bar{a}g^{p-2}.$$

Recall that  $\bar{a} > 0$ ,  $\bar{c} \geq 0$ ,  $p > 1$  and  $g > 0$ . We have assumed that  $p \leq 3$ , and therefore deduce from the above relation that  $\frac{\partial \mathcal{A}_2}{\partial \bar{c}} \geq 0$ . The structure assumption (2.22) hence holds in that case.

**Remark 2.14.** *The argument above also shows that the case*

$$W(y, \omega, \xi) = a(y, \omega) \frac{|\xi|^p}{p},$$

*along with assumption (2.54), falls into our framework for any  $p \geq 2$ .*

*It is likely that other settings, such as*

$$W(y, \omega, \xi) = (a(y, \omega) + c(y, \omega)) \frac{|\xi|^p}{p} + c(y, \omega) \frac{|\xi|^2}{2},$$

*along with assumption (2.54), where  $a_k$  and  $c_k$  are all independent random variables, also fall into our framework. We do not pursue in this direction here.*

## 2.3 Numerical results

Our numerical experiments are presented in Section 2.3.2, and discussed in details in the subsequent sections. In Section 2.3.1, we first discuss the algorithm we used to solve the variational problem (2.9) that defines the apparent homogenized energy density.

### 2.3.1 Newton algorithm to solve the truncated corrector problem

As mentioned above, the corrector problem (2.9) is a convex minimization problem, which has been well studied in the literature (see e.g. [BL93, Cho89, GM75, LT94]). We explain here how we proceed in practice to solve this problem, assuming that  $W$  is not only strictly convex, but actually  $\alpha$ -convex (i.e. satisfies (2.6)).

To simplify our exposition, we use the notation of the  $Q$ -periodic case, where the corrector problem is (2.34). We introduce some basis functions  $\{\varphi_i\}_{i \in I}$  (e.g. finite element

functions) where  $\varphi_i \in W_{\#}^{1,p}(Q)$ , and the finite dimensional space  $V_h = \text{Span}\{\varphi_i, i \in I\}$ . Consider the functional

$$J(w) = J(\{w_i\}_{i \in I}) = \int_Q W(y, \xi + \nabla w(y)) dy$$

defined on  $V_h$ , with

$$w(y) = \sum_{i \in I} w_i \varphi_i(y),$$

and the variational problem

$$\inf_{v_h \in V_h} J(v_h). \quad (2.56)$$

This problem has a unique solution (denoted  $w_h \in V_h$ ) up to the addition of a constant. The quantity  $\nabla w_h$  is well-defined, and is the finite-dimensional approximation of  $\nabla w$ , where  $w$  is the solution to (2.34).

In practice, problem (2.56) is solved using a Newton algorithm. We see that

$$\frac{\partial J}{\partial w_j}(w) = D_w(\varphi_j) \quad \text{and} \quad \frac{\partial^2 J}{\partial w_j \partial w_k}(w) = H_w(\varphi_j, \varphi_k)$$

where

$$D_w(\varphi) = \int_Q \nabla \varphi(y) \cdot \partial_{\xi} W(y, \xi + \nabla w(y)) dy$$

and

$$H_w(\varphi, \psi) = \int_Q (\nabla \varphi(y))^T \partial_{\xi}^2 W(y, \xi + \nabla w(y)) \nabla \psi(y) dy.$$

The Newton algorithm consists in defining  $w_h^{m+1} \in V_h$  from  $w_h^m \in V_h$  by the following linear elliptic problem: find  $w_h^{m+1} \in V_h$  such that

$$\forall \theta \in V_h, \quad H_{w_h^m}(w_h^{m+1} - w_h^m, \theta) = -D_{w_h^m}(\theta).$$

Again,  $w_h^{m+1}$  is uniquely defined up to the addition of a constant.

The finite-dimensional problem (2.56) is  $\alpha$ -convex, and  $W$  is smooth with respect to  $\xi$ : the Newton algorithm hence locally converges (quadratically), and  $\lim_{m \rightarrow \infty} \nabla w_h^m = \nabla w_h$ .

In practice, we consider a sequence  $\mathcal{T}_h$  of meshes on  $Q$ , and set  $V_h = \mathbb{P}_h^1(Q) = \{v_h \in C(Q) \text{ s.t. } \forall T \in \mathcal{T}_h, v_h \text{ is affine on } T\}$ . By classical finite element results, we know that  $\lim_{h \rightarrow 0} \|\nabla w_h - \nabla w\|_{L^p(Q)} = 0$  (see e.g. [BL93] and also [AV12, Tho97]).

### 2.3.2 Overview of numerical results

We have considered three test-cases of the form (2.53)–(2.54), namely

$$W(y, \omega, \xi) = a(y, \omega) \frac{|\xi|^p}{p} + c(y, \omega) \frac{|\xi|^2}{2}$$

with  $a(y, \omega) = \sum_{k \in \mathbb{Z}^d} \mathbf{1}_{Q+k}(y) a_k(\omega)$  and  $c(y, \omega) = \sum_{k \in \mathbb{Z}^d} \mathbf{1}_{Q+k}(y) c_k(\omega)$ ,

with  $p = 4$ , in dimension  $d = 2$ . The random variables  $a_k$  follow a Bernoulli distribution:  $\mathbb{P}(a_k = \alpha) = \mathbb{P}(a_k = \beta) = 1/2$ , with  $\alpha = 3$  and  $\beta = 23$ . The value of the field  $c$  is chosen as follows:

- Test Case 1: in this first test case,  $c(y, \omega) = 0$ . The problem is thus strictly convex but not  $\alpha$ -convex. In addition, the energy density is positively homogeneous of degree  $p$ , hence Remarks 2.10 and 2.13 apply.
- Test Case 2: the second test case corresponds to  $c(y, \omega) = 1$ . The problem is then  $\alpha$ -convex, and highly oscillatory only in its non-harmonic component.
- Test Case 3: for the third test case, we work with  $c(y, \omega)$  chosen according to (2.54), where  $\mathbb{P}(c_k = \gamma) = \mathbb{P}(c_k = \delta) = 1/2$ , with  $\gamma = 1$  and  $\delta = 3$ . The problem is thus highly oscillatory both in its non-harmonic and its harmonic components.

We take the meshsize  $h = 0.2$ . The Newton algorithm is initialized with the solution  $w_0$  to

$$-\operatorname{div} [(a(y, \omega) + c(y, \omega))(\xi + \nabla w_0)] = 0 \quad \text{in } Q_N, \quad w_0 \text{ is } Q_N\text{-periodic,}$$

and the iterations stop when  $\frac{\|w_h^{n+1} - w_h^n\|_{W^{1,p}}}{\|w_h^n\|_{W^{1,p}}} \leq \mathbf{tol}$ . If  $\mathbf{tol}$  is chosen too large, then (2.56)

is inaccurately solved, and the variance reduction is not very good. For our numerical tests, we set  $\mathbf{tol} = 10^{-5}$ : the discrete problem (2.56) is accurately solved, while only a limited number of iterations (in practice, around 5 iterations) are needed.

For the numerical tests, we adopt the convention that  $Q_N = (-N, N)^2$ . For each  $Q_N$ , the standard Monte Carlo results have been obtained using  $2M = 100$  realizations (from which we build the empirical estimator (2.14)). For the antithetic variable approach, we have also solved  $2M$  corrector problems, from which we build the empirical estimator (2.15). Therefore, in all what follows, we compare the accuracy of the Monte Carlo approach (MC) and the Antithetic Variable approach (AV) at *equal computational cost*.

### 2.3.3 Test Case 1

In this test case, the energy density is positively homogeneous. We therefore know, from Proposition 2.1 and Remark 2.13, that our approach yields estimations of the expectation of  $W_N^*(\omega, \xi)$ ,  $\xi \cdot \partial_\xi W_N^*(\omega, \xi)$  and  $\xi^T \partial_\xi^2 W_N^*(\omega, \xi) \xi$  with a smaller variance than the standard Monte Carlo approach. Our aim here is to *quantify* the efficiency gain. Note also that we have *not* taken into account, in our implementation, the fact that  $W_N^*(\omega, \xi)$ ,  $\xi \cdot \partial_\xi W_N^*(\omega, \xi)$  and  $\xi^T \partial_\xi^2 W_N^*(\omega, \xi) \xi$  are here proportional to one another.

To begin with, we show on Figure 2.1 the estimation by empirical means (along with a 95 % confidence interval) of three quantities (the homogenized energy density, its first derivative with respect to  $\xi_1$  and its second derivative with respect to  $\xi_1$  and  $\xi_2$ ; we refer to [LM13] for more comprehensive numerical results). We observe that the variance of all quantities decreases when the size of  $Q_N$  increases, and that confidence intervals obtained with the antithetic variable approach are smaller than those obtained with a standard Monte Carlo approach, for an equal computational cost.

We now turn to a more quantitative analysis of the variance. Figure 2.2 shows the variances

$$V_{\text{MC}} = \frac{1}{2} \operatorname{Var} [W_N^*(\cdot, \xi)] \quad \text{and} \quad V_{\text{AV}} = \operatorname{Var} [\widetilde{W}_N^*(\cdot, \xi)] \quad (2.57)$$

as a function of  $N$  (note the factor  $1/2$  in the definition of  $V_{\text{MC}}$ , consistent with (2.13), (2.14) and (2.15)). We observe that the variance of any of our quantities of interest (obtained either with the Monte Carlo approach or the Antithetic Variable approach) decreases at the rate  $1/|Q_N|$  as  $N$  increases (as expected if one could use the Central Limit Theorem).

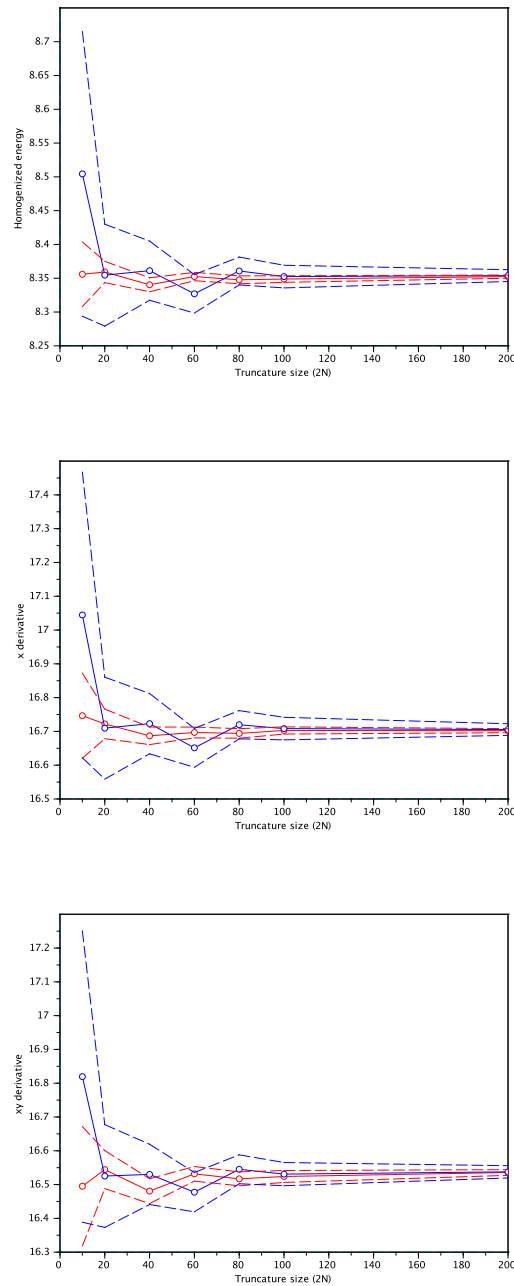


Figure 2.1 – Test-Case 1: Homogenized quantities as a function of  $2N$ , for the vector  $\xi = (1, 1)^T$  (Blue: Monte Carlo results; Red: Antithetic variable approach; Dashed lines: 95% confidence interval, equating the cost of the two approaches). From top to bottom: estimation of  $\mathbb{E}[W_N^*(\cdot, \xi)]$ ,  $\mathbb{E}[\partial_{\xi_1} W_N^*(\cdot, \xi)]$  and  $\mathbb{E}[\partial_{\xi_1 \xi_2}^2 W_N^*(\cdot, \xi)]$ .

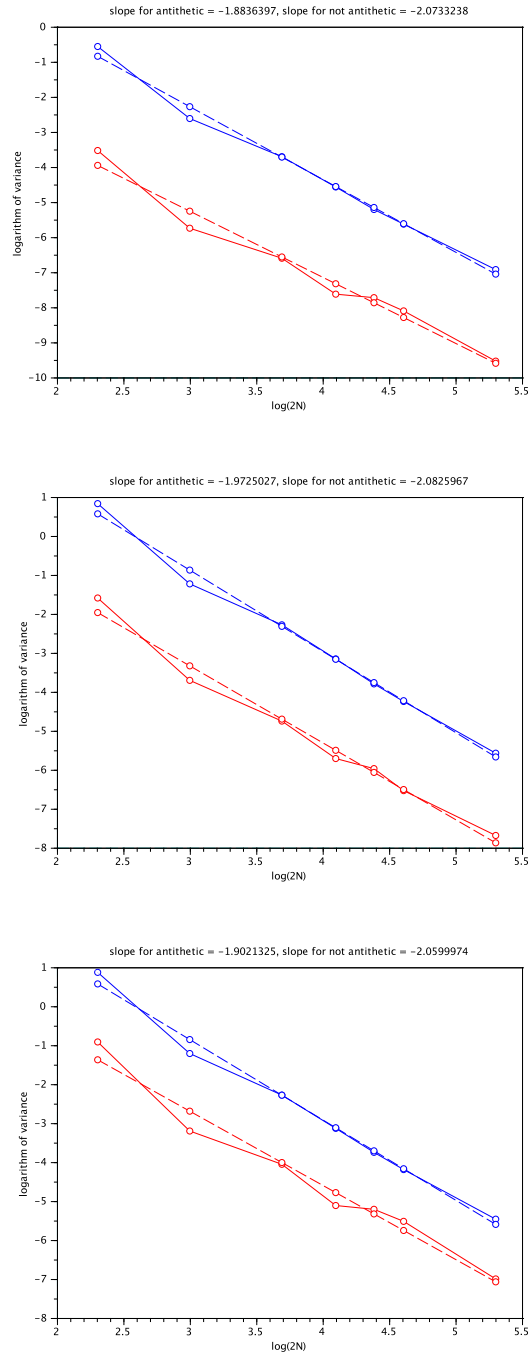


Figure 2.2 – Test-Case 1: Variances (2.57) of the same quantities of interest as on Figure 2.1, as a function of  $2N$  (Blue: Monte Carlo approach; Red: Antithetic Variable approach; Natural logarithm plot). Solid line: actual results. Dashed line: linear regression fit.

We also observe that the variance obtained with our approach is systematically smaller than the Monte Carlo variance, in the sense that  $V_{AV} \leq V_{MC}$ .

We next report on Table 2.1 the variance reduction ratio

$$R = \frac{V_{MC}}{V_{AV}} = \frac{\text{Var}[W_N^*(\cdot, \xi)]}{2\text{Var}[\widetilde{W}_N^*(\cdot, \xi)]}, \quad (2.58)$$

which measures the gain in computational cost at equal accuracy, or the square of the accuracy gain at equal computational cost. We report this ratio for several quantities of interest. Although this ratio somewhat varies with  $N$ , we observe that it is of the order of 10 for all quantities of interest, except for  $\partial_{\xi_1 \xi_2}^2 W_N^*$ , for which it is always larger than 4. In particular, even if  $N$  is not large (because we cannot afford to work on a large domain  $Q_N$ ), we still observe variance reduction.

$2N$	$W_N^*$	$\partial_{\xi_1} W_N^*$	$\partial_{\xi_2} W_N^*$	$\partial_{\xi_1 \xi_1}^2 W_N^*$	$\partial_{\xi_1 \xi_2}^2 W_N^*$	$\partial_{\xi_2 \xi_2}^2 W_N^*$	$\xi \cdot \partial_{\xi} W_N^*$	$\xi^T \partial_{\xi}^2 W_N^* \xi$
10	19.41	11.26	13.86	9.846	5.966	13.34	19.39	19.41
20	22.82	11.89	13.03	9.865	7.306	9.096	22.77	22.83
40	18.08	11.82	9.816	9.576	5.904	8.831	18.03	18.11
60	21.26	12.89	12.98	10.57	7.247	10.73	21.24	21.28
80	12.36	8.798	9.050	10.05	4.316	8.454	12.31	12.37
100	11.88	9.856	8.412	11.10	3.775	10.24	11.82	11.88
200	13.60	8.261	11.52	8.057	4.636	12.62	13.54	13.61

Table 2.1 – Test-Case 1: Variance reduction ratios (2.58).

**Remark 2.15.** *Similar variance reduction ratios are obtained in the case when the corrector problem is supplemented with homogeneous Dirichlet boundary conditions on the boundary on  $Q_N$  (in the spirit of (2.7)), rather than periodic boundary conditions as used here following (2.9) (results not shown).*

### 2.3.4 Test Case 2

We now consider a test-case for which the energy density is not positively homogeneous. From our results of Section 2.2.1, we know that our approach yields variance reduction for the estimation of  $\mathbb{E}[W_N^*(\cdot, \xi)]$ . Our aim here is two-fold: we first quantify the efficiency gain, and we next verify (and this will indeed be the case) that we also obtain a gain in efficiency for quantities of interest (such as the first or second derivatives of  $W_N^*(\omega, \xi)$  with respect to  $\xi$ ) for which we do not have theoretical results in the two-dimensional case.

We show on Figure 2.3 the variances (2.57) of the same quantities of interest as on Figures 2.1 and 2.2 (obtained either with the Monte Carlo approach or the Antithetic Variable approach). As for the previous test-case, we observe that all variances decrease at the rate  $1/|Q_N|$  as  $N$  increases. In addition, we observe that the variance obtained with our approach is systematically smaller than the Monte Carlo variance, in the sense that  $V_{AV} \leq V_{MC}$ .

On Table 2.2, we report the variance reduction ratios (2.58) (with the same convention as in Table 2.1). We observe an efficiency gain of more than 10 for all quantities of interest, except again the cross derivative  $\partial_{\xi_1 \xi_2}^2 W_N^*$ , for which the gain is smaller, and of the order of 4.

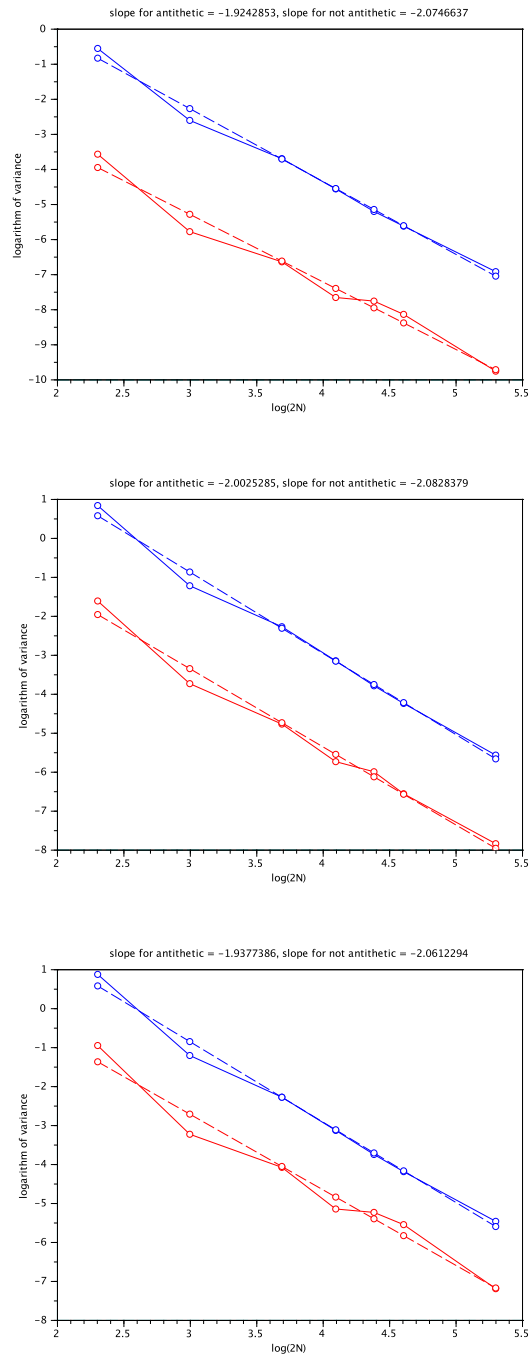


Figure 2.3 – Test-Case 2: Variances (2.57) as a function of  $2N$  (Blue: Monte Carlo approach; Red: Antithetic Variable approach; Natural logarithm plot). Solid line: actual results. Dashed line: linear regression fit. The quantities of interest are the same as on Figure 2.1.

$2N$	$W_N^*$	$\partial_{\xi_1} W_N^*$	$\partial_{\xi_2} W_N^*$	$\partial_{\xi_1 \xi_1}^2 W_N^*$	$\partial_{\xi_1 \xi_2}^2 W_N^*$	$\partial_{\xi_2 \xi_2}^2 W_N^*$	$\xi \cdot \partial_{\xi} W_N^*$	$\xi^T \partial_{\xi}^2 W_N^* \xi$
10	20.38	11.57	14.14	9.940	6.206	13.28	19.89	19.57
20	23.86	12.34	13.32	9.993	7.548	9.265	23.33	23.00
40	18.94	12.16	10.16	9.726	6.060	8.902	18.50	18.24
60	22.11	13.30	13.35	10.73	7.513	10.88	21.68	21.41
80	12.89	9.080	9.295	10.09	4.420	8.598	12.61	12.45
100	12.37	10.17	8.635	11.21	3.896	10.24	12.12	11.96
200	17.07	9.708	9.864	7.731	5.631	8.284	16.71	16.49

Table 2.2 – Test-Case 2: Variance reduction ratios (2.58).

### 2.3.5 Test Case 3

We eventually turn to our final test-case, where both coefficients  $a$  and  $c$  do depend on the space variable.

We show on Figure 2.4 the variances (2.57). Again, we observe that they all decrease at the rate  $1/|Q_N|$  as  $N$  increases, and that the variance obtained with our approach is systematically smaller than the Monte Carlo variance.

On Table 2.3, we report the variance reduction ratios (2.58) (with the same convention as in Table 2.1). Results are quantitatively similar to the ones obtained on Table 2.2: we do observe a robust variance reduction, even in cases for which theoretical support is still currently missing.

$2N$	$W_N^*$	$\partial_{\xi_1} W_N^*$	$\partial_{\xi_2} W_N^*$	$\partial_{\xi_1 \xi_1}^2 W_N^*$	$\partial_{\xi_1 \xi_2}^2 W_N^*$	$\partial_{\xi_2 \xi_2}^2 W_N^*$	$\xi \cdot \partial_{\xi} W_N^*$	$\xi^T \partial_{\xi}^2 W_N^* \xi$
10	14.26	12.69	10.00	12.38	8.333	10.65	14.76	19.37
20	10.82	8.166	7.669	8.304	7.730	8.827	11.29	18.11
40	7.014	7.077	5.613	10.28	6.776	7.310	7.731	14.32
60	10.45	10.84	8.666	11.72	8.896	9.524	11.82	19.01
80	6.961	5.880	7.250	8.800	4.646	8.996	7.522	11.10
100	8.543	6.780	7.970	8.873	4.669	10.26	8.798	11.66
200	7.589	7.362	6.816	9.457	5.373	9.328	8.391	13.14

Table 2.3 – Test Case 3: Variance reduction ratios (2.58).

## Acknowledgments

The work of FL and WM is partially supported by ONR under Grant N00014-09-1-0470. WM gratefully acknowledges the support from Labex MMCD (Multi-Scale Modelling & Experimentation of Materials for Sustainable Construction) under contract ANR-11-LABX-0022. We also wish to thank Claude Le Bris for enlightening discussions.

## 2.4 Appendix: Proof of (2.8)

In all what follows,  $\xi \in \mathbb{R}^d$  is fixed. We introduce the quantity

$$W_N^{*,D}(\omega, \xi) := \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} W(y, \omega, \xi + \nabla w(y)) dy, w \in W_0^{1,p}(Q_N) \right\}$$



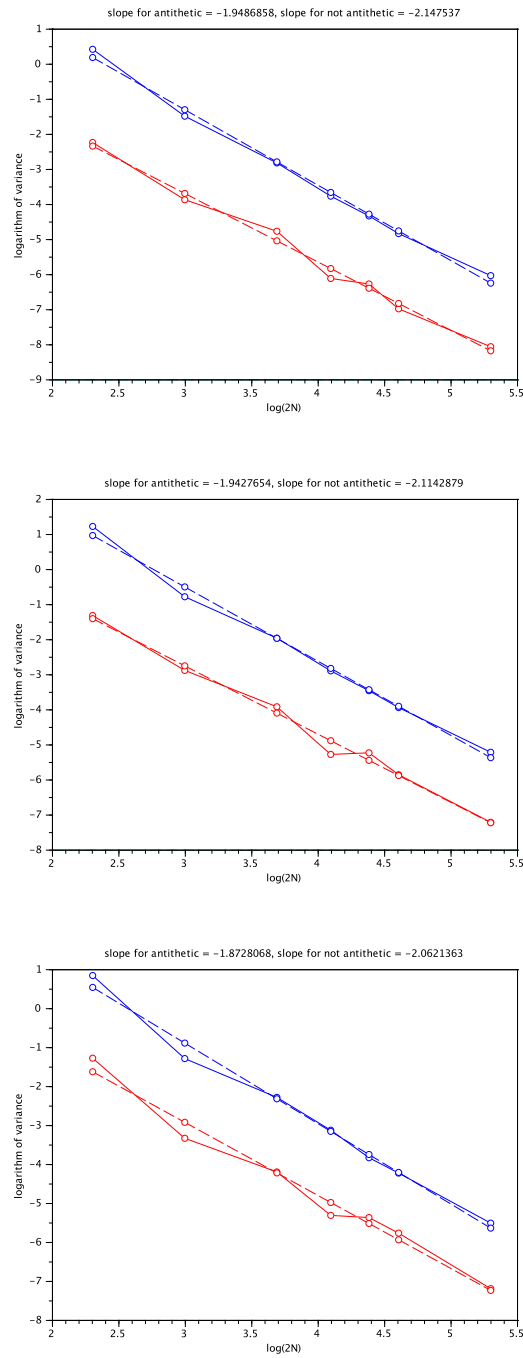


Figure 2.4 – Test-Case 3: Variances (2.57) as a function of  $2N$  (Blue: Monte Carlo approach; Red: Antithetic Variable approach; Natural logarithm plot). Solid line: actual results. Dashed line: linear regression fit. The quantities of interest are the same as on Figure 2.1.

where we emphasize in the notation that we work with homogeneous Dirichlet boundary conditions. We know from [DMM86a, DMM86b] (see also [GN11, Theorem 3.1]) that  $W_N^{*,D}(\omega, \xi)$  almost surely converges to a deterministic limit, and that this limit is the homogenized energy density:

$$W^*(\xi) = \lim_{N \rightarrow \infty} W_N^{*,D}(\omega, \xi) \quad \text{a.s.} \quad (2.59)$$

We now minimize the same functional in the space of functions satisfying periodic boundary conditions rather than homogeneous Dirichlet boundary conditions: following (2.9), we introduce

$$W_N^*(\omega, \xi) := \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} W(y, \omega, \xi + \nabla w(y)) \, dy, \quad w \in W_{\#}^{1,p}(Q_N) \right\}. \quad (2.60)$$

We see that  $W_0^{1,p}(Q_N) \subset W_{\#}^{1,p}(Q_N)$ , hence  $W_N^*(\omega, \xi) \leq W_N^{*,D}(\omega, \xi)$ . Passing to the limit  $N \rightarrow \infty$  and using (2.59), we have

$$\limsup_{N \rightarrow \infty} W_N^*(\omega, \xi) \leq W^*(\xi) \quad \text{a.s.} \quad (2.61)$$

We prove below the following result, claimed in (2.8):

**Lemma 2.16.** *Under the assumptions of Section 2.1.1, we have*

$$\lim_{N \rightarrow \infty} W_N^*(\omega, \xi) = W^*(\xi) \quad \text{a.s.} \quad (2.62)$$

The proof goes as follows. Let  $w_N$  be the minimizer of (2.60) of mean zero:

$$W_N^*(\omega, \xi) = \frac{1}{|Q_N|} \int_{Q_N} W(y, \omega, \xi + \nabla w_N(y, \omega)) \, dy, \quad \int_{Q_N} w_N(y, \omega) = 0.$$

We introduce

$$w_0^N(x, \omega) := \frac{1}{N} w^N(Nx, \omega)$$

and recognize that

$$w_0^N(\cdot, \omega) = \arg \inf \left\{ \int_Q W(Nx, \omega, \xi + \nabla \theta(x)) \, dx, \quad \theta \in W_{\#}^{1,p}(Q), \quad \int_Q \theta = 0 \right\} \quad (2.63)$$

and

$$W_N^*(\omega, \xi) = \int_Q W(Nx, \omega, \xi + \nabla w_0^N(x, \omega)) \, dx. \quad (2.64)$$

Using the test function  $\theta \equiv 0$  in (2.63) and using (2.4), we obtain

$$c_1 \int_Q |\xi + \nabla w_0^N|^p \leq W_N^*(\omega, \xi) \leq \int_Q W(Nx, \omega, \xi) \leq c_2(1 + |\xi|^p).$$

Hence there exists a constant  $C$  independent of  $N$  and  $\omega$  such that  $\|\nabla w_0^N(\cdot, \omega)\|_{L^p(Q)} \leq C$ . A standard Poincaré-Wirtinger argument yields  $\|w_0^N(\cdot, \omega)\|_{W^{1,p}(Q)} \leq C$ . Up to the extraction of a subsequence, we thus have

$$w_0^N(\cdot, \omega) \rightharpoonup w_0^\infty(\cdot, \omega) \quad \text{in } W^{1,p}(Q) \quad \text{almost surely}$$

where, almost surely,  $w_0^\infty(\cdot, \omega) \in W_{\#}^{1,p}(Q)$  and satisfies  $\int_Q w_0^\infty(\cdot, \omega) = 0$ .

We are now in position to apply [JKO94, Lemma 15.3 p. 444], from which we deduce that

$$\liminf_{N \rightarrow \infty} \int_Q W(N\cdot, \omega, \xi + \nabla w_0^N) \geq \int_Q W^*(\xi + \nabla w_0^\infty). \quad (2.65)$$

Using Jensen inequality and (2.64), we infer from (2.65) that

$$\liminf_{N \rightarrow \infty} W_N^*(\omega, \xi) \geq W^* \left( \int_Q \xi + \nabla w_0^\infty \right) = W^*(\xi). \quad (2.66)$$

Collecting (2.61) and (2.66), we get

$$W^*(\xi) \leq \liminf_{N \rightarrow \infty} W_N^*(\omega, \xi) \leq \limsup_{N \rightarrow \infty} W_N^*(\omega, \xi) \leq W^*(\xi).$$

This concludes the proof of Lemma 2.16.

## Chapter 3

# A control variate approach based on a defect-type theory for variance reduction in stochastic homogenization

Ce **Chapitre** reprend l'intégralité d'un article publié dans le journal *Multiscale Modeling & Simulation* [LM15a], et écrit en collaboration avec Frédéric Legoll.

Par la méthode des variables de contrôle, nous construisons un estimateur de la matrice homogénéisée (pour un problème d'homogénéisation aléatoire) de moindre variance. Pour construire notre estimateur, nous exploitons certains éléments d'une théorie perturbative, que nous utilisons en tant que variable de contrôle, et non dans le régime perturbatif. Nous utilisons pour calculer notre variable de contrôle une méthode de bases réduites. L'efficacité de l'approche est illustrée par des expérimentations numériques en dimension 2. De plus, l'approche est analysée dans certains cas simples.

## A control variate approach based on a defect-type theory for variance reduction in stochastic homogenization

Frédéric Legoll<sup>1,3</sup> and William Minvielle<sup>2,3</sup>

legoll@lami.enpc.fr, william.minvielle@cermics.enpc.fr

<sup>1</sup> Laboratoire Navier, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal, 77455 Marne-La-Vallée Cedex 2, France,

<sup>2</sup> CERMICS, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal, 77455 Marne-La-Vallée Cedex 2, France;

<sup>3</sup> INRIA Rocquencourt, MATHERIALS research-team, Domaine de Voluceau, B.P. 105, 78153 Le Chesnay Cedex, France.

**Abstract.** *We consider a variance reduction approach for the stochastic homogenization of divergence form linear elliptic problems. Although the exact homogenized coefficients are deterministic, their practical approximations are random. We introduce a control variate technique to reduce the variance of the computed approximations of the homogenized coefficients. Our approach is based on a surrogate model inspired by a defect-type theory, where a perfect periodic material is perturbed by rare defects. This model has been introduced in [ALB10] in the context of weakly random models. In this work, we address the fully random case, and show that the perturbative approaches proposed in [ALB10, ALB11] can be turned into an efficient control variable.*

*We theoretically demonstrate the efficiency of our approach in simple cases. We next provide illustrating numerical results and compare our approach with other variance reduction strategies. We also show how to use the Reduced Basis approach proposed in [LBT12] so that the cost of building the surrogate model remains limited.*

### 3.1 Introduction

In this work, we introduce a variance reduction approach based on the control variate technique for the homogenization of the following stochastic, elliptic, linear problem:

$$-\operatorname{div}\left(A\left(\frac{x}{\varepsilon}, \omega\right) \nabla u^\varepsilon\right) = f \text{ in } \mathcal{D}, \quad u^\varepsilon(\cdot, \omega) = 0 \text{ on } \partial\mathcal{D}, \quad (3.1)$$

set on a bounded domain  $\mathcal{D}$  in  $\mathbb{R}^d$ , where  $f$  is a deterministic function in  $L^2(\mathcal{D})$ . The random matrix  $A$  is assumed to be uniformly elliptic, bounded and stationary in a sense made precise below.

It is well-known that, in the limit when  $\varepsilon$  goes to 0, the above problem converges to the homogenized problem

$$-\operatorname{div}(A^* \nabla u^*) = f \text{ in } \mathcal{D}, \quad u^* = 0 \text{ on } \partial\mathcal{D}, \quad (3.2)$$

where the homogenized matrix  $A^*$  is deterministic, and given by an expectation of an integral involving the so-called corrector function, that solves a random auxiliary problem set

on the *entire* space. In practice, the corrector problem is approximated by a problem set on a *bounded* domain  $Q_N$  (see Section 3.1.2 below for details). A by-product of this truncation procedure is that the *deterministic* matrix  $A^*$  is in practice approximated by a *random*, apparent homogenized matrix  $A_N^*(\omega)$ . Randomness therefore comes again into the picture. In this work, we introduce a variance reduction approach to obtain practical approximations of  $A^*$  with a smaller variance. Our approach is a control variate technique, which is based on a surrogate random model, simple enough to allow for easier computations, and close enough to the reference model to eventually improve the accuracy.

We mention that, in our previous works [BCLBL12b, BCLBL12a, CLBL10], we have already proposed variance reduction approaches to compute better approximations of  $A^*$ . We used there the technique of antithetic variables, which is a generic variance reduction approach. In addition, we have shown in [LM15b] that this technique carries over to nonlinear stochastic homogenization problems, when the problem at hand is formulated as a variational convex problem. In this work, we return to the linear equation (3.1), and design an approach based on the control variate technique, where a surrogate model is used to improve the computational efficiency. Our approach here is therefore much more specific to the problem at hand than the antithetic variable approaches proposed previously. We therefore expect this technique to provide better results. This is indeed the case, as discussed along the numerical examples of Section 3.5.1.

Generally speaking, control variate approaches are based on using surrogate models as a kind of preconditioner (see Section 3.1.3 below for more details). In this work, the surrogate model that we use is inspired by a defect-type model, introduced in [ALB10, ALB12, ALB11] in the context of weakly random models. The model considered there is that of a perfect periodic material perturbed by rare defects. These defects may introduce a significant change in the local properties of the random matrix  $A(x, \omega)$ . However they only occur with a small probability  $\eta$ . In that setting, when  $\eta$  is small, the authors of [ALB10, ALB12, ALB11] have shown that a good approximation of the homogenized properties can be obtained by only solving deterministic problems rather than random problems, as usually required in stochastic homogenization. In this work, we build our surrogate model upon the ideas of [ALB10, ALB12, ALB11]. However, we address the regime when  $\eta$  is not small, hence perturbative approaches are not accurate enough.

Our article is organized as follows. In the sequel of this introduction, we present in more details some basic elements of stochastic homogenization, situate the questions under consideration in a more general setting, and introduce the control variate approach in a general setting (see Section 3.1.3). In Section 3.2, we recall the weakly stochastic model introduced in [ALB10, ALB12, ALB11].

Next, in Section 3.3, we describe how to use this weakly stochastic model to build surrogate models that can be used in the “fully random” (non perturbative) regime. We introduce two control variate approaches. The first approach (see Section 3.3.1) is based on a first-order weakly stochastic approach, where defects are considered as *isolated* from one another. The second one (see Section 3.3.2) is based on a second-order weakly stochastic approach, where *pairs of defects* are considered. The main qualitative difference between these two control variate approaches is that the second one takes into account the geometry, whereas the first one essentially only depends on  $\int_{Q_N} A(x, \omega) dx$ . It is well known that, in dimension  $d \geq 2$ , geometry – i.e. the way different materials are located one with respect to the other – matters in the homogenization process. The fact that our second approach takes into account the geometry is thus a very interesting feature.

We next collect in Section 3.4 some elements of theoretical analysis. We first consider

the one-dimensional case (Section 3.4.1) and provide there a complete analysis of our approach (see Propositions 3.11 and 3.13). We show that the variance of the apparent homogenized coefficient scales as  $N^{-1}$  (where  $N$  is the size of the large domain on which, in practice, the corrector problem is solved), while it is decreased to  $N^{-2}$  (resp.  $N^{-3}$ ) when using our first-order (resp. second-order) control variate approach. In Section 3.4.2, we next turn to the multi-dimensional case. Our main result is Lemma 3.14.

Section 3.5 is devoted to numerical experiments. We quantitatively demonstrate the efficiency of our approach on two test cases in Sections 3.5.1 and 3.5.2. As pointed out above, our second approach is based on considering pairs of defects. In order to keep limited the offline cost associated to building the surrogate model, we show in Section 3.5.3 that it is possible to use the Reduced Basis approach introduced in [LBT12]: the precomputation cost is then dramatically decreased, while the gain in variance with respect to a Monte Carlo approach remains similar.

### 3.1.1 Homogenization theoretical setting

To begin with, we introduce the basic setting of stochastic homogenization we employ. We refer to [PV81] for some seminal contribution, to [ES08] for a general, numerically oriented presentation, and to [BLP78, CD99, JKO94] for classical textbooks. We also refer to [LB10] and the review article [ACLB<sup>+</sup>12] (and the extensive bibliography contained therein) for a presentation of our particular setting. Throughout this article,  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space and we denote by  $\mathbb{E}(X) = \int_{\Omega} X(\omega) d\mathbb{P}(\omega)$  the expectation of any random variable  $X \in L^1(\Omega, d\mathbb{P})$ . We next fix  $d \in \mathbb{N}^*$  (the ambient physical dimension), and assume that the group  $(\mathbb{Z}^d, +)$  acts on  $\Omega$ . We denote by  $(\tau_k)_{k \in \mathbb{Z}^d}$  this action, and assume that it preserves the measure  $\mathbb{P}$ , that is, for all  $k \in \mathbb{Z}^d$  and all  $A \in \mathcal{F}$ ,  $\mathbb{P}(\tau_k A) = \mathbb{P}(A)$ . We assume that the action  $\tau$  is *ergodic*, that is, if  $A \in \mathcal{F}$  is such that  $\tau_k A = A$  for any  $k \in \mathbb{Z}^d$ , then  $\mathbb{P}(A) = 0$  or 1. In addition, we define the following notion of stationarity (see [BLBL06, BLBL07]): a function  $F \in L^1_{\text{loc}}(\mathbb{R}^d, L^1(\Omega))$  is *stationary* if

$$\forall k \in \mathbb{Z}^d, \quad F(x+k, \omega) = F(x, \tau_k \omega) \quad \text{a.e. in } x \text{ and a.s.} \quad (3.3)$$

In this setting, the ergodic theorem [Kre85, Shi84, Tem72] can be stated as follows: *Let  $F \in L^\infty(\mathbb{R}^d, L^1(\Omega))$  be a stationary random variable in the above sense. For  $k = (k_1, k_2, \dots, k_d) \in \mathbb{Z}^d$ , we set  $|k|_\infty = \sup_{1 \leq i \leq d} |k_i|$ . Then*

$$\frac{1}{(2N+1)^d} \sum_{|k|_\infty \leq N} F(x, \tau_k \omega) \xrightarrow{N \rightarrow \infty} \mathbb{E}(F(x, \cdot)) \quad \text{in } L^\infty(\mathbb{R}^d), \text{ almost surely.}$$

*This implies (denoting by  $Q$  the unit cube in  $\mathbb{R}^d$ ) that*

$$F\left(\frac{x}{\varepsilon}, \omega\right) \xrightarrow[\varepsilon \rightarrow 0]{*} \mathbb{E}\left(\int_Q F(x, \cdot) dx\right) \quad \text{in } L^\infty(\mathbb{R}^d), \text{ almost surely.}$$

Besides technicalities, the purpose of the above setting is simply to formalize that, even though realizations may vary, the function  $F$  at point  $x \in \mathbb{R}^d$  and the function  $F$  at point  $x+k$ ,  $k \in \mathbb{Z}^d$ , share the same law. In the homogenization context we now turn to, this means that the local, microscopic environment (encoded in the matrix field  $A$  in (3.1)) is everywhere the same *on average*. From this, homogenized, macroscopic properties will follow. In addition, and this is evident reading the above setting, the

microscopic environment has a relation to an underlying *periodic* structure (thus the integer shifts  $k$  in (3.3)).

We consider problem (3.1), where  $\mathcal{D}$  is an open, bounded domain of  $\mathbb{R}^d$  and where  $f \in L^2(\mathcal{D})$  is deterministic. The random matrix  $A$  is assumed stationary in the sense of (3.3). We also assume that  $A$  is bounded and that, in the sense of quadratic forms,  $A$  is positive and almost surely bounded away from zero: there exist deterministic constants  $c$  and  $C$  such that, almost surely,

$$\|A(\cdot, \omega)\|_{L^\infty(\mathbb{R}^d)} \leq C \quad \text{and} \quad \forall \xi \in \mathbb{R}^d, \quad \xi^T A(x, \omega) \xi \geq c \xi^T \xi \quad \text{a.e.} \quad (3.4)$$

In this specific setting, the solution  $u^\varepsilon(\cdot, \omega)$  to (3.1) converges (when  $\varepsilon$  goes to 0) to the solution  $u^*$  to the homogenized problem (3.2) almost surely, weakly in  $H^1(\mathcal{D})$  and strongly in  $L^2(\mathcal{D})$ . The homogenized matrix  $A^*$  that appears in (3.2) reads

$$\forall p \in \mathbb{R}^d, \quad A^* p = \mathbb{E} \left[ \int_Q A(x, \cdot) (\nabla w_p(x, \cdot) + p) dx \right], \quad Q = (0, 1)^d, \quad (3.5)$$

where, for any vector  $p \in \mathbb{R}^d$ , the *corrector*  $w_p$  is the solution (unique up to the addition of a random constant) to the following corrector problem:

$$\begin{cases} -\operatorname{div} [A(\nabla w_p + p)] = 0 & \text{in } \mathbb{R}^d \text{ a.s.}, \\ \nabla w_p \text{ is stationary in the sense of (3.3),} & \int_Q \mathbb{E}(\nabla w_p) = 0. \end{cases} \quad (3.6)$$

### 3.1.2 Practical approximation of the homogenized matrix

The corrector problem (3.6) is set on the entire space  $\mathbb{R}^d$ , and is therefore challenging to solve. Approximations are in order. In practice, the deterministic matrix  $A^*$  is approximated by the random matrix  $A_N^*(\omega)$  defined by

$$\forall p \in \mathbb{R}^d, \quad A_N^*(\omega) p = \frac{1}{|Q_N|} \int_{Q_N} A(x, \omega) (p + \nabla w_p^N(x, \omega)) dx, \quad (3.7)$$

which is obtained by solving the corrector problem on a *truncated* domain, say the cube  $Q_N = (-N/2, N/2)^d$ :

$$-\operatorname{div} (A(\cdot, \omega) (p + \nabla w_p^N(\cdot, \omega))) = 0, \quad w_p^N(\cdot, \omega) \text{ is } Q_N\text{-periodic.} \quad (3.8)$$

As briefly explained above, although  $A^*$  itself is a deterministic object, its practical approximation  $A_N^*$  is random. It is only in the limit of infinitely large domains  $Q_N$  that the deterministic value is attained. Indeed, as shown in [BP04], we have

$$\lim_{N \rightarrow \infty} A_N^*(\omega) = A^* \quad \text{almost surely.}$$

Many studies have been recently devoted to establishing sharp estimates on the convergence of the random apparent homogenized quantities (computed on  $Q_N$ ) to the exact deterministic homogenized quantities. We refer e.g. to [BP04, GNO14, Nol14, Yur86] and to the comprehensive discussion of [BCLBL12b, Section 1.2]. We take here the problem from a slightly different perspective. We observe that the error

$$A^* - A_N^*(\omega) = \left( A^* - \mathbb{E}[A_N^*] \right) + \left( \mathbb{E}[A_N^*] - A_N^*(\omega) \right)$$



is the sum of a systematic error and of a statistical error (the first and second terms in the above right-hand side, respectively). We focus here on the *statistical error*, and propose approaches to reduce the confidence interval of empirical means approximating  $\mathbb{E}[A_N^*]$ , for a given truncated domain  $Q_N$ . Optimal estimates on the variance of  $A_N^*$  have been established in [Nol14, Theorem 1.3 and Proposition 1.4]. For a setting slightly different from ours (namely for homogenization problems set on random *lattices*), optimal estimates on the systematic and statistical errors have been established in [GNO14, Theorem 2]. The authors noted there that “the systematic error is much smaller than the statistical error”, in the sense that the latter decays with a slower rate with respect to  $N$  than the former. For large values of  $N$ , the statistical error (that we address in this work) is therefore dominating over the systematic error.

A standard technique to compute an approximation of  $\mathbb{E}[(A_N^*)_{ij}]$  (for any entry  $ij$ ) is to consider  $M$  independent and identically distributed realizations of the field  $A$ , solve for each of them the corrector problem (3.8) (thereby obtaining i.i.d. realizations  $A_N^{*,m}(\omega)$ ) and proceed following a Monte Carlo approach:

$$\mathbb{E}[(A_N^*)_{ij}] \approx I_M^{\text{MC}} := \frac{1}{M} \sum_{m=1}^M (A_N^{*,m}(\omega))_{ij}. \quad (3.9)$$

In view of the Central Limit Theorem, we know that our quantity of interest  $\mathbb{E}[(A_N^*)_{ij}]$  asymptotically lies in the confidence interval

$$\left[ I_M^{\text{MC}} - 1.96 \frac{\sqrt{\text{Var}[(A_N^*)_{ij}]}}{\sqrt{M}}, I_M^{\text{MC}} + 1.96 \frac{\sqrt{\text{Var}[(A_N^*)_{ij}]}}{\sqrt{M}} \right]$$

with a probability equal to 95 %.

In this article, we show that, using a control variate approach, we can design a practical approach that, for any finite  $N$ , allows to compute a better approximation of  $\mathbb{E}[(A_N^*)_{ij}]$  than  $I_M^{\text{MC}}$ . Otherwise stated, for an equal computational cost, we obtain a more accurate (i.e. with a smaller confidence interval) approximation.

### 3.1.3 Control variate approach

Before presenting our specific approach, we describe here the control variate approach in a general context (see [Fis96, page 277]). Consider a general probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and a scalar random variable  $X \in L^2(\Omega, \mathbb{R})$ . Our aim is to compute its expectation  $\mathbb{E}(X)$ . In the sequel, we will use that approach for the random variable  $(A_N^*(\omega))_{ij}$ , for any entry  $1 \leq i, j \leq d$ .

As always, a first possibility is to resort to  $M$  i.i.d. realizations of  $X$ , denoted  $X^m(\omega)$  for  $1 \leq m \leq M$ . The expectation is then approximated by the Monte Carlo empirical mean

$$I_M^{\text{MC}} := \frac{1}{M} \sum_{m=1}^M X^m(\omega)$$

and we know that, with a probability equal to 95 %,  $\mathbb{E}[X]$  asymptotically lies in the confidence interval

$$\left[ I_M^{\text{MC}} - 1.96 \frac{\sqrt{\text{Var}[X]}}{\sqrt{M}}, I_M^{\text{MC}} + 1.96 \frac{\sqrt{\text{Var}[X]}}{\sqrt{M}} \right]. \quad (3.10)$$

To reduce the variance of the estimation, consider now a random variable  $Y \in L^2(\Omega, \mathbb{R})$ , the expectation of which is analytically known. Then, for any scalar deterministic parameter  $\rho$  to be fixed later, we consider the *controlled variable*

$$D_\rho(\omega) = X(\omega) - \rho(Y(\omega) - \mathbb{E}[Y]). \quad (3.11)$$

Since  $\mathbb{E}[Y]$  is known exactly, sampling realizations of  $D_\rho$  amounts to sampling realizations of  $X$  and  $Y$ . We obviously have  $\mathbb{E}[D_\rho] = \mathbb{E}[X]$ . To approximate  $\mathbb{E}[X]$ , the control variate approach consists in performing a standard Monte Carlo approximation on  $D_\rho$ . We hence consider  $M$  i.i.d. realizations of  $D_\rho$ , denoted  $D_\rho^m(\omega)$ , introduce the empirical mean

$$I_M^{\text{CV}} := \frac{1}{M} \sum_{m=1}^M D_\rho^m(\omega)$$

and write that, with a probability equal to 95 %,  $\mathbb{E}[D_\rho] = \mathbb{E}[X]$  asymptotically lies in the confidence interval

$$\left[ I_M^{\text{CV}} - 1.96 \frac{\sqrt{\text{Var}[D_\rho]}}{\sqrt{M}}, I_M^{\text{CV}} + 1.96 \frac{\sqrt{\text{Var}[D_\rho]}}{\sqrt{M}} \right]. \quad (3.12)$$

If  $\rho$  and  $Y$  are such that  $\text{Var}[D_\rho] < \text{Var}[X]$ , then the width of the above confidence interval is smaller than that of (3.10), and hence we have built a more accurate approximation of  $\mathbb{E}[X]$ .

We now detail how to choose  $\rho$  and  $Y$  in (3.11). Suppose for now that  $Y$  is given. We wish to pick  $\rho$  such that the variance of  $D_\rho$  is minimal. Writing that

$$\text{Var}[D_\rho] = \text{Var}[X] - 2\rho \text{Cov}[X, Y] + \rho^2 \text{Var}[Y],$$

we see that the optimal value of  $\rho$  reads

$$\rho^* = \text{argmin} \text{Var}[D_\rho] = \frac{\text{Cov}[X, Y]}{\text{Var}[Y]}. \quad (3.13)$$

For this choice, we have, using the Cauchy-Schwarz inequality,

$$\text{Var}[D_{\rho^*}] = \text{Var}[X] \left( 1 - \frac{(\text{Cov}[X, Y])^2}{\text{Var}[X]\text{Var}[Y]} \right) \leq \text{Var}[X].$$

We thus observe that, for any choice of  $Y$ , we can choose  $\rho$  such that the variance of  $D_\rho$  is indeed smaller than that of  $X$ . Of course, the ratio of variances  $\frac{\text{Var}[D_{\rho^*}]}{\text{Var}[X]}$ , which is directly related to the gain in accuracy, depends on  $Y$ , and more precisely on the value of  $\frac{(\text{Cov}[X, Y])^2}{\text{Var}[X]\text{Var}[Y]}$ . The larger the correlation between  $X$  and  $Y$ , the better. In contrast to the choice of  $\rho$ , the choice of  $Y$  is problem dependent. In addition, the control variable  $Y$  needs to be *random*.

**Remark 3.1.** *In practice, we do not have access to the optimal value (3.13), which involves exact expectations. One possibility (which is the one we adopt in this work) is to replace (3.13) by the empirical estimator*

$$\rho^* \approx \frac{\sum_{m=1}^M (X^m(\omega) - \mu_M(X)) (Y^m(\omega) - \mathbb{E}[Y])}{\sum_{m=1}^M (Y^m(\omega) - \mathbb{E}[Y])^2},$$

where  $\mu_M(X) = \frac{1}{M} \sum_{m=1}^M X^m(\omega)$ . This choice corresponds to minimizing with respect to  $\rho$  the empirical variance of  $D_\rho$  defined as  $\frac{1}{M} \sum_{m=1}^M (D_\rho^m(\omega) - \mu_M(X))^2$ , where  $D_\rho^m(\omega) = X^m(\omega) - \rho(Y^m(\omega) - \mathbb{E}[Y])$ .

## 3.2 A weakly random setting: rare defects in a periodic structure

As pointed out above, the surrogate model that we use to build our controlled variable is inspired by a defect-type model, introduced in [ALB10, ALB12, ALB11] in the context of weakly random models, and that we describe now.

### 3.2.1 Presentation of the model

Assume that, in (3.1), the random matrix  $A$  is of the form

$$A(x, \omega) = A_\eta(x, \omega) = A_{\text{per}}(x) + b_\eta(x, \omega) \left( C_{\text{per}}(x) - A_{\text{per}}(x) \right) \quad (3.14)$$

where  $A_{\text{per}}$  and  $C_{\text{per}}$  are  $\mathbb{Z}^d$ -periodic matrices that are bounded and positive in the sense of (3.4), and

$$b_\eta(x, \omega) = \sum_{k \in \mathbb{Z}^d} \mathbf{1}_{Q+k}(x) B_k^\eta(\omega), \quad (3.15)$$

where  $(B_k^\eta)_{k \in \mathbb{Z}^d}$  are i.i.d. scalar random variables. The matrix  $A$  is indeed stationary in the sense of (3.3). We furthermore assume that  $B_k^\eta$  follows a Bernoulli law of parameter  $\eta \in (0, 1)$ :

$$\mathbb{P}(B_k^\eta = 1) = \eta, \quad \mathbb{P}(B_k^\eta = 0) = 1 - \eta. \quad (3.16)$$

The matrix  $A(x, \omega)$  then satisfies assumption (3.4).

In each cell  $Q + k$ , the field  $A$  is equal to  $A_{\text{per}}$  with the probability  $1 - \eta$ , and equal to  $C_{\text{per}}$  with the probability  $\eta$ . When  $\eta$  is small, then (3.14)–(3.15)–(3.16) models a periodic material (described by  $A_{\text{per}}$ ) that is randomly perturbed (and then described by  $C_{\text{per}}$ ). The perturbation is rare when  $\eta$  is small (therefore the material is described by  $A_{\text{per}}$  “most of the time”), and thus it can be considered as a defect. However, the perturbation is not small in  $L^\infty$  norm:  $\|C_{\text{per}} - A_{\text{per}}\|_{L^\infty}$  is not assumed to be small. We refer to [ALB11] for practical examples motivating this framework.

On Fig. 3.1, we show two realizations of the field  $A_\eta(x, \omega)$  (on the domain  $Q_N$  for  $N = 20$ ) for some specific choices of  $A_{\text{per}}$  and  $C_{\text{per}}$  (see [ALB11, Fig. 4.2] for more details). On the right part of that figure, we set  $\eta = 0.4$ , which is close to the value  $\eta = 1/2$ , when defects are as frequent as non-defects.

Note that specifying  $A_\eta(x, \omega)$  on  $Q_N$  simply amounts to specifying the values of  $B_k^\eta(\omega)$  for all  $k$  such that  $k + Q \subset Q_N$ .

The above setting is actually quite general. Consider for instance a classical test-case, the random checkerboard case:

$$A(x, \omega) = \sum_{k \in \mathbb{Z}^d} \mathbf{1}_{Q+k}(x) X_k(\omega),$$

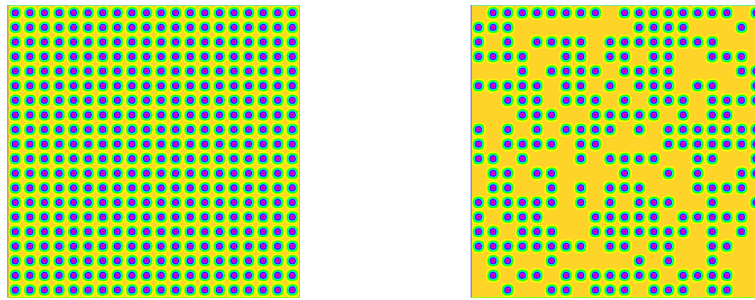


Figure 3.1 – Two instances of material (3.14). Left ( $\eta = 0$ ): perfect material with circular inclusions located on a periodic network. Right ( $\eta = 0.4$ ): perturbed material (each inclusion is deleted with a probability equal to 0.4). Courtesy A. Anantharaman and C. Le Bris.

where  $X_k$  are i.i.d. random variables satisfying  $\mathbb{P}(X_k = \alpha) = \mathbb{P}(X_k = \beta) = 1/2$ . This model falls into the framework (3.14)–(3.15)–(3.16) with

$$A_{\text{per}} = \alpha \text{Id}, \quad C_{\text{per}} = \beta \text{Id}, \quad \eta = 1/2.$$

An alternate choice (corresponding to choosing a different reference periodic materials) is

$$A_{\text{per}} = \beta \text{Id}, \quad C_{\text{per}} = \alpha \text{Id}, \quad \eta = 1/2.$$

In this work, we restrict our attention to the case (3.14)–(3.15)–(3.16), i.e. when  $B_k^\eta$  are i.i.d. Bernoulli random variables. This is the case specifically studied in [ALB11]. See [ALB10, ALB12] for more general settings.

### 3.2.2 Weakly-random homogenization result

Consider the model (3.14)–(3.15)–(3.16). The random variable  $B_k^\eta(\omega)$  can take only two values, 0 or 1. Therefore, on the domain  $Q_N$ , there are only a finite number of realizations of  $A_\eta(x, \omega)$ . The realizations with the highest probability are as follows.

With probability  $(1 - \eta)^{|Q_N|}$ , there are no defects in  $Q_N$ , and the realization actually corresponds to the perfect periodic situation. We introduce the periodic corrector  $w_p^0$ , solution to

$$-\text{div} (A_{\text{per}} (p + \nabla w_p^0)) = 0, \quad w_p^0 \text{ is } Q\text{-periodic}, \quad (3.17)$$

and the associated matrix  $A_{\text{per}}^*$ , obtained by periodic homogenization:

$$\forall p \in \mathbb{R}^d, \quad A_{\text{per}}^* p = \int_Q A_{\text{per}} (p + \nabla w_p^0). \quad (3.18)$$

With probability  $\eta(1 - \eta)^{|Q_N|-1}$ , there is a unique defect in  $Q_N$ , located, say, in the cell  $k + Q$  (see Fig. 3.2). Let us define

$$A_1^k = A_{\text{per}} + 1_{k+Q} (C_{\text{per}} - A_{\text{per}}), \quad (3.19)$$

the associated corrector  $w_p^{1,k,N}$ , solution to

$$-\text{div} (A_1^k (p + \nabla w_p^{1,k,N})) = 0, \quad w_p^{1,k,N} \text{ is } Q_N\text{-periodic}, \quad (3.20)$$

and the homogenized matrix  $A_{1,k,N}^*$ , given by

$$\forall p \in \mathbb{R}^d, \quad A_{1,k,N}^* p = \frac{1}{|Q_N|} \int_{Q_N} A_1^k \left( p + \nabla w_p^{1,k,N} \right). \quad (3.21)$$

With probability  $\eta^2(1-\eta)^{|Q_N|-2}$ , there are two defects in  $Q_N$ , located, say, in the cells  $k+Q$  and  $l+Q$  (see Fig. 3.2). Let us define

$$A_2^{k,l} = A_{\text{per}} + \left( 1_{k+Q} + 1_{l+Q} \right) \left( C_{\text{per}} - A_{\text{per}} \right), \quad (3.22)$$

the associated corrector  $w_p^{2,k,l,N}$ , solution to

$$-\operatorname{div} \left( A_2^{k,l} \left( p + \nabla w_p^{2,k,l,N} \right) \right) = 0, \quad w_p^{2,k,l,N} \text{ is } Q_N\text{-periodic}, \quad (3.23)$$

and the homogenized matrix  $A_{2,k,l,N}^*$ , given by

$$\forall p \in \mathbb{R}^d, \quad A_{2,k,l,N}^* p = \frac{1}{|Q_N|} \int_{Q_N} A_2^{k,l} \left( p + \nabla w_p^{2,k,l,N} \right). \quad (3.24)$$

All the other configurations (with three defects or more) have a smaller probability.

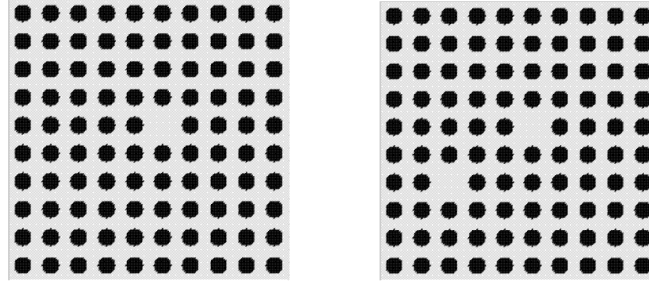


Figure 3.2 – Left: material modelled by  $A_1^k$ , with a single defect. Right: material modelled by  $A_2^{k,l}$ , with two defects (Courtesy A. Anantharaman and C. Le Bris).

Let us define

$$\mathcal{I}_N := \left\{ k \in \mathbb{Z}^d; Q + k \subset Q_N \right\}.$$

As shown in [ALB11], we then have the following result:

**Proposition 3.2** ([ALB11], Section 3.2). *Let  $A_{\eta,N}^*(\omega)$  be the apparent homogenized matrix defined by (3.7), where  $A \equiv A_\eta$  is given by (3.14)–(3.15)–(3.16). Then*

$$\mathbb{E} [A_{\eta,N}^*] = A_{\text{per}}^* + \eta \bar{A}_1^N + \eta^2 \bar{A}_2^N + O_N(\eta^3), \quad (3.25)$$

where  $O_N(\eta^3)$  is a quantity of the order of  $\eta^3$  with a prefactor that may depend on  $N$ ,  $A_{\text{per}}^*$  is given by (3.18) and

$$\begin{aligned} \bar{A}_1^N &= \sum_{k \in \mathcal{I}_N} \left( A_{1,k,N}^* - A_{\text{per}}^* \right), \\ \bar{A}_2^N &= \frac{1}{2} \sum_{k,l \in \mathcal{I}_N, k \neq l} \left( A_{2,k,l,N}^* - A_{1,k,N}^* - A_{1,l,N}^* + A_{\text{per}}^* \right). \end{aligned}$$

We note that

$$\overline{A}_1^N = \sum_{k \in \mathcal{I}_N} \overline{A}_{1 \text{ def}}^{k,N}, \quad \text{and} \quad \overline{A}_2^N = \frac{1}{2} \sum_{k \in \mathcal{I}_N} \sum_{l \in \mathcal{I}_N, l \neq k} \overline{A}_{2 \text{ def}}^{k,l,N}, \quad (3.26)$$

where  $\overline{A}_{1 \text{ def}}^{k,N}$  (resp.  $\overline{A}_{2 \text{ def}}^{k,l,N}$ ) is the marginal contribution to the homogenized matrix from a configuration with a single defect in  $k + Q$  (resp. two defects in  $k + Q$  and  $l + Q$ ):

$$\overline{A}_{1 \text{ def}}^{k,N} = A_{1,k,N}^* - A_{\text{per}}^*, \quad (3.27)$$

$$\overline{A}_{2 \text{ def}}^{k,l,N} = A_{2,k,l,N}^* - A_{1,k,N}^* - A_{1,l,N}^* + A_{\text{per}}^*. \quad (3.28)$$

**Remark 3.3.** *Passing to the limit  $N \rightarrow \infty$  in (3.25) is not easy. We refer to [ALB11, Section 3.2] and [Mou15].*

When  $\eta$  is small, the advantage of (3.25) over the approach recalled in Section 3.1.2 is evident. Rather than solving the random problem (3.8) (for several realizations of  $A_\eta$ ), it is enough to solve the deterministic problems (3.17), (3.20) and (3.23) to infer an accurate approximation of  $\mathbb{E} \left[ A_{\eta,N}^* \right]$ . We refer to [ALB11] for illustrative numerical results.

Furthermore, due to periodic boundary conditions (3.20), that are reminiscent of the periodic boundary conditions in (3.8), we have that

$$A_{1,k,N}^* \text{ does not depend on } k. \quad (3.29)$$

Likewise,  $A_{2,k,l,N}^*$  depends only on  $k - l$ . Thus, there is only *one* problem (3.20) to be solved (say for  $k = 0$ ). Likewise, there are  $|\mathcal{I}_N| - 1$  problems (3.23) to be solved (say for  $k = 0$  and  $l \neq 0$ ), and not  $|\mathcal{I}_N|(|\mathcal{I}_N| - 1)$ . Noticing that (3.23) is a problem parameterized by  $l$ , the authors of [LBT12] have shown how to use a Reduced Basis approach to further speed-up the computation of  $\overline{A}_2^N$ . In practice, one can still obtain a good approximation of  $\overline{A}_2^N$  without solving all the  $|\mathcal{I}_N| - 1$  problems (3.23). We return to this specific question in Section 3.5.3.

### 3.3 Control variate approaches for stochastic homogenization

We now introduce, for the model (3.14)–(3.15)–(3.16), a control variate approach. Our aim is now to address the regime when  $\eta$  is not close to 0 or 1 (the approximation (3.25) is therefore not accurate enough). Recall also that, in view of the discussion at the end of Section 3.1.3, we need a *random* surrogate model to build our controlled variable. In what follows, we first build an approximate model based on configurations with a single defect (see Section 3.3.1), and next turn to building a better approximate model that also uses configurations with two defects (see Section 3.3.2). As will be seen below, this second approximate model not only depends on the quantity of defects, but also on their geometry, that is on where the defects are located in  $Q_N$ .

#### 3.3.1 A first-order model

Introduce

$$A_1^{\eta,N}(\omega) = \sum_{k \in \mathcal{I}_N} B_k^\eta(\omega) \overline{A}_{1 \text{ def}}^{k,N}, \quad (3.30)$$

where  $\bar{A}_{1\text{def}}^{k,N}$ , defined by (3.27), is the marginal contribution to the homogenized matrix coming from the configuration with a single defect located in  $k+Q$ . In view of (3.26), we notice that

$$\mathbb{E} \left[ A_1^{\eta,N} \right] = \sum_{k \in \mathcal{I}_N} \mathbb{E} \left[ B_k^\eta \right] \bar{A}_{1\text{def}}^{k,N} = \eta \sum_{k \in \mathcal{I}_N} \bar{A}_{1\text{def}}^{k,N} = \eta \bar{A}_1^N,$$

which is the first order correction in the expansion (3.25). When  $\eta$  is small, the *expectation* of  $A_{\text{per}}^* + A_1^{\eta,N}(\omega)$  is a good approximation of the expectation of  $A_{\eta,N}^*(\omega)$ , accurate up to an error of the order of  $\eta^2$ . The following observation provides additional motivation for our choice (3.30). It turns out that the *law* of the random variable  $A_{\text{per}}^* + A_1^{\eta,N}(\omega)$  is a good approximation of that of  $A_{\eta,N}^*(\omega)$ :

**Lemma 3.4.** *For any deterministic and continuous function  $\varphi$ , we have*

$$\mathbb{E} \left[ \varphi \left( A_{\eta,N}^* \right) \right] = \mathbb{E} \left[ \varphi \left( A_{\text{per}}^* + A_1^{\eta,N} \right) \right] + O_N(\eta^2).$$

The proof of Lemma 3.4 is postponed until Section 3.4.2.

We thus think that  $A_{\text{per}}^* + A_1^{\eta,N}(\omega)$  is a good surrogate model for  $A_{\eta,N}^*(\omega)$ . As shown by Lemma 3.4, this is the case when  $\eta \ll 1$ , which is however not the regime we address. One-dimensional computations presented in Section 3.4.1 and numerical observations reported in Section 3.5 (for two-dimensional test-cases) confirm that it is indeed the case, even when  $\eta$  is not small.

Following (3.11), we now introduce our controlled variable as

$$\begin{aligned} D_\rho^{1,\eta}(\omega) &= A_{\eta,N}^*(\omega) - \rho \left( A_{\text{per}}^* + A_1^{\eta,N}(\omega) - \mathbb{E} \left[ A_{\text{per}}^* + A_1^{\eta,N} \right] \right) \\ &= A_{\eta,N}^*(\omega) - \rho \left( A_1^{\eta,N}(\omega) - \eta \bar{A}_1^N \right). \end{aligned} \quad (3.31)$$

In view of (3.30), (3.27) and (3.29), we recast (3.31) as

$$D_\rho^{1,\eta}(\omega) = A_{\eta,N}^*(\omega) - \rho \left[ \left( \sum_{k \in \mathcal{I}_N} B_k^\eta(\omega) \right) - \eta |\mathcal{I}_N| \right] \bar{A}_{1\text{def}}^{0,N}. \quad (3.32)$$

**Remark 3.5.** *Note that, in (3.32),  $A_{\eta,N}^*(\omega)$  and  $\sum_{k \in \mathcal{I}_N} B_k^\eta(\omega)$  are correlated. Indeed, in practice, we start by drawing a realization of the random variables  $B_k^\eta(\omega)$  for all  $k \in \mathcal{I}_N$ . This determines first  $\sum_{k \in \mathcal{I}_N} B_k^\eta(\omega)$ , and second the field  $A(x, \omega)$  on  $Q_N$ , from which we compute the associated  $A_{\eta,N}^*(\omega)$  following (3.7)–(3.8).*

Computing  $M$  realizations of  $D_\rho^{1,\eta}(\omega)$  therefore amounts to:

- offline stage: determine  $\bar{A}_{1\text{def}}^{0,N}$  by solving the problem (3.17)–(3.18) on  $Q$  and solving only once the problem (3.20)–(3.21) on  $Q_N$  (say for  $k=0$ ).
- online stage: solve  $M$  corrector problems (3.7)–(3.8) on  $Q_N$  (for  $M$  i.i.d. realizations of  $A$  on  $Q_N$ ), and evaluate  $D_\rho^{1,\eta}(\omega)$  according to (3.32).

Let  $\mathcal{C}_N$  be the cost to solve a single corrector problem on  $Q_N$ . The Monte Carlo empirical estimator and the Control Variate empirical estimator, defined respectively by

$$I_M^{\text{MC}} = \frac{1}{M} \sum_{m=1}^M A_{\eta,N}^{*,m}(\omega) \quad \text{and} \quad I_M^{\text{CV}} = \frac{1}{M} \sum_{m=1}^M D_\rho^{1,\eta,m}(\omega)$$



therefore share the same cost ( $M C_N$  for the former,  $(1+M) C_N$  for the latter). To minimize the variance of  $D_\rho^{1,\eta}$ , the parameter  $\rho$  in (3.31) is chosen following (3.13).

Notice that, in the above construction, we have considered as reference configuration the defect-free material, i.e. that for  $\eta = 0$ . Since, in the regime we focus on,  $\eta$  is not small, there is no reason to favor the defect-free configuration ( $\eta = 0$ ) rather than the full defect configuration ( $\eta = 1$ ), which corresponds to the periodic matrix  $C_{\text{per}}$ . We therefore introduce (compare with (3.27))

$$\overline{C}_{1\text{def}}^{k,N} = C_{1,k,N}^* - C_{\text{per}}^*,$$

where  $C_{1,k,N}^*$  is the homogenized matrix corresponding to a unique defect with respect to the periodic configuration  $C_{\text{per}}$  (compare with (3.19), (3.20) and (3.21)):

$$\forall p \in \mathbb{R}^d, \quad C_{1,k,N}^* p = \frac{1}{|Q_N|} \int_{Q_N} C_1^k \left( p + \nabla v_p^{1,k,N} \right), \quad (3.33)$$

where, for any  $p$ , the corrector  $v_p^{1,k,N}$  is a solution to

$$-\operatorname{div} \left( C_1^k \left( p + \nabla v_p^{1,k,N} \right) \right) = 0, \quad v_p^{1,k,N} \text{ is } Q_N\text{-periodic,}$$

where  $C_1^k = C_{\text{per}} - 1_{k+Q} \left( C_{\text{per}} - A_{\text{per}} \right)$ . In the spirit of (3.32), we introduce the controlled variable

$$\widehat{D}_\rho^{1,\eta}(\omega) = A_{\eta,N}^*(\omega) - \widehat{\rho} \left[ \left( \sum_{k \in \mathcal{I}_N} (1 - B_k^\eta(\omega)) \right) - (1 - \eta) |\mathcal{I}_N| \right] \overline{C}_{1\text{def}}^{0,N},$$

that we recast as

$$\widehat{D}_\rho^{1,\eta}(\omega) = A_{\eta,N}^*(\omega) + \widehat{\rho} \left[ \left( \sum_{k \in \mathcal{I}_N} B_k^\eta(\omega) \right) - \eta |\mathcal{I}_N| \right] \overline{C}_{1\text{def}}^{0,N}.$$

Consider now any entry  $1 \leq i, j \leq d$  of the homogenized matrix. Assuming that our control variate model is non trivial (i.e. that  $\left[ \overline{A}_{1\text{def}}^{0,N} \right]_{ij} \neq 0$ ), we see that, for any deterministic  $\widehat{\rho}$ , there exists a deterministic parameter  $\rho$  such that  $\left[ \widehat{D}_\rho^{1,\eta}(\omega) \right]_{ij} = \left[ D_\rho^{1,\eta}(\omega) \right]_{ij}$  a.s. Working with the controlled variable  $D_\rho^{1,\eta}(\omega)$  is hence equivalent to working with the controlled variable  $\widehat{D}_\rho^{1,\eta}(\omega)$ . In the sequel, we only consider the former.

**Remark 3.6.** *The situation is different in the second order model, where taking  $A_{\text{per}}$  or  $C_{\text{per}}$  as reference is not equivalent. See Section 3.3.2 below.*

**Remark 3.7.** *In view of (3.32), we see that our first order control variable only depends on  $\sum_{k \in \mathcal{I}_N} B_k^\eta(\omega)$ , which is the number of defects in the material. This approach can thus be extended to any two-phase materials, say of the type  $A(x, \omega) = A_1 + \chi(x, \omega) A_2$ , where  $\chi$  is stationary and equal to 0 or 1. In this case, the control variable reads  $\int_{Q_N} \chi(x, \omega) dx$ . We refer to [BL] for works in that direction.*



### 3.3.2 A second-order model

We now introduce a model that not only takes into account the contributions from single defects (through  $\overline{A}_{1\text{def}}^{k,N}$ , see (3.30)) but also contributions from pairs of defects. To that aim, we introduce

$$A_2^{\eta,N}(\omega) = \frac{1}{2} \sum_{k \in \mathcal{I}_N} \sum_{l \in \mathcal{I}_N, l \neq k} B_k^\eta(\omega) B_l^\eta(\omega) \overline{A}_{2\text{def}}^{k,l,N}, \quad (3.34)$$

where  $\overline{A}_{2\text{def}}^{k,l,N}$ , defined by (3.28), is the marginal contribution to the homogenized matrix associated to the configuration with two defects located in  $k + Q$  and  $l + Q$ . In view of (3.26), we notice that

$$\mathbb{E} \left[ A_2^{\eta,N} \right] = \frac{1}{2} \sum_{k \in \mathcal{I}_N} \sum_{l \in \mathcal{I}_N, l \neq k} \mathbb{E} \left[ B_k^\eta B_l^\eta \right] \overline{A}_{2\text{def}}^{k,l,N} = \frac{\eta^2}{2} \sum_{k \in \mathcal{I}_N} \sum_{l \in \mathcal{I}_N, l \neq k} \overline{A}_{2\text{def}}^{k,l,N} = \eta^2 \overline{A}_2^N,$$

which is the second order correction in the expansion (3.25). When  $\eta$  is small, the *expectation* of  $A_{\text{per}}^* + A_1^{\eta,N}(\omega) + A_2^{\eta,N}(\omega)$  is a good approximation of the expectation of  $A_{\eta,N}^*(\omega)$ , accurate up to an error of the order of  $\eta^3$ . Furthermore, we have the following result (compare with Lemma 3.4), the proof of which follows the same lines as that of Lemma 3.4 and is therefore omitted:

**Lemma 3.8.** *For any deterministic and continuous function  $\varphi$ , we have*

$$\mathbb{E} \left[ \varphi \left( A_{\eta,N}^* \right) \right] = \mathbb{E} \left[ \varphi \left( A_{\text{per}}^* + A_1^{\eta,N} + A_2^{\eta,N} \right) \right] + O_N(\eta^3).$$

In a way similar to (3.31), we now introduce our second-order controlled variable as

$$D_{\rho_1, \rho_2}^{2,\eta}(\omega) = A_{\eta,N}^*(\omega) - \rho_1 \left( A_1^{\eta,N}(\omega) - \eta \overline{A}_1^N \right) - \rho_2 \left( A_2^{\eta,N}(\omega) - \eta^2 \overline{A}_2^N \right). \quad (3.35)$$

We have introduced two deterministic parameters  $\rho_1$  and  $\rho_2$ , which need not be equal. For any choice of these parameters, we have  $\mathbb{E} \left[ D_{\rho_1, \rho_2}^{2,\eta} \right] = \mathbb{E} \left[ A_{\eta,N}^* \right]$ .

To evaluate (3.35), we first have to precompute the deterministic matrices

$$\overline{A}_{1\text{def}}^{k,N} = \overline{A}_{1\text{def}}^{0,N} \quad \text{and} \quad \overline{A}_{2\text{def}}^{k,l,N} = \overline{A}_{2\text{def}}^{0,l-k,N}.$$

Computing  $M$  realizations of  $D_{\rho_1, \rho_2}^{2,\eta}(\omega)$  therefore amounts to:

- offline stage: (i) determine  $\overline{A}_{1\text{def}}^{0,N}$  by solving the problem (3.17)–(3.18) on  $Q$  and by solving only once the problem (3.20)–(3.21) on  $Q_N$  (say for  $k = 0$ ); (ii) determine  $\overline{A}_{2\text{def}}^{0,l,N}$  by solving  $|\mathcal{I}_N| - 1$  problems (3.23)–(3.24) on  $Q_N$  (for  $k = 0$  and  $l \in \mathcal{I}_N, l \neq 0$ ).
- online stage: solve  $M$  corrector problems (3.7)–(3.8) on  $Q_N$  (for  $M$  i.i.d. realizations of  $A$  on  $Q_N$ ), and evaluate  $D_{\rho_1, \rho_2}^{2,\eta}(\omega)$  according to (3.35).

Questions related to the cost for evaluating  $\overline{A}_{2\text{def}}^{0,l,N}$  are discussed at the end of this section.

As pointed out in Section 3.3.1, in our regime of interest, there is no reason to favor the defect-free configuration rather than the full defect configuration, which corresponds to the periodic matrix  $C_{\text{per}}$ . We have shown there that there is no use to introduce the terms

representing the first order correction with respect to  $C_{\text{per}}$ . We therefore solely introduce the second order correction (compare with (3.28)):

$$\bar{C}_{2 \text{ def}}^{k,l,N} = C_{2,k,l,N}^* - C_{1,k,N}^* - C_{1,l,N}^* + C_{\text{per}}^*, \quad (3.36)$$

where  $C_{1,k,N}^*$  is defined by (3.33) and  $C_{2,k,l,N}^*$  is defined by (compare with (3.22), (3.23) and (3.24)):

$$\forall p \in \mathbb{R}^d, \quad C_{2,k,l,N}^* p = \frac{1}{|Q_N|} \int_{Q_N} C_2^{k,l} \left( p + \nabla v_p^{2,k,l,N} \right), \quad (3.37)$$

where, for any  $p \in \mathbb{R}^d$ , the corrector  $v_p^{2,k,l,N}$  is a solution to

$$-\operatorname{div} \left( C_2^{k,l} \left( p + \nabla v_p^{2,k,l,N} \right) \right) = 0, \quad v_p^{2,k,l,N} \text{ is } Q_N\text{-periodic,}$$

where  $C_2^{k,l} = C_{\text{per}} - (1_{k+Q} + 1_{l+Q}) (C_{\text{per}} - A_{\text{per}})$ . As in (3.34), we introduce

$$C_2^{\eta,N}(\omega) = \frac{1}{2} \sum_{k \in \mathcal{I}_N} \sum_{l \in \mathcal{I}_N, l \neq k} \left( 1 - B_k^\eta(\omega) \right) \left( 1 - B_l^\eta(\omega) \right) \bar{C}_{2 \text{ def}}^{k,l,N}, \quad (3.38)$$

where  $\bar{C}_{2 \text{ def}}^{k,l,N}$  is defined by (3.36), and its expectation reads

$$\begin{aligned} \bar{C}_2^{\eta,N} &:= \mathbb{E} \left[ C_2^{\eta,N} \right] = \frac{1}{2} \sum_{k \in \mathcal{I}_N} \sum_{l \in \mathcal{I}_N, l \neq k} \mathbb{E} \left[ (1 - B_k^\eta) (1 - B_l^\eta) \right] \bar{C}_{2 \text{ def}}^{k,l,N} \\ &= \frac{1}{2} \sum_{k \in \mathcal{I}_N} \sum_{l \in \mathcal{I}_N, l \neq k} (1 - \eta)^2 \bar{C}_{2 \text{ def}}^{k,l,N}. \end{aligned}$$

We eventually introduce the controlled variable (compare with (3.35))

$$\begin{aligned} D_{\rho_1, \rho_2, \rho_3}^{3,\eta}(\omega) &= A_{\eta,N}^*(\omega) - \rho_1 \left( A_1^{\eta,N}(\omega) - \eta \bar{A}_1^N \right) \\ &\quad - \rho_2 \left( A_2^{\eta,N}(\omega) - \eta^2 \bar{A}_2^N \right) - \rho_3 \left( C_2^{\eta,N}(\omega) - \bar{C}_2^{\eta,N} \right). \end{aligned} \quad (3.39)$$

Consider now a specific entry  $1 \leq i, j \leq d$  of the homogenized matrix. The control variate approach consists in approximating  $\mathbb{E} \left[ \left( A_{\eta,N}^* \right)_{ij} \right]$  by considering a Monte Carlo estimator for  $\mathbb{E} \left[ \left( D_{\rho_1, \rho_2, \rho_3}^{3,\eta} \right)_{ij} \right]$ . The deterministic parameters  $\rho_1$ ,  $\rho_2$  and  $\rho_3$  are chosen to minimize the variance of  $\left( D_{\rho_1, \rho_2, \rho_3}^{3,\eta}(\omega) \right)_{ij}$ . They are thus the solution of the following  $3 \times 3$  linear system (we drop the subscript  $i, j$  for conciseness):

$$\begin{aligned} \operatorname{Var}[A_1^{\eta,N}] \rho_1 + \operatorname{Cov}[A_1^{\eta,N}, A_2^{\eta,N}] \rho_2 + \operatorname{Cov}[A_1^{\eta,N}, C_2^{\eta,N}] \rho_3 &= \operatorname{Cov}[A_{\eta,N}^*, A_1^{\eta,N}] \\ \operatorname{Cov}[A_2^{\eta,N}, A_1^{\eta,N}] \rho_1 + \operatorname{Var}[A_2^{\eta,N}] \rho_2 + \operatorname{Cov}[A_2^{\eta,N}, C_2^{\eta,N}] \rho_3 &= \operatorname{Cov}[A_{\eta,N}^*, A_2^{\eta,N}] \\ \operatorname{Cov}[C_2^{\eta,N}, A_1^{\eta,N}] \rho_1 + \operatorname{Cov}[C_2^{\eta,N}, A_2^{\eta,N}] \rho_2 + \operatorname{Var}[C_2^{\eta,N}] \rho_3 &= \operatorname{Cov}[A_{\eta,N}^*, C_2^{\eta,N}] \end{aligned} \quad (3.40)$$

depending on the covariances between the entries  $ij$  of  $A_{\eta,N}^*$ ,  $A_1^{\eta,N}$ ,  $A_2^{\eta,N}$  and  $C_2^{\eta,N}$ . In practice, these covariances are approximated by empirical estimators (see Remark 3.1).

In practice, computing the matrices  $\bar{A}_{2 \text{ def}}^{0,l,N}$  (and likewise  $\bar{C}_{2 \text{ def}}^{0,l,N}$ ) is rather expensive (because each problem is set on the large domain  $Q_N$ , and the number of these problems

increases when  $N$  increases). It is therefore useful to approximate them using the Reduced Basis strategy introduced in [LBT12], which dramatically decreases the computational cost. The procedure is essentially as follows. We first solve the single defect problem (3.20) for  $k = 0$ , and solve (3.23) for a limited number of locations of the defect pairs, say  $k = 0$  and  $l$  close to  $k$ . On the basis of these computations, we are then in position to obtain very efficient approximations of the matrices  $\overline{A}_{2, \text{def}}^{0, l, N}$  for all  $l \in \mathcal{I}_N$ ,  $l \neq 0$ . Evaluating (3.34) is thus inexpensive. Thus, up to a limited offline cost (i.e. the cost for solving the few problems (3.23) that we have to consider), the Monte Carlo empirical estimator and the Control Variate empirical estimator, defined respectively by

$$I_M^{\text{MC}} = \frac{1}{M} \sum_{m=1}^M A_{\eta, N}^{*, m}(\omega) \quad \text{and} \quad I_M^{\text{CV}} := \frac{1}{M} \sum_{m=1}^M D_{\rho_1, \rho_2, \rho_3}^{3, \eta, m}(\omega)$$

share the same cost. We refer to Section 3.5.3 for numerical experiments using this procedure.

**Remark 3.9.** *In sharp contrast to the first order control variable, the second order control variable not only depends on the number of defects in the materials, i.e.  $\sum_{k \in \mathcal{I}_N} B_k^\eta(\omega)$ , but also on their location. The specific geometry of the materials, which is ignored in (3.32), is taken into account in (3.39).*

### 3.4 Elements of theoretical analysis

This section is devoted to establishing estimates on the gain provided by our approach. We proceed in two directions. First, in Section 3.4.1, we consider the one-dimensional case. Our main results are Propositions 3.11 and 3.13. We consider the large  $N$  regime, and estimate the variance (in terms of  $N$ ) of  $A_{\eta, N}^*$ , the controlled variables  $D_\rho^{1, \eta}$  defined by (3.31) and  $D_{\rho_1, \rho_2, \rho_3}^{3, \eta}$  defined by (3.39). We show that they are of the order of  $N^{-1}$ ,  $N^{-2}$  and  $N^{-3}$ , respectively. Note that, in this section, we do not assume  $\eta$  to be close to 0 or 1, i.e. we are in a fully random case.

In Section 3.4.2, we turn to the multi-dimensional case. Our main result is Lemma 3.14. We consider the regime when  $\eta$  is small, and estimate the variance (in terms of  $\eta$ ) of  $A_{\eta, N}^*$  and of the controlled variables  $D_\rho^{1, \eta}$  defined by (3.31) and  $D_{\rho_1, \rho_2}^{2, \eta}$  defined by (3.35). We show that the control variate approach using the first order (resp. second order) surrogate model allows to decrease the variance from  $O(\eta)$  to  $O(\eta^2)$  (resp. from  $O(\eta)$  to  $O(\eta^3)$ ).

Still in the regime  $\eta \ll 1$ , we show in Section 3.4.2 that, for an equal computational cost, the weakly stochastic approach proposed in [ALB11] (which directly compute  $\mathbb{E}(A_{\eta, N}^*)$  as in series in powers of  $\eta$ ) is more accurate than the control variate approach proposed in this work. The regime of interest for our approach is therefore when  $\eta$  is neither close to 0 nor to 1. This is the regime we consider in the numerical experiments of Section 3.5.

#### 3.4.1 One-dimensional case

In the one-dimensional case, we know that

$$A_{\eta, N}^*(\omega) = \left( \frac{1}{N} \int_0^N \frac{1}{A_\eta(x, \omega)} \right)^{-1},$$

where, for ease of notation, we set  $Q_N = (0, N)$  rather than  $Q_N = (-N/2, N/2)$  as before. In view of (3.14)–(3.15)–(3.16), we thus have

$$\frac{1}{A_{\eta, N}^*(\omega)} = \frac{1}{N} \sum_{k=0}^{N-1} \int_k^{k+1} \frac{dx}{A_{\text{per}}(x) + B_k^\eta(\omega) (C_{\text{per}}(x) - A_{\text{per}}(x))}.$$

Introducing the functions

$$f(x) = \frac{1}{x} \quad \text{and} \quad \phi(b) = \int_0^1 \frac{dx}{A_{\text{per}}(x) + b(C_{\text{per}}(x) - A_{\text{per}}(x))},$$

we thus see that

$$A_{\eta, N}^*(\omega) = f\left(\frac{1}{N} \sum_{k=0}^{N-1} \phi(B_k^\eta(\omega))\right).$$

Since  $B_k^\eta(\omega)$  are equal to 0 or 1, we can write  $\phi(B_k^\eta(\omega)) = \phi(0) + B_k^\eta(\omega)(\phi(1) - \phi(0))$ , and thus

$$A_{\eta, N}^*(\omega) = g\left(\frac{1}{N} \sum_{k=0}^{N-1} B_k^\eta(\omega)\right) \quad (3.41)$$

where the smooth function  $g$  is defined by  $g(b) = f(\phi(0) + b(\phi(1) - \phi(0)))$ .

### First order model

In view of (3.31), (3.30) and (3.27), the first-order surrogate model is given by  $A_{\text{per}}^* + A_1^{\eta, N}(\omega)$ , with

$$A_1^{\eta, N}(\omega) = \sum_{k=0}^{N-1} B_k^\eta(\omega) \bar{A}_1^{k, N} = \bar{A}_1^{0, N} \sum_{k=0}^{N-1} B_k^\eta(\omega). \quad (3.42)$$

We first state the following general result, the proof of which is postponed until Section 3.4.1.

**Lemma 3.10.** *Let*

$$X(\omega) = g\left(\frac{1}{N} \sum_{k=0}^{N-1} B_k(\omega)\right)$$

where  $B_k(\omega)$  are i.i.d. random variables valued in  $[0, 1]$  and  $g$  is a function in  $C^3(\mathbb{R})$ . Then

$$\text{Var}(X) = \frac{(g'(\eta))^2 \sigma^2}{N} + O\left(\frac{1}{N^2}\right) \quad (3.43)$$

with  $\eta = \mathbb{E}(B_0)$  and  $\sigma = \sqrt{\text{Var}(B_0)}$ .

For any  $\rho$ , introduce

$$D_\rho(\omega) = X(\omega) - \rho(Y_1(\omega) - \mathbb{E}[Y_1]) \quad \text{where} \quad Y_1(\omega) = \sum_{k=0}^{N-1} B_k(\omega). \quad (3.44)$$

There exists a constant  $C$  independent of  $N$  and some deterministic parameter  $\rho_N$  such that

$$\text{Var}(D_{\rho_N}) \leq \frac{C}{N^2}. \quad (3.45)$$

The following proposition, of direct interest to us, directly falls from the above lemma.

**Proposition 3.11.** *Consider the model (3.14)–(3.15)–(3.16). Let  $A_{\eta,N}^*$  be the apparent homogenized matrix defined by (3.7)–(3.8) and  $D_\rho^{1,\eta}$  be the first-order controlled variable defined by (3.31). In the one-dimensional case, we have*

$$\mathbb{V}\text{ar}(A_{\eta,N}^*) = \frac{C}{N} + O\left(\frac{1}{N^2}\right) \quad (3.46)$$

and, for the optimal value of the deterministic parameter  $\rho$ ,

$$\min_{\rho} \mathbb{V}\text{ar}\left(D_\rho^{1,\eta}\right) = \mathbb{V}\text{ar}\left(D_{\rho^*}^{1,\eta}\right) = O\left(\frac{1}{N^2}\right). \quad (3.47)$$

Using the control variate approach based on the first-order model, the variance is thus improved by at least one order in terms of  $N$ . Note in particular that, in the above results, we have not assumed  $\eta$  to be small.

*Proof of Proposition 3.11.* The proof of (3.46) falls from (3.41) and (3.43). We now prove (3.47). In view of (3.31), (3.41), (3.42) and (3.44), we see that

$$D_\rho^{1,\eta} = X(\omega) - \rho \bar{A}_{1\text{def}}^{0,N} \left( Y_1(\omega) - \mathbb{E}[Y_1] \right).$$

Using (3.45), we thus have

$$\min_{\rho} \mathbb{V}\text{ar}\left(D_\rho^{1,\eta}\right) \leq \mathbb{V}\text{ar}\left(D_{\rho_N}\right) \leq \frac{C}{N^2},$$

which concludes the proof of Proposition 3.11.  $\square$

### Second order model

In view of (3.39), (3.30), (3.34) and (3.38), the second-order controlled variable reads

$$\begin{aligned} D_{\rho_1,\rho_2,\rho_3}^{3,\eta}(\omega) &= A_{\eta,N}^*(\omega) - \rho_1 \bar{A}_{1\text{def}}^{0,N} \sum_{k=0}^{N-1} \left( B_k^\eta(\omega) - \eta \right) \\ &\quad - \rho_2 \bar{A}_{2\text{def}}^{0,1,N} \sum_{k \neq l}^{N-1} \left( B_k^\eta(\omega) B_l^\eta(\omega) - \eta^2 \right) \\ &\quad - \rho_3 \bar{C}_{2\text{def}}^{0,1,N} \sum_{k \neq l}^{N-1} \left( (1 - B_k^\eta(\omega))(1 - B_l^\eta(\omega)) - (1 - \eta)^2 \right) \end{aligned}$$

where we have used (3.29) and the fact that, in the one-dimensional case,  $\bar{A}_{2\text{def}}^{k,l,N}$  and  $\bar{C}_{2\text{def}}^{k,l,N}$  are independent of  $k$  and  $l$ . We hence obtain that

$$D_{\rho_1,\rho_2,\rho_3}^{3,\eta}(\omega) = A_{\eta,N}^*(\omega) - \bar{\rho}_1 \sum_{k=0}^{N-1} \left( B_k^\eta(\omega) - \eta \right) - \bar{\rho}_2 \sum_{k \neq l}^{N-1} \left( B_k^\eta(\omega) B_l^\eta(\omega) - \eta^2 \right) \quad (3.48)$$

with

$$\bar{\rho}_1 = \rho_1 \bar{A}_{1\text{def}}^{0,N} - 2(N-1)\rho_3 \bar{C}_{2\text{def}}^{0,1,N}, \quad \bar{\rho}_2 = \rho_2 \bar{A}_{2\text{def}}^{0,1,N} + \rho_3 \bar{C}_{2\text{def}}^{0,1,N}.$$

We first state the following general result, the proof of which is postponed until Section 3.4.1.

**Lemma 3.12.** *Let*

$$X(\omega) = g\left(\frac{1}{N} \sum_{k=0}^{N-1} B_k(\omega)\right)$$

where  $g$  is a function in  $C^3(\mathbb{R})$  and  $B_k(\omega)$  are i.i.d. random variables taking values in  $\{0, 1\}$ . Let  $Y_1$  be defined by (3.44) and  $Y_2$  be defined by

$$Y_2(\omega) = \sum_{k=0}^{N-1} \sum_{l=0, l \neq k}^{N-1} B_k(\omega) B_l(\omega). \quad (3.49)$$

There exists a constant  $C$  independent of  $N$  and some deterministic parameters  $\bar{\rho}_1$  and  $\bar{\rho}_2$  (that depend on  $N$ ) such that

$$\text{Var}\left(\bar{D}_{\bar{\rho}_1, \bar{\rho}_2}\right) \leq \frac{C}{N^3} \quad (3.50)$$

where  $\bar{D}_{\bar{\rho}_1, \bar{\rho}_2}(\omega) = X(\omega) - \bar{\rho}_1(Y_1(\omega) - \mathbb{E}(Y_1)) - \bar{\rho}_2(Y_2(\omega) - \mathbb{E}(Y_2))$ .

The following proposition directly falls from the above lemma.

**Proposition 3.13.** *Consider the model (3.14)–(3.15)–(3.16). Let  $A_{\eta, N}^*$  be the apparent homogenized matrix defined by (3.7)–(3.8) and  $D_{\rho_1, \rho_2, \rho_3}^{3, \eta}(\omega)$  be the second-order controlled variable defined by (3.39). In the one-dimensional case, for the optimal value of the deterministic parameters  $\rho_1$ ,  $\rho_2$  and  $\rho_3$ , we have*

$$\min_{\rho_1, \rho_2, \rho_3} \text{Var}\left(D_{\rho_1, \rho_2, \rho_3}^{3, \eta}\right) = O\left(\frac{1}{N^3}\right). \quad (3.51)$$

We recall that

$$\text{Var}\left(A_{\eta, N}^*\right) = \frac{C}{N} + O\left(\frac{1}{N^2}\right).$$

Thus, using the control variate approach based on the second-order model, the variance is improved by at least two orders in terms of  $N$ . This result is to be compared with Proposition 3.11.

*Proof of Proposition 3.13.* In view of (3.48), (3.41), (3.44) and (3.49), we see that

$$D_{\rho_1, \rho_2, \rho_3}^{3, \eta}(\omega) = X(\omega) - \bar{\rho}_1(Y_1(\omega) - \mathbb{E}[Y_1]) - \bar{\rho}_2(Y_2(\omega) - \mathbb{E}[Y_2]).$$

Using (3.50), we thus have

$$\min_{\rho_1, \rho_2, \rho_3} \text{Var}\left(D_{\rho_1, \rho_2, \rho_3}^{3, \eta}\right) \leq \text{Var}\left(\bar{D}_{\bar{\rho}_1, \bar{\rho}_2}\right) \leq \frac{C}{N^3},$$

which concludes the proof of Proposition 3.13.  $\square$

### Proofs of Lemmas 3.10 and 3.12

*Proof of Lemma 3.10.* Introducing the centered random variables

$$d_k(\omega) = B_k(\omega) - \eta$$

and a smooth function  $h$  on  $[0, 1]$ , we write

$$\begin{aligned} h\left(\frac{1}{N}\sum_{k=0}^{N-1}B_k(\omega)\right) &= h\left(\eta + \frac{1}{N}\sum_{k=0}^{N-1}d_k(\omega)\right) \\ &= h(\eta) + \frac{h'(\eta)}{N}\sum_{k=0}^{N-1}d_k(\omega) + \frac{h''(\eta)}{2}\left(\frac{1}{N}\sum_{k=0}^{N-1}d_k(\omega)\right)^2 \\ &\quad + \frac{h'''(\theta_3^N(\omega))}{6}\left(\frac{1}{N}\sum_{k=0}^{N-1}d_k(\omega)\right)^3 \end{aligned} \quad (3.52)$$

for some  $\theta_3^N(\omega) \in [0, 1]$ . Recall now that any i.i.d. variables  $d_k$  with mean value zero satisfy the following bounds:

$$\forall p \in \mathbb{N}^*, \exists C_p > 0, \quad \left| \mathbb{E}\left[\left(\frac{1}{N}\sum_{k=0}^{N-1}d_k\right)^p\right]\right| \leq \begin{cases} \frac{C_p}{N^{p/2}} & \text{if } p \text{ is even;} \\ \frac{C_p}{N^{(p+1)/2}} & \text{if } p \text{ is odd.} \end{cases} \quad (3.53)$$

This is proved by developing the power  $p$  of the sum, and then using the fact that the variables are i.i.d and have mean value zero. Taking expectations in (3.52), we thus deduce that

$$\mathbb{E}\left[h\left(\frac{1}{N}\sum_{k=0}^{N-1}B_k(\omega)\right)\right] = h(\eta) + \frac{h''(\eta)}{2N}\sigma^2 + O\left(\frac{1}{N^2}\right),$$

where  $\sigma^2 = \mathbb{E}[d_0^2] = \text{Var}(B_0)$ . Choosing  $h(x) = g(x)$  and  $h(x) = (g(x))^2$ , we obtain (3.43).

We next turn to proving (3.45). As in (3.52), we have

$$\begin{aligned} X(\omega) &= g\left(\frac{1}{N}\sum_{k=0}^{N-1}B_k(\omega)\right) \\ &= g(\eta) + \frac{g'(\eta)}{N}\sum_{k=0}^{N-1}d_k(\omega) + \frac{g''(\theta_2^N(\omega))}{2}S_N(\omega) \\ &= g(\eta) + \frac{g'(\eta)}{N}\left(Y_1(\omega) - \mathbb{E}[Y_1]\right) + \frac{g''(\theta_2^N(\omega))}{2}S_N(\omega) \end{aligned}$$

for some  $\theta_2^N(\omega) \in [0, 1]$ , where  $S_N(\omega) = \left(\frac{1}{N}\sum_{k=0}^{N-1}d_k(\omega)\right)^2$ . Set  $\rho_N = \frac{g'(\eta)}{N}$ . Then

$$D_{\rho_N}(\omega) = X(\omega) - \rho_N\left(Y_1(\omega) - \mathbb{E}[Y_1]\right) = g(\eta) + \frac{g''(\theta_2^N(\omega))}{2}S_N(\omega).$$

Using (3.53), we thus obtain that

$$\text{Var}(D_{\rho_N}) \leq \mathbb{E}\left[\left(\frac{g''(\theta_2^N(\omega))}{2}S_N(\omega)\right)^2\right] \leq C\mathbb{E}[(S_N)^2] \leq \frac{C}{N^2}$$

which is the claimed bound (3.45). This concludes the proof of Lemma 3.10.  $\square$

*Proof of Lemma 3.12.* We follow the same lines as in the proof of Lemma 3.10. Introducing the centered random variables

$$d_k(\omega) = B_k(\omega) - \eta,$$

we write, as in (3.52), that

$$\begin{aligned} X(\omega) &= g\left(\frac{1}{N}\sum_{k=0}^{N-1} B_k(\omega)\right) \\ &= g(\eta) + \frac{g'(\eta)}{N}\sum_{k=0}^{N-1} d_k(\omega) + \frac{g''(\eta)}{2N^2}\left(\sum_{k=0}^{N-1} d_k(\omega)\right)^2 + \frac{g'''(\theta_3^N(\omega))}{6}S_N(\omega) \end{aligned} \quad (3.54)$$

for some  $\theta_3^N(\omega) \in [0, 1]$ , where  $S_N(\omega) = \left(\frac{1}{N}\sum_{k=0}^{N-1} d_k(\omega)\right)^3$ . We now recall that  $\sum_{k=0}^{N-1} d_k(\omega) = Y_1 - \mathbb{E}(Y_1)$ . Furthermore, we compute that

$$\left(\sum_{k=0}^{N-1} d_k(\omega)\right)^2 = N^2\eta^2 + (1 - 2N\eta)Y_1(\omega) + Y_2(\omega).$$

We thus recast (3.54) as

$$X(\omega) = \mathcal{C} + \frac{g'(\eta)}{N}Y_1(\omega) + \frac{g''(\eta)}{2N^2}\left((1 - 2N\eta)Y_1(\omega) + Y_2(\omega)\right) + \frac{g'''(\theta_3^N(\omega))}{6}S_N(\omega)$$

where  $\mathcal{C}$  is a deterministic quantity.

Set  $\bar{\rho}_1 = \frac{g'(\eta)}{N} + \frac{g''(\eta)}{2N^2}(1 - 2N\eta)$  and  $\bar{\rho}_2 = \frac{g''(\eta)}{2N^2}$ . Then

$$\bar{D}_{\bar{\rho}_1, \bar{\rho}_2}(\omega) = X(\omega) - \bar{\rho}_1\left(Y_1(\omega) - \mathbb{E}(Y_1)\right) - \bar{\rho}_2\left(Y_2(\omega) - \mathbb{E}(Y_2)\right) = \mathcal{C} + \frac{g'''(\theta_3^N(\omega))}{6}S_N(\omega).$$

Using (3.53), we thus obtain that

$$\text{Var}\left(\bar{D}_{\bar{\rho}_1, \bar{\rho}_2}\right) \leq \mathbb{E}\left[\left(\frac{g'''(\theta_3^N(\omega))}{6}S_N(\omega)\right)^2\right] \leq C\mathbb{E}\left[(S_N)^2\right] \leq \frac{C}{N^3}$$

which is the claimed bound (3.50). This concludes the proof of Lemma 3.12.  $\square$

### 3.4.2 Multi-dimensional case

#### Proof of Lemma 3.4

The proof follows the same lines as that of (3.25). It falls by enumerating the possible configurations according to the number of defects they include. We thus have, following Section 3.2.2,

$$\mathbb{E}\left[\varphi(A_{\eta, N}^*)\right] = (1 - \eta)^{|Q_N|}\varphi(A_{\text{per}}^*) + \sum_{k \in \mathcal{I}_N} \eta(1 - \eta)^{|Q_N|-1}\varphi(A_{1, k, N}^*) + O_N(\eta^2). \quad (3.55)$$

On the other hand, using (3.30) and (3.27), we write

$$\begin{aligned} &\mathbb{E}\left[\varphi\left(A_{\text{per}}^* + A_1^{\eta, N}\right)\right] \\ &= (1 - \eta)^{|Q_N|}\varphi(A_{\text{per}}^*) + \sum_{k \in \mathcal{I}_N} \eta(1 - \eta)^{|Q_N|-1}\varphi\left(A_{\text{per}}^* + \bar{A}_1^{k, N}\right) + O_N(\eta^2) \\ &= (1 - \eta)^{|Q_N|}\varphi(A_{\text{per}}^*) + \sum_{k \in \mathcal{I}_N} \eta(1 - \eta)^{|Q_N|-1}\varphi(A_{1, k, N}^*) + O_N(\eta^2). \end{aligned}$$



We deduce from the above relation and (3.55) the claimed result.

### Estimates of the variances as a function of $\eta$

Lemmas 3.4 and 3.8 show that our surrogate model is a good approximation (in terms of its law) of the random variable  $A_{\eta,N}^*$ . The lemma below shows, again in the regime  $\eta \ll 1$ , that variance is indeed decreased.

Consider any entry  $ij$  of the homogenized matrix. The estimation of  $\mathbb{E} \left[ \left( A_{\eta,N}^* \right)_{ij} \right]$  can be done by a Monte Carlo empirical mean on  $\left( A_{\eta,N}^*(\omega) \right)_{ij}$ ,  $\left( D_{\rho=1}^{1,\eta}(\omega) \right)_{ij}$  (see Section 3.3.1) or  $\left( D_{\rho_1,\rho_2}^{2,\eta}(\omega) \right)_{ij}$  (see Section 3.3.2).

**Lemma 3.14.** *For any entry  $ij$  of the homogenized matrix, we have*

$$\text{Var} \left[ \left( A_{\eta,N}^* \right)_{ij} \right] = \eta C_N^0 + O_N(\eta^2), \quad (3.56)$$

$$\text{Var} \left[ \left( D_{\rho=1}^{1,\eta} \right)_{ij} \right] = O_N(\eta^2), \quad (3.57)$$

$$\text{Var} \left[ \left( D_{\rho_1=\rho_2=1}^{2,\eta} \right)_{ij} \right] = O_N(\eta^3), \quad (3.58)$$

where  $C_N^0$  is a positive constant.

In practice, we would not necessarily work with  $\rho = 1$ , but with the optimal parameter  $\rho^*$ . A direct consequence of (3.57) is of course that

$$\text{Var} \left[ \left( D_{\rho^*}^{1,\eta} \right)_{ij} \right] = \inf_{\rho} \text{Var} \left[ \left( D_{\rho}^{1,\eta} \right)_{ij} \right] = O_N(\eta^2).$$

**Remark 3.15.** *Even though the variance of  $D_{\rho^*}^{1,\eta}$  is much smaller than that of  $A_{\eta,N}^*$ , we will see in Section 3.4.2 below that, in the regime  $\eta \ll 1$ , the weakly stochastic approximation described in Section 3.2.2 is even more efficient.*

*Proof.* We infer from (3.55) and (3.29) that, for any function  $\varphi$ ,

$$\mathbb{E} [\varphi (A_{\eta,N}^*)] = \varphi (A_{\text{per}}^*) + \eta |Q_N| \left( \varphi (A_{1,0,N}^*) - \varphi (A_{\text{per}}^*) \right) + O_N(\eta^2).$$

Taking  $\varphi(M) = M_{ij}$  and  $\varphi(M) = M_{ij}^2$ , we obtain (3.56).

We next turn to proving (3.57). For any function  $\varphi$ , we write, using (3.31) and (3.27), that

$$\begin{aligned} & \mathbb{E} \left[ \varphi \left( D_{\rho=1}^{1,\eta} \right) \right] \\ &= (1-\eta)^{|Q_N|} \varphi \left( A_{\text{per}}^* + \eta \bar{A}_1^N \right) \\ & \quad + \sum_{k \in \mathcal{I}_N} \eta (1-\eta)^{|Q_N|-1} \varphi \left( A_{1,k,N}^* - \bar{A}_1^{k,N} + \eta \bar{A}_1^N \right) + O_N(\eta^2) \\ &= (1-\eta)^{|Q_N|} \varphi \left( A_{\text{per}}^* + \eta \bar{A}_1^N \right) + \sum_{k \in \mathcal{I}_N} \eta (1-\eta)^{|Q_N|-1} \varphi \left( A_{\text{per}}^* + \eta \bar{A}_1^N \right) + O_N(\eta^2) \\ &= \varphi \left( A_{\text{per}}^* + \eta \bar{A}_1^N \right) + O_N(\eta^2). \end{aligned}$$

Taking  $\varphi(M) = M_{ij}$  and  $\varphi(M) = M_{ij}^2$ , we obtain (3.57). The proof of (3.58) follows the same lines.  $\square$

### Comparison to a weakly stochastic approach

In the regime  $\eta \ll 1$ , we have three approaches at our disposal to estimate  $\mathbb{E} [A_{\eta,N}^*]$ : the standard Monte Carlo approach, the control variate approach, and the weakly stochastic approach described in Section 3.2.2. We compare here their efficiency. Let  $\mathcal{C}_N$  be the cost to solve a single corrector problem on  $Q_N$ .

The standard Monte Carlo approach amounts to writing

$$\mathbb{E} [A_{\eta,N}^*] \approx \frac{1}{M} \sum_{m=1}^M A_{\eta,N}^{*,m}(\omega).$$

In the above approximation, the error on the entry  $ij$  is controlled by  $\sqrt{\text{Var} \left[ \left( A_{\eta,N}^* \right)_{ij} \right] / M}$ . In view of (3.56), it is thus of the order of  $\sqrt{\eta/M}$ . The cost is  $M \mathcal{C}_N$ .

The control variate approach (say using the first order surrogate model) amounts to writing

$$\mathbb{E} [A_{\eta,N}^*] \approx \frac{1}{M} \sum_{m=1}^M D_{\rho}^{1,\eta,m}(\omega),$$

where  $D_{\rho}^{1,\eta}(\omega)$  is defined by (3.31). The error is of the order of  $\sqrt{\eta^2/M}$  in view of (3.57). The cost is that of solving  $M$  corrector problems and that of determining  $\bar{A}_{1\text{def}}^{0,N}$ , namely  $(1+M)\mathcal{C}_N$ .

Using the same kind of information as in the above control variate approach, the weakly stochastic approximation (3.25) reads

$$\mathbb{E} [A_{\eta,N}^*] \approx A_{\text{per}}^* + \eta \bar{A}_1^N.$$

The error is of the order of  $\eta^2$ . The cost is that of determining  $\bar{A}_{1\text{def}}^{0,N}$ , i.e.  $\mathcal{C}_N$ .

Obviously, the control variate approach is always more efficient than the Monte Carlo approach. However, to reach the same accuracy as the weakly stochastic approach, one would need to take  $M = \eta^{-2}$  realizations, leading to a cost much larger than with the weakly stochastic approach. The same observation holds when using the control variate approach using the second order surrogate model. Therefore, in the regime  $\eta \ll 1$ , the weakly stochastic approach (3.25) is the most efficient one.

## 3.5 Numerical results

We consider the so-called random checkerboard case, in dimension  $d = 2$  (see Fig. 3.3). It falls into the framework (3.14)–(3.15)–(3.16) with

$$A_{\text{per}}(x) = \alpha \text{Id}_2 \quad \text{and} \quad C_{\text{per}}(x) = \beta \text{Id}_2. \quad (3.59)$$

In what follows, we choose  $\alpha = 3$  and  $\beta = 23$  (in Section 3.5.1) or  $\beta = 103$  (in Section 3.5.2). All variances are estimated on the basis of  $M = 100$  independent realizations.

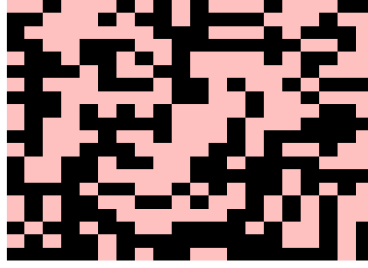


Figure 3.3 – A typical realization of the checkerboard test-case with  $\eta = 1/2$ .

### 3.5.1 Low contrast test-case

We choose here  $(\alpha, \beta) = (3, 23)$ . The motivation for this choice is that we already considered this test-case in [BCLBL12b, BCLBL12a, CLBL10] when introducing an antithetic variable approach. We are thus in position to compare the results obtained here with our previous results.

On Fig. 3.4, we plot as a function of  $\eta \in (0, 1)$  three quantities:

- the first entry of the matrix  $\mathbb{E} \left[ A_{\eta, N}^* \right]$  (obtained in practice by an expensive Monte Carlo estimation);
- the weakly stochastic approximation (3.25), which is an approximation of  $\mathbb{E} \left[ A_{\eta, N}^* \right]$  with an error of the order of  $O_N(\eta^3)$ ;
- the weakly stochastic approximation obtained in the regime  $(1 - \eta) \ll 1$ , which is an approximation of  $\mathbb{E} \left[ A_{\eta, N}^* \right]$  with an error of the order of  $O_N((1 - \eta)^3)$ .

In all cases, we work with  $N = 10$ , and the following observations are also valid for larger values of  $N$ . We see on Fig. 3.4 that, when  $\eta \leq 0.4$ , the deterministic expansion (3.25) is a very accurate approximation of  $\mathbb{E} \left[ \left( A_{\eta, N}^* \right)_{11} \right]$ . This approximation is inexpensive to compute. The same observation holds in the regime  $\eta \geq 0.7$ , where the deterministic expansion around  $\eta = 1$  provides a satisfying approximation. However, we note that none of the two weakly stochastic expansions are accurate when  $0.4 \leq \eta \leq 0.7$ . In that regime, one has to compute  $\mathbb{E} \left[ \left( A_{\eta, N}^* \right)_{11} \right]$  by considering several realizations of (3.7)–(3.8). In that regime, considering a variance reduction approach is useful.

In the regime we have identified, we show on Fig. 3.5 the ratios of variance

$$R_{\eta, N} = \frac{\text{Var} \left( \left[ A_{\eta, N}^* \right]_{11} \right)}{\text{Var}(D)}, \quad (3.60)$$

where  $D$  is either the first-order controlled variable  $D_\rho^{1, \eta}(\omega)$  defined by (3.31), or the second-order controlled variable  $D_{\rho_1, \rho_2}^{2, \eta}(\omega)$  defined by (3.35), or the controlled variable  $D_{\rho_1, \rho_2, \rho_3}^{3, \eta}(\omega)$  defined by (3.39). The parameter  $\rho$  (resp.  $(\rho_1, \rho_2)$  and  $(\rho_1, \rho_2, \rho_3)$ ) is chosen to minimize the variance of the estimator. In this section, we exactly compute (up to finite element errors) the quantities  $\bar{A}_{2, \text{def}}^{k, l, N}$  needed to build the controlled variables (3.35) and (3.39). In

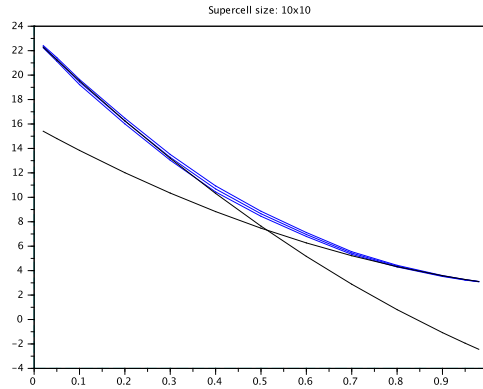


Figure 3.4 –  $\mathbb{E} \left[ \left( A_{\eta, N}^* \right)_{11} \right]$  as a function of  $\eta$ , for  $N = 10$ . Black curves: weakly stochastic approximations. Blue curve: Monte Carlo standard estimator.

Section 3.5.3 below, we approximate them using a Reduced Basis approach. We postpone until that section the discussion on computational costs and only focus here on accuracy.

**Remark 3.16.** *The second-order controlled variable  $D_{\rho_1, \rho_2}^{2, \eta}(\omega)$  defined by (3.35) is built by considering  $A_{\text{per}}$  as the reference. One could alternatively build a second-order controlled variable considering  $C_{\text{per}}$  as the reference. Numerical results obtained with such a controlled variable are similar to those obtained with  $D_{\rho_1, \rho_2}^{2, \eta}(\omega)$  (results not shown).*

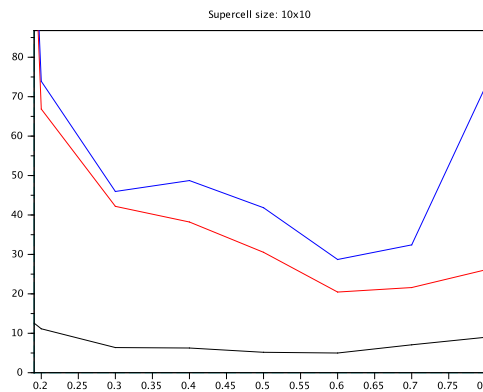


Figure 3.5 – Ratio  $R_{\eta, N}$  defined by (3.60) as a function of  $\eta$  ( $N = 10$ ). Black curve: controlled variable  $D_{\rho}^{1, \eta}(\omega)$ . Red curve: controlled variable  $D_{\rho_1, \rho_2}^{2, \eta}(\omega)$ . Blue curve: controlled variable  $D_{\rho_1, \rho_2, \rho_3}^{3, \eta}(\omega)$ .

We observe on Fig. 3.5 that, for  $\eta = 1/2$ , the approach using the first-order controlled variable (3.31) provides a variance reduction ratio (3.60) close to 6. This gain is close to the gain obtained using an antithetic variable approach (see [CLBL10, Table 2]). In contrast, when using the controlled variable (3.39) taking into account first order and second order corrections with respect to both the cases  $\eta = 0$  and  $\eta = 1$ , we obtain a gain close to 40.

We now monitor how the gain depends on the size of the domain  $Q_N$ . To that aim, we

show on Table 3.1 the ratio (3.60) as a function of  $N$ , for  $\eta = 1/2$ . We observe that the gain is essentially independent of  $N$ .

	$N = 6$	$N = 10$	$N = 20$	$N = 30$	$N = 50$
First order	7.57	5.18	6.55	8.51	7.34
Second order	35.9	41.8	37.6	35.6	40.4

Table 3.1 – Ratio  $R_{\eta,N}$  defined by (3.60) as a function of  $N$  ( $\eta = 1/2$ ). First order: controlled variable  $D_{\rho}^{1,\eta}(\omega)$ . Second order: controlled variable  $D_{\rho_1,\rho_2,\rho_3}^{3,\eta}(\omega)$ .

**Remark 3.17.** *In the one-dimensional case, we have shown that the variance ratio is proportional to  $N$  or  $N^2$  (see Propositions 3.11 and 3.13). In the two-dimensional case, we do not observe such an excellent behavior for our approach. The gain rather seems to be independent of  $N$  (see also Fig. 3.9). Nevertheless, the variance ratio is significantly higher than 1, making the approach definitely superior to the standard Monte Carlo approach.*

### 3.5.2 High contrast test-case

We now turn to a test-case with a larger contrast and set  $(\alpha, \beta) = (3, 103)$  in (3.59). On Fig. 3.6, we plot as a function of  $\eta \in (0, 1)$  the same three quantities as on Fig. 3.4 (again with  $N = 10$ ). We again see that, when  $0.3 \leq \eta \leq 0.7$ , none of the two weakly stochastic expansions are accurate. This is the regime we focus on.

We also show on Fig. 3.6 the ratios of variance (3.60) for the same three control variate approaches as on Fig. 3.5. We observe that, for  $\eta = 1/2$ , the approach using the controlled variable (3.39) provides a gain close to 6.7. This gain is smaller than in the case of Section 3.5.1 (the contrast is now larger), but still significant. As in the low-contrast test-case, the gain is essentially independent of  $N$ , as shown in Table 3.2.

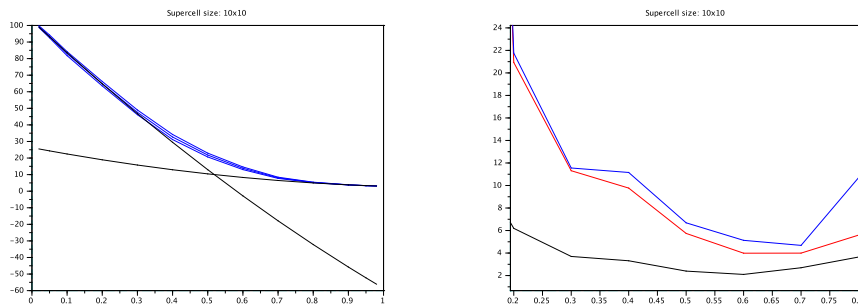


Figure 3.6 – Left:  $\mathbb{E} \left[ (A_{\eta,N}^*)_{11} \right]$  as a function of  $\eta$ , for  $N = 10$ . Blue curve: standard Monte Carlo estimator. Black curves: weakly stochastic approximations. Right: Ratio  $R_{\eta,N}$  defined by (3.60) as a function of  $\eta$  ( $N = 10$ ). Black curve: controlled variable  $D_{\rho}^{1,\eta}(\omega)$ . Red curve: controlled variable  $D_{\rho_1,\rho_2}^{2,\eta}(\omega)$ . Blue curve: controlled variable  $D_{\rho_1,\rho_2,\rho_3}^{3,\eta}(\omega)$ .

### 3.5.3 Using a Reduced Basis (RB) approach

In Sections 3.5.1 and 3.5.2, we have used the second-order surrogate model (3.39), which takes into account the contributions from pairs of defects located at any site  $k$  and  $l$ , namely

	$N = 10$	$N = 30$	$N = 50$
First order	2.40	3.62	3.87
Second order	6.69	6.32	5.82

Table 3.2 – Ratio  $R_{\eta,N}$  defined by (3.60) as a function of  $N$  ( $\eta = 1/2$ ). First order: controlled variable  $D_{\rho}^{1,\eta}(\omega)$ . Second order: controlled variable  $D_{\rho_1,\rho_2,\rho_3}^{3,\eta}(\omega)$ .

$A_{2,k,l,N}^*$  defined by (3.24) and  $C_{2,k,l,N}^*$  defined by (3.37). These quantities are deterministic, and computed beforehand. However, in practice, computing these quantities is expensive, because we have to consider all possible configurations of pairs of defects.

This high computational cost can be decreased by using the Reduced Basis (RB) approach proposed in [LBT12]. This approach amounts to solving the one-defect problem (3.20), and *a few* two-defects problems (3.23), for  $k = 0$  and  $l$  in some set  $\mathcal{N}_N \subset \mathcal{I}_N \setminus \{0\}$  (in practice, we solve (3.23) for some  $l$  close to  $k$ ). Then, it turns out that the solutions to the other two-defects problems, i.e.  $w_p^{2,k,l,N}$  for  $k = 0$  and  $l \notin \mathcal{N}_N$ , can be well-approximated on the basis of  $w_p^{1,0,N}$  and  $\left\{w_p^{2,k,l,N}\right\}_{k=0, l \in \mathcal{N}_N}$ .

In the sequel, we consider the low-contrast test-case (i.e.  $(\alpha, \beta) = (3, 23)$  in (3.59)), set  $\eta = 1/2$ , and use this RB approach in order to decrease the offline cost of our control variate approach.

### Robutness with respect to the RB basis set

First, we evaluate the robustness of the gain in variance when we approximate the quantities  $A_{2,k,l,N}^*$  and  $C_{2,k,l,N}^*$  by the above RB approach, in contrast to computing them exactly (i.e., up to a small Finite Element error). To do so, we fix  $N$  and monitor the variance ratio for the sets  $\mathcal{N}_N$  shown on Fig. 3.7. Results are given in Table 3.3. We see that the gain in variance is independent of the set  $\mathcal{N}_N$ : we can use the RB approach with a very small set of configurations for which the correctors  $w_p^{2,k,l,N}$  are exactly computed (thereby dramatically decreasing the offline computational cost), and still retain an excellent variance reduction.

	$N = 6$	$N = 20$
$\mathcal{N}_N = \mathcal{I}_N \setminus \{0\}$	35.9	37.6
Card $\mathcal{N}_N = 20$	36.1	37.6
Card $\mathcal{N}_N = 12$	35.7	37.0
Card $\mathcal{N}_N = 8$	36.6	36.5
Card $\mathcal{N}_N = 4$	36.6	37.6

Table 3.3 – Ratio  $R_{\eta,N}$  defined by (3.60) for two values of  $N$  ( $\eta = 1/2$ ), using the second order model  $D_{\rho_1,\rho_2,\rho_3}^{3,\eta}(\omega)$  defined by (3.39). The first line corresponds to the reference computation of  $A_{2,k,l,N}^*$  and  $C_{2,k,l,N}^*$ . The subsequent lines correspond to using a RB approach to compute  $A_{2,k,l,N}^*$  and  $C_{2,k,l,N}^*$ , with a decreasing set  $\mathcal{N}_N$ .

Following the above idea, we have also tested the approach when we set  $\overline{A}_{2\text{def}}^{k,l,N} = \overline{C}_{2\text{def}}^{k,l,N} = \text{Id}$  in (3.34) and (3.38) for any  $k \neq l$  (which amounts to setting  $A_{2,k,l,N}^* = \text{Id} + 2A_{1,0,N}^* - A_{\text{per}}^*$ , see (3.28)). We do not expect (and this is indeed the case) to obtain

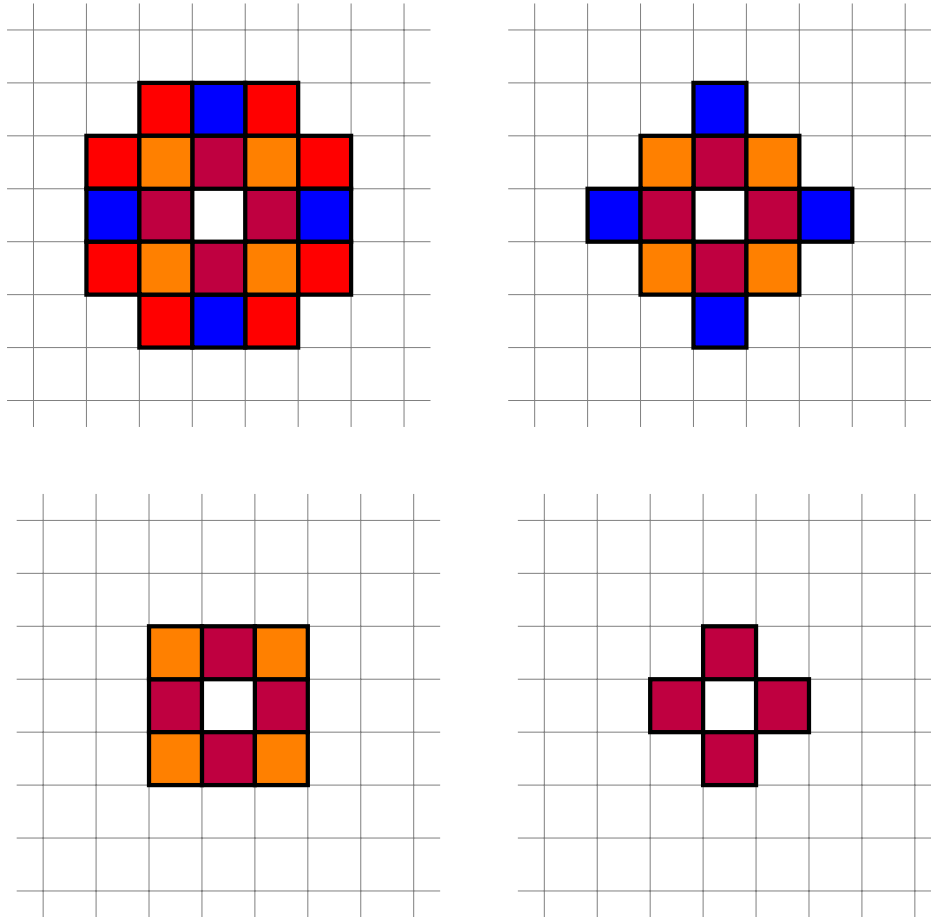


Figure 3.7 – Sets  $\mathcal{N}_N$  of position of second defect that we consider to build the RB basis set (the first defect is always in the central white cell). Top left: Card  $\mathcal{N}_N = 20$ . Top right: Card  $\mathcal{N}_N = 12$ . Bottom left: Card  $\mathcal{N}_N = 8$ . Bottom right: Card  $\mathcal{N}_N = 4$ .

good results. The controlled variable reads

$$\begin{aligned} D_{\rho_1, \rho_2, \rho_3}^{3, \eta, \text{approx}}(\omega) &= A_{\eta, N}^*(\omega) - \rho_1 \left( A_1^{\eta, N}(\omega) - \eta \bar{A}_1^N \right) \\ &\quad - \frac{\rho_2}{2} \sum_{k \neq l \in \mathcal{I}_N} \left( B_k^\eta(\omega) B_l^\eta(\omega) - \mathbb{E} [B_k^\eta B_l^\eta] \right) \\ &\quad - \frac{\rho_3}{2} \sum_{k \neq l \in \mathcal{I}_N} \left( (1 - B_k^\eta(\omega))(1 - B_l^\eta(\omega)) - \mathbb{E} [(1 - B_k^\eta)(1 - B_l^\eta)] \right) \end{aligned} \quad (3.61)$$

instead of (3.39). Computing the second order surrogate model is then extremely cheap, and as expensive as computing the first order surrogate model: one only has to solve the one-defect problem (3.20). In that case, for  $N = 20$  and  $\eta = 1/2$ , the variance ratio is equal to 6.96, which is extremely close to the variance ratio obtained by simply using the first order model (see Table 3.1), which is equal to 6.55. Considering the last two lines in (3.61) therefore does not improve the efficiency.

The above results show that it is not needed to compute with a high accuracy the quantities  $A_{2, k, l, N}^*$  and  $C_{2, k, l, N}^*$  to obtain a significant variance reduction. Using a RB approach with a very small set  $\mathcal{N}_N$  is sufficient and the gain (in terms of variance reduction) is essentially the same as that if  $A_{2, k, l, N}^*$  and  $C_{2, k, l, N}^*$  are exactly computed. However, even though the approach is quite flexible, it still requires approximations of  $A_{2, k, l, N}^*$  and  $C_{2, k, l, N}^*$  with a reasonable accuracy. Otherwise, the efficiency significantly drops down, as shown by our last test.

### Results as a function of $N$

We now fix the RB basis set corresponding to  $\text{Card } \mathcal{N}_N = 12$  on Fig. 3.7, and compare the Monte Carlo results with our control variate results, using the controlled variable (3.39). To evaluate the Monte Carlo estimator

$$I_M^{\text{MC}} = \frac{1}{M} \sum_{m=1}^M A_{\eta, N}^{*, m}(\omega),$$

we need to solve  $M$  corrector problems. In contrast, to evaluate the Control Variate estimator

$$I_M^{\text{CV}} := \frac{1}{M} \sum_{m=1}^M D_{\rho_1, \rho_2, \rho_3}^{3, \eta, m}(\omega),$$

we need to solve first the problem (3.20) and the problems (3.23) for  $k = 0$  and  $l \in \mathcal{N}_N$ , and second  $M$  corrector problems. Let  $\mathcal{C}_N$  be the cost to solve a single corrector problem on  $Q_N$ . Then the Monte Carlo cost is  $M \mathcal{C}_N$ , the Control Variate offline cost is  $(1 + \mathcal{N}_N) \mathcal{C}_N = 13 \mathcal{C}_N$ , and its online cost is  $M \mathcal{C}_N$ . In the sequel, we work with  $M = 100$ , therefore the Control Variate cost is just 13% higher than the Monte Carlo cost.

First, we plot on Fig. 3.8 the confidence intervals obtained for the Monte Carlo approach and the Control Variate approach based on (3.39). The latter confidence interval width is dramatically smaller than the former.

We next show on Fig. 3.9 the variance ratios (3.60). They somewhat vary with  $N$ . Recall that these ratios are computed on the basis of  $M = 100$  i.i.d. realizations. From one set of i.i.d. realizations to another, results may slightly vary, although qualitative conclusions remain alike. For the first order method based on (3.31), the variance ratio



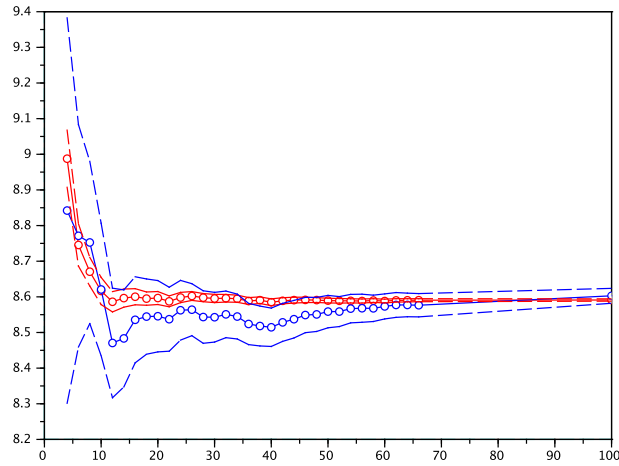


Figure 3.8 – Estimation of  $\mathbb{E} \left( \left[ A_{\eta, N}^* \right]_{11} \right)$  as a function of  $N$ . Blue: standard Monte-Carlo estimator. Red: Control Variate estimator based on (3.39). In both cases, estimators are built using  $M = 100$  i.i.d. realizations.

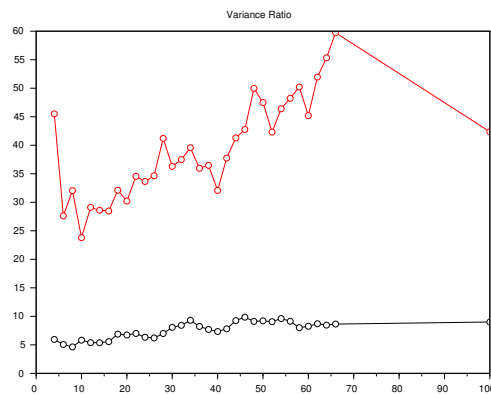


Figure 3.9 – Variance ratio (3.60) as a function of  $N$ . Black curve: using the first order controlled variable (3.31). Red curve: using the second order controlled variable (3.39). We have considered all values  $N \in \{4, 6, \dots, 66\}$  as well as  $N = 100$ .

is between 5 and 10, whereas it is around 30 or more for the second order method based on (3.39).

We plot on Fig. 3.10 the optimal values of  $\rho_1$ ,  $\rho_2$  and  $\rho_3$ , solution to (3.40). None of these parameters is close to 0: all random variables  $A_1^{\eta,N}(\omega)$ ,  $A_2^{\eta,N}(\omega)$  and  $C_2^{\eta,N}(\omega)$  are useful in (3.39) to decrease the variance.

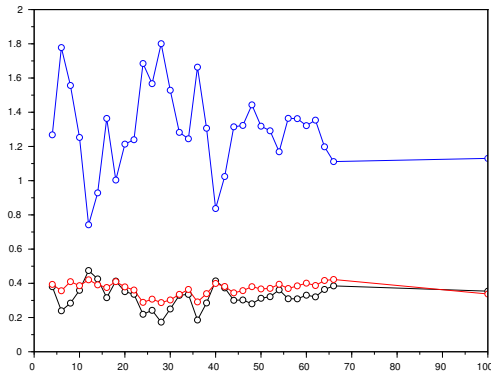


Figure 3.10 – Optimal values of  $\rho_1$  (black),  $\rho_2$  (red) and  $\rho_3$  (blue) for the controlled variable (3.39) as a function of  $N$ . We have considered all values  $N \in \{4, 6, \dots, 66\}$  as well as  $N = 100$ .

On Fig. 3.11, we eventually plot the complete errors, that is

$$e_{N,M}^{\text{MC}} = \left| \frac{1}{M} \sum_{m=1}^M A_{\eta,N}^{*,m}(\omega) - A_{\eta}^* \right|, \quad e_{N,M}^{\text{CV}} = \left| \frac{1}{M} \sum_{m=1}^M D_{\rho_1, \rho_2, \rho_3}^{3, \eta, m}(\omega) - A_{\eta}^* \right|, \quad (3.62)$$

where the exact value  $A_{\eta}^*$  is actually approximated using  $M_{\text{ref}}$  realizations on a large domain  $Q_{N_{\text{ref}}}$ . These errors are a sum of:

- the bias error  $\mathbb{E} \left[ A_{\eta,N}^* \right] - A_{\eta}^*$ ,
- the statistical error, which scales as  $\sqrt{\text{Var} \left( A_{\eta,N}^* \right) / M}$  for the Monte-Carlo approach and  $\sqrt{\text{Var} \left( D_{\rho_1, \rho_2, \rho_3}^{3, \eta} \right) / M}$  for the Control Variate approach.

When  $d \geq 3$ , the variance of  $A_{\eta,N}^*$  has been shown to scale as  $N^{-d}$  in [Nol14, Theorem 1.3 and Proposition 1.4]. For homogenization problems set on random *lattices*, optimal estimates on the above two errors have been established in [GNO14, Theorem 2] for any  $d \geq 2$ : the former scales  $N^{-d}(\ln N)^d$  while  $\text{Var} \left( A_{\eta,N}^* \right)$  scales as  $N^{-d}$ .

In the standard Monte Carlo approach, for large values of  $N$ , we expect the statistical error to dominate, and thus the error to be of the order of  $N^{-d/2}$ . This is indeed what we observe on the blue curve of Fig. 3.11. For the Control Variate approach, we observe that the error decreases as  $N^{-d}$  (see red curve of Fig. 3.11). This is consistent with the fact that, for the values of  $N$  we consider, the statistical error has been dramatically decreased and is now smaller than the bias error.

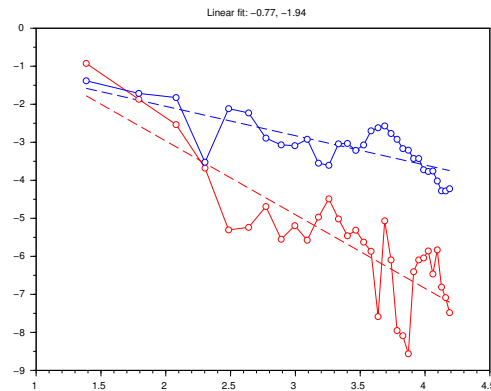


Figure 3.11 – Errors (3.62) as a function of  $N$  ( $M = 100$ ; log-log plot). Blue curve (with slope  $-0.77$ ): Monte-Carlo approach. Red curve (with slope  $-1.94$ ): Control Variate approach using (3.39). The reference value has been computed using  $N_{\text{ref}} = 100$  and  $M_{\text{ref}} = 100$ .

## Acknowledgments

The work of FL and WM is partially supported by ONR under Grant N00014-12-1-0383 and EOARD under grant FA8655-13-1-3061. WM gratefully acknowledges the support from Labex MMCD (Multi-Scale Modelling & Experimentation of Materials for Sustainable Construction) under contract ANR-11-LABX-0022. We also wish to thank Claude Le Bris and Xavier Blanc for enlightening discussions.

## Chapter 4

# Special Quasirandom Structures: a selection approach for stochastic homogenization

Ce **Chapitre** reprend l'intégralité d'un manuscrit écrit en collaboration avec Claude Le Bris et Frédéric Legoll.

Nous adaptons et étudions une approche de réduction de variance au contexte de l'homogénéisation stochastique. L'approche initiale, utilisée en physique du solide (voir [[vPDFN10](#), [WFBZ90](#), [ZWFB90](#)]) consiste à sélectionner les réalisations aléatoires qui satisfont au mieux certaines propriétés statistiques (comme la proportion volumique atteignant le taux de présence de chaque phase dans un matériau biphasique) que l'on attend qu'asymptotiquement.

Nous étudions l'approche théoriquement dans des cas simplifiés (uni dimensionnel, perturbatif en dimension plus grande) et nous démontrons son efficacité dans des cas plus généraux numériquement.

## Special Quasirandom Structures: a selection approach for stochastic homogenization

Claude Le Bris<sup>1,3</sup>, Frédéric Legoll<sup>2,3</sup> and William Minvielle<sup>1,3</sup>

legoll@lami.enpc.fr, {lebris, william.minvielle}@cermics.enpc.fr

<sup>1</sup> CERMICS, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal, 77455 Marne-La-Vallée Cedex 2, France;

<sup>2</sup> Laboratoire Navier, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal, 77455 Marne-La-Vallée Cedex 2, France ;

<sup>3</sup> INRIA Rocquencourt, MATHERIALS research-team, Domaine de Voluceau, B.P. 105, 78153 Le Chesnay Cedex, France.

**Abstract.** We adapt and study a variance reduction approach for the homogenization of elliptic equations in divergence form. The approach, borrowed from atomistic simulations and solid-state science [vPDFN10, WFBZ90, ZWFB90], consists in selecting random realizations that best satisfy some statistical properties (such as the volume fraction of each phase in a composite material) usually only obtained asymptotically.

We study the approach theoretically in some simplified settings (one-dimensional setting, perturbative setting in higher dimensions), and numerically demonstrate its efficiency in more general cases.

## 4.1 Introduction

### 4.1.1 Overview

In this article, we adapt, theoretically study and numerically test a specific variance reduction approach for the numerical homogenization of an elliptic equation with heterogeneous coefficients.

The equation we consider is the following scalar elliptic equation in divergence form

$$-\operatorname{div} \left( A \left( \frac{\cdot}{\varepsilon}, \omega \right) \nabla u^\varepsilon(\cdot, \omega) \right) = f \quad \text{in } \mathcal{D}, \quad u^\varepsilon(\cdot, \omega) = 0 \quad \text{on } \partial\mathcal{D}, \quad (4.1)$$

set on a bounded regular domain  $\mathcal{D}$  in  $\mathbb{R}^d$  (for some  $d \geq 1$ ), with a function  $f \in H^{-1}(\mathcal{D})$  in the right-hand side. The field  $A$  is a fixed matrix-valued random field. It is assumed to be uniformly elliptic, uniformly bounded and stationary in a discrete sense. All this is made precise in Section 4.1.2. Since the coefficient  $\varepsilon$  in (4.1) is assumed small, the coefficient  $A \left( \frac{\cdot}{\varepsilon}, \omega \right)$  is oscillatory and (4.1) is challenging to solve numerically. On the other hand, the problem is theoretically well understood, as is recalled below.

In the numerical practice, the traditional approach to approximate the solution  $u^\varepsilon(\cdot, \omega)$  to (4.1) is to consider (for any  $p \in \mathbb{R}^d$ ), and solve, the so-called corrector problem

$$\begin{cases} -\operatorname{div} [A(p + \nabla w_p)] = 0 & \text{in } \mathbb{R}^d \text{ almost surely,} \\ \int_Q \mathbb{E}(\nabla w_p) = 0, & \nabla w_p \text{ is stationary in the sense of (4.5) below,} \end{cases} \quad (4.2)$$

associated to (4.1). The solution to (4.2) gives the deterministic and constant coefficient  $A^*$  of the homogenized equation that in turn serves for the approximation of (4.1). We refer to Section 4.1.2 below for details.

Since (4.2) is a problem set on the entire space  $\mathbb{R}^d$ , it is necessary to truncate it on a bounded domain, and to complement it with appropriate boundary conditions. In practice, it is standard to consider the problem

$$-\operatorname{div}(A(\cdot, \omega)(p + \nabla w_p^N(\cdot, \omega))) = 0, \quad w_p^N(\cdot, \omega) \text{ is } Q_N\text{-periodic}, \quad (4.3)$$

where, say,  $Q_N = (0, N)^d$ . The deterministic homogenized matrix  $A^*$  is then approximated by the random variable  $A_N^*(\omega)$  defined by

$$\forall p \in \mathbb{R}^d, \quad A_N^*(\omega) p = \frac{1}{|Q_N|} \int_{Q_N} A(\cdot, \omega)(p + \nabla w_p^N(\cdot, \omega)). \quad (4.4)$$

This approximate homogenized coefficient  $A_N^*(\omega)$  is then evaluated using the Monte-Carlo method. Random realizations of the environment, namely the matrix coefficient  $A(y, \omega)$ , are considered within the truncated domain  $Q_N$ . For each of these environments, (4.3) is solved, and  $A_N^*(\omega)$  is computed using (4.4). The homogenized coefficient  $A^*$  is eventually approximated as an empirical mean over several realizations of  $A_N^*(\omega)$ . More details are given below in Section 4.1.3.

The purpose of this article is to reduce the variance of the approximation of  $A^*$ .

For this purpose, we borrow a variance reduction approach originally introduced in a completely different context, namely that of atomistic simulations for microscopic solid state science. In the series of articles [vPDFN10, WFBZ90, ZWFB90], an approach is indeed described that selects some particular random realizations of the environment, based on some selection criteria derived from asymptotic properties. Intuitively, the approach aims at considering only realizations that, for  $N$  fixed, *already* satisfy properties that are usually only obtained in the asymptotic limit  $N \rightarrow \infty$ . The approach carries the name SQS, abbreviation of *Special Quasirandom Structures*.

We aim at adapting this approach to our context, at studying it theoretically in some simple situations, and testing it numerically in more general situations.

For the sake of completeness, we mention that we have already studied the theoretical properties and the practical performance of several variance reduction methods for numerical random homogenization in some previous works of ours. The classical approach of *antithetic variables*, an approach that is quite generic and does not require nor exploit knowledge of the specific structure of the random problem at hand, has been considered in [BCLBL12a, BCLBL12b, CLBL10, LM15b]. The significantly more elaborate (and thus efficient) approach of *control variates* is the subject of [LM15a]. That approach requires a better knowledge of the problem considered, and is not always amenable to fully generic situations.

Our article is articulated as follows.

In the remainder of this introductory section, we present the basics of the theoretical setting (in Section 4.1.2) and of the numerical approximation method (in Section 4.1.3) for the homogenization of the random equation (4.1).

In Section 4.2, we introduce the variance reduction approach we consider. For pedagogic purposes, we first briefly expose the approach in the context of solid state physics it has originally been introduced in. This is the purpose of Section 4.2.1. In Section 4.2.2, we formally derive the specifics of our variance reduction approach using a perturbative

setting. This formal derivation provides the motivation for the general so-called SQS conditions that we use in the sequel of the work. Section 4.2.3 presents how we compute these conditions in practice. Section 4.2.4 contains the pseudo-code of our approach, along with some comments.

The theoretical analysis of the approach is the substance of Section 4.3. We begin by proving, in a fairly general situation (in any ambient dimension), that the approximation provided by our approach (at least the simplest variant of our approach) converges to the homogenized coefficient  $A^*$  when the truncated domain converges to the whole space (see Theorem 4.8 in Section 4.3.1). Next, in Section 4.3.2, we investigate more thoroughly particular and simple situations (such as the one-dimensional setting), where we can indeed completely analyze our approach and actually prove its efficiency.

Our final Section 5.4 contains numerical tests. First, it is often necessary to enforce the desired conditions up to some tolerance (see Remark 4.3 below). In Section 4.4.1, we investigate how this tolerance affects the quality of the approximation and the efficiency of the approach. We observe here that the approach is robust in this respect.

In Section 4.4.2, we illustrate on a prototypical situation the efficiency of our approach. The systematic error is not deteriorated by the approach (it might even be reduced), while the variance is reduced by several orders of magnitude. Such an efficiency is achieved at almost no additional cost with respect to the classical Monte Carlo algorithm.

### 4.1.2 Theoretical setting

To begin with, we introduce the basic setting of stochastic homogenization. We refer to the seminal works [Koz79, PV81], to [ES08] for a general, numerically oriented presentation and to [BLP78, CD99, JKO94] for classical textbooks. We also refer to [LB10] and the review article [ACLB<sup>+</sup>12] (and the extensive bibliography therein) for a presentation of our particular setting.

Throughout this article,  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space and we denote by  $\mathbb{E}(X) = \int_{\Omega} X(\omega) d\mathbb{P}(\omega)$  the expectation of any random variable  $X \in L^1(\Omega, d\mathbb{P})$ . We next fix  $d \in \mathbb{N}^*$  (the ambient physical dimension), and assume that the group  $(\mathbb{Z}^d, +)$  acts on  $\Omega$ . We denote by  $(\tau_k)_{k \in \mathbb{Z}^d}$  this action, and assume that it preserves the measure  $\mathbb{P}$ , that is, for all  $k \in \mathbb{Z}^d$  and all  $E \in \mathcal{F}$ ,  $\mathbb{P}(\tau_k E) = \mathbb{P}(E)$ . We assume that the action  $\tau$  is *ergodic*, that is, if  $E \in \mathcal{F}$  is such that  $\tau_k E = E$  for any  $k \in \mathbb{Z}^d$ , then  $\mathbb{P}(E) = 0$  or 1. In addition, we define the following notion of stationarity (see [LB10]): a function  $F \in L^1_{\text{loc}}(\mathbb{R}^d, L^1(\Omega))$  is *stationary* if

$$\forall k \in \mathbb{Z}^d, \quad F(x+k, \omega) = F(x, \tau_k \omega) \quad \text{a.e. in } x \text{ and a.s.} \quad (4.5)$$

In this setting, the ergodic theorem [Shi84] can be stated as follows:

Let  $F \in L^{\infty}(\mathbb{R}^d, L^1(\Omega))$  be a stationary random variable in the above sense. For  $k = (k_1, k_2, \dots, k_d) \in \mathbb{Z}^d$ , we set  $|k|_{\infty} = \max_{1 \leq i \leq d} |k_i|$ . Then

$$\frac{1}{(2N+1)^d} \sum_{|k|_{\infty} \leq N} F(x, \tau_k \omega) \xrightarrow{N \rightarrow \infty} \mathbb{E}(F(x, \cdot)) \quad \text{in } L^{\infty}(\mathbb{R}^d), \text{ almost surely.}$$

This implies (denoting by  $Q = (0, 1)^d$  the unit cube in  $\mathbb{R}^d$ ) that

$$F\left(\frac{x}{\varepsilon}, \omega\right) \xrightarrow{\varepsilon \rightarrow 0} \mathbb{E}\left(\int_Q F(x, \cdot) dx\right) \quad \text{in } L^{\infty}(\mathbb{R}^d), \text{ almost surely.}$$

Besides technicalities, the purpose of the above setting is simply to formalize that, even though realizations may vary, the function  $F$  at point  $x \in \mathbb{R}^d$  and the function  $F$  at point  $x + k$ ,  $k \in \mathbb{Z}^d$ , share the same law. In the homogenization context, this means that the local, microscopic environment (encoded in the matrix field  $A$  in (4.1)) is everywhere the same *on average*. From this, homogenized, macroscopic properties follow.

We consider problem (4.1), which we recall here for convenience:

$$-\operatorname{div} \left( A \left( \frac{\cdot}{\varepsilon}, \omega \right) \nabla u^\varepsilon(\cdot, \omega) \right) = f \quad \text{in } \mathcal{D}, \quad u^\varepsilon(\cdot, \omega) = 0 \quad \text{on } \partial\mathcal{D}.$$

The random matrix  $A$  is assumed stationary in the sense of (4.5). We also assume that  $A$  is bounded and coercive, that is, there exist two scalars  $0 < c \leq C < \infty$  such that, almost surely,

$$\|A(\cdot, \omega)\|_{L^\infty(\mathbb{R}^d)} \leq C \quad \text{and} \quad \forall \xi \in \mathbb{R}^d, \quad \xi^T A(x, \omega) \xi \geq c \xi^T \xi \quad \text{a.e.}$$

In this specific setting, the solution  $u^\varepsilon(\cdot, \omega)$  to (4.1) almost surely converges (when  $\varepsilon$  goes to 0) to the solution  $u^*$  to the homogenized problem

$$-\operatorname{div} (A^* \nabla u^*) = f \quad \text{in } \mathcal{D}, \quad u^* = 0 \quad \text{on } \partial\mathcal{D}. \quad (4.6)$$

The convergence of  $u^\varepsilon(\cdot, \omega)$  to  $u^*$  holds weakly in  $H^1(\mathcal{D})$  and strongly in  $L^2(\mathcal{D})$ .

The homogenized matrix  $A^*$  in (4.6) is deterministic, and given by an expectation of an integral involving the so-called corrector function, that solves a random auxiliary problem set on the *entire* space. It is given by

$$\forall p \in \mathbb{R}^d, \quad A^* p = \mathbb{E} \left[ \int_Q A(x, \cdot) (p + \nabla w_p(x, \cdot)) dx \right], \quad (4.7)$$

where we recall that  $Q = (0, 1)^d$  and where, for any vector  $p \in \mathbb{R}^d$ , the *corrector*  $w_p$  is the unique solution (up to the addition of a random constant) in  $L^2(\Omega; L^2_{\text{loc}}(\mathbb{R}^d))$  with gradient in  $L^2(\Omega; L^2_{\text{unif}}(\mathbb{R}^d))^d$  of the corrector problem (4.2). We have used the notation  $L^2_{\text{unif}}(\mathbb{R}^d)$  for the *uniform*  $L^2$  space, that is the space of functions for which, say, the  $L^2$  norm on a ball of unit size is bounded from above independently of the center of the ball.

### 4.1.3 Numerical approximation of the homogenized matrix

As briefly mentioned above, the corrector problem (4.2) is set on the *entire* space  $\mathbb{R}^d$ , and is therefore challenging to solve. Approximations are in order. In practice, the deterministic matrix  $A^*$  is approximated by the random matrix  $A_N^*(\omega)$  defined by (4.4), which is obtained by solving the corrector problem (4.3) on a *truncated* domain, say the cube  $Q_N = (0, N)^d$ . Although  $A^*$  itself is a deterministic object, its practical approximation  $A_N^*$  is random. It is only in the limit of infinitely large domains  $Q_N$  that the deterministic value is attained. As shown in [BP04], we indeed have

$$\lim_{N \rightarrow \infty} A_N^*(\omega) = A^* \quad \text{almost surely.} \quad (4.8)$$

As usual, the error  $A^* - A_N^*(\omega)$  may be expanded as

$$A^* - A_N^*(\omega) = \left( A^* - \mathbb{E}[A_N^*] \right) + \left( \mathbb{E}[A_N^*] - A_N^*(\omega) \right), \quad (4.9)$$

that is the sum of a *systematic* error and of a *statistical* error (the first and second terms in the above right-hand side, respectively).



A standard technique to compute an approximation of  $\mathbb{E}[A_N^*]$  is to consider  $M$  independent and identically distributed realizations of the field  $A$ , solve for each of them the corrector problem (4.3) (thereby obtaining i.i.d. realizations  $A_N^{*,m}(\omega)$ , for  $1 \leq m \leq M$ ) and compute the Monte Carlo approximation

$$\mathbb{E}[(A_N^*)_{ij}] \approx I_M^{\text{MC}}(\omega) := \frac{1}{M} \sum_{m=1}^M (A_N^{*,m}(\omega))_{ij} \quad (4.10)$$

for any  $1 \leq i, j \leq d$ . In view of the Central Limit Theorem, we know that  $\mathbb{E}[(A_N^*)_{ij}]$  asymptotically lies within the confidence interval

$$\left[ I_M^{\text{MC}} - 1.96 \frac{\sqrt{\text{Var}[(A_N^*)_{ij}]}}{\sqrt{M}}, I_M^{\text{MC}} + 1.96 \frac{\sqrt{\text{Var}[(A_N^*)_{ij}]}}{\sqrt{M}} \right]$$

with a probability equal to 95 %.

For simplicity, and because this is overwhelmingly the case in the numerical practice, we have considered in (4.3) *periodic* boundary conditions. These are the conditions we adopt throughout our study. It is to be remarked, however, that other boundary conditions may be employed. Likewise, other slightly modified forms of equation (4.3) may be considered. The specific choice of approximation technique is motivated by considerations about the decrease of the systematic error in (4.9). Several recent mathematical studies have clarified this issue. In addition, in the particular case of periodic boundary conditions (4.3), it has been recently established in [GNO15, Theorem 2] that the statistical error in (4.9) decays like  $N^{-d/2}$  while the systematic error in (4.9) scales as  $N^{-d}(\log N)^d$ . Both estimates have been established for the *discrete variant* of the problem. A similar decay of the statistical error has also been established for the continuous case we consider in the present article (see [GO15, Theorem 1] and [Nol14, Theorem 1.3 and Proposition 1.4]).

## 4.2 Variance reduction approach

### 4.2.1 Original formulation of the SQS approach

The variance reduction approach we elaborate upon in this article has been originally introduced for a slightly different purpose in atomistic solid-state science [vPDFN10, WFBZ90, ZWFB90].

In order to convey to the reader the intuition of the original approach, we consider here a simple one-dimensional setting, which nevertheless includes all the difficulties of a generic problem. We consider a linear chain of atomistic sites of two species  $A$  and  $B$  which interact by an interaction potential  $V_{AA}$ ,  $V_{AB}$  and  $V_{BB}$  with obvious notation. For simplicity we consider only nearest neighbour interaction. The atomic sites are occupied by a single species randomly chosen between  $A$  and  $B$ . A typical random configuration of the “material” therefore reads as an infinite sequence of the type  $\cdots ABBAABBBAAAA \cdots$

In order to compute the energy per unit particle of that atomistic system, one has to consider all possible such infinite sequences, and for each of them its normalized energy

$$\lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{i=-N}^N V_{X_{i+1}X_i}, \quad (4.11)$$

where  $X_i$  denotes the species present at the  $i$ -th site for that particular configuration ( $X_i \equiv A$  or  $B$ ). The “energy” of the system is then defined as the *expectation* of (4.11) over all possible configurations. Other quantities than (4.11) may be considered, or may be simultaneously considered.

In practice, one considers a presumably extremely large, finite  $N$ , truncates the infinite sequence over the finite length  $2N + 1$ , and compute

$$\frac{1}{2N + 1} \sum_{i=-N}^N V_{X_{i+1}X_i}$$

for many (say  $M$ , where  $M$  is also presumably large) configurations.

The approach introduced in [vPDFN10, WFBZ90, ZWFB90] consists in *selecting* specific configurations  $(X_i)_{-N \leq i \leq N}$  of atomic sites that satisfy statistical properties usually obtained only in the limit of infinitely large  $N$ .

The first such statistical property is the volume fraction, namely the proportion of species ( $A, B$ ) present on average. If the sites are all occupied randomly with probability  $1/2$  of  $A$  and  $1/2$  of  $B$  (and assuming that all these random variables are independent), then obviously the volume fraction of  $A$  is  $1/2$  and so is that of  $B$ . Then, one only consider truncated sequences  $(X_i)_{-N \leq i \leq N}$  that *exactly* reproduce that volume fraction.

Similarly, again for such an evenly distributed proportion of  $A$  and  $B$ , the energy of the entire infinite system evidently reads as

$$\mathcal{E} = \frac{1}{4} [V_{AA} + 2V_{AB} + V_{BB}]$$

(recall that we only consider nearest-neighbour interactions). Thus, one only considers truncated sequences  $(X_i)_{-N \leq i \leq N}$  which, in addition to exhibiting the exact volume fraction, have an average energy  $\frac{1}{2N + 1} \sum_{i=-N}^N V_{X_{i+1}X_i}$  which is *equal* to  $\mathcal{E}$ . And so on and so forth for other quantities of interest.

Mathematically, this *selection* of suitable configurations among all the possible configurations classically considered in a Monte-Carlo sample amounts to replacing the computation of an expectation by that of a *conditional expectation*.

The simplistic model we have just considered for pedagogic purposes can of course be replaced by more elaborate models, with more sophisticated quantities to compute, and more demanding statistical quantities to condition the computations with. The bottom line of the approach remains the same, and we adapt it to design a variance reduction approach for numerical random homogenization.

In the next section, we derive the appropriate conditions, which we call the SQS conditions, for our specific context.

### 4.2.2 Formal derivation of the SQS conditions using a perturbative setting

The purpose of this Section is to formally derive the SQS conditions that we use in the sequel. Such conditions can be easily intuitively understood. We however believe it is interesting to (formally) *derive* them in a particular case. The case we proceed with is a perturbative setting (although, we emphasize it, the conditions will be employed even in the full general, not necessarily perturbative, setting).

We assume throughout this section that the matrix valued coefficient  $A$  in (4.1) reads as

$$A_\eta(x, \omega) = C_0(x, \omega) + \eta \chi(x, \omega) C_1(x, \omega) \quad (4.12)$$

for some presumably small scalar coefficient  $\eta$ , where

- $C_0$  and  $C_1$  are two stationary, uniformly bounded matrix fields,
- $C_0 - C_1$  and  $C_0 + C_1$  are coercive,
- $\chi$  is a stationary scalar field with values in  $[-1, 1]$ .

Under these assumptions, for any  $\eta \in (-1, 1)$ , the matrix  $A_\eta$  is stationary, bounded and coercive. Intuitively, when  $\eta$  is small,  $A_\eta$  is a perturbation of the matrix-valued field  $C_0(x, \omega)$ .

**Remark 4.1.** *The expression (4.12) models e.g. a two-phase composite material, where the phases are modelled by the coefficients  $C_0$  and  $C_1$ , while  $\chi$  is the indicator function of the first phase.*

Let  $p \in \mathbb{R}^d$ . The corrector problem (4.2) reads, in this particular setting, as

$$\begin{cases} -\operatorname{div} [(C_0 + \eta\chi C_1)(p + \nabla w_\eta)] = 0 & \text{in } \mathbb{R}^d, \\ \mathbb{E} \int_Q \nabla w_\eta = 0, \quad \nabla w_\eta \text{ is stationary in the sense of (4.5),} \end{cases} \quad (4.13)$$

and the homogenized matrix (4.7) is given by

$$\forall p \in \mathbb{R}^d, \quad A_\eta^* p = \mathbb{E} \int_Q A_\eta(p + \nabla w_\eta). \quad (4.14)$$

Note that, for the sake of clarity, we omit to write the dependency of  $w_\eta$  with respect to  $p$ .

The truncated version of (4.13) on the supercell  $Q_N$  is

$$\begin{cases} -\operatorname{div} [(C_0 + \eta\chi C_1)(p + \nabla w_\eta^N)] = 0 & \text{in } Q_N, \\ w_\eta^N \text{ is } Q_N\text{-periodic,} \end{cases} \quad (4.15)$$

and we approach the homogenized matrix (4.14) by

$$\forall p \in \mathbb{R}^d, \quad A_\eta^{*,N}(\omega) p = \frac{1}{|Q_N|} \int_{Q_N} A_\eta(\cdot, \omega) (p + \nabla w_\eta^N(\cdot, \omega)). \quad (4.16)$$

### Expansion in powers of $\eta$

As  $\eta$  goes to 0, we may now expand  $A_\eta^{*,N}$  and  $A_\eta^*$  in powers of  $\eta$ . This expansion is classical (see for instance [BCLBL12b, Cos12]). We only provide it here for the sake of consistency. The corrector expands as

$$\nabla w_\eta = \nabla w_0 + \eta \nabla u_1 + \eta^2 \nabla u_2 + o(\eta^2). \quad (4.17)$$

This expansion holds in  $L^2(\Omega; L^2_{\text{unif}}(\mathbb{R}^d))$ . The functions  $w_0$ ,  $u_1$  and  $u_2$  appearing in the expansion are respectively defined by the following systems of equations:

$$\begin{cases} -\operatorname{div} [C_0(p + \nabla w_0)] = 0 & \text{in } \mathbb{R}^d, \\ \mathbb{E} \int_Q \nabla w_0 = 0, \quad \nabla w_0 \text{ is stationary,} \end{cases} \quad (4.18)$$

$$\begin{cases} -\operatorname{div} [C_0 \nabla u_1] = \operatorname{div} [\chi C_1 (p + \nabla w_0)] & \text{in } \mathbb{R}^d, \\ \mathbb{E} \int_Q \nabla u_1 = 0, \quad \nabla u_1 \text{ is stationary,} \end{cases} \quad (4.19)$$

and

$$\begin{cases} -\operatorname{div} [C_0 \nabla u_2] = \operatorname{div} [\chi C_1 \nabla u_1] & \text{in } \mathbb{R}^d, \\ \mathbb{E} \int_Q \nabla u_2 = 0, \quad \nabla u_2 \text{ is stationary.} \end{cases}$$

Inserting the expansion (4.12) of  $A_\eta$  and (4.17) of  $w_\eta$  in (4.14), we obtain

$$A_\eta^* = A_0^* + \eta A_1^* + \eta^2 A_2^* + o(\eta^2), \quad (4.20)$$

with, for any  $p \in \mathbb{R}^d$ ,

$$\begin{aligned} A_0^* p &= \mathbb{E} \left[ \int_Q C_0 (p + \nabla w_0) \right], \\ A_1^* p &= \mathbb{E} \left[ \int_Q \chi C_1 (p + \nabla w_0) \right] + \mathbb{E} \left[ \int_Q C_0 \nabla u_1 \right], \\ A_2^* p &= \mathbb{E} \left[ \int_Q \chi C_1 \nabla u_1 \right] + \mathbb{E} \left[ \int_Q C_0 \nabla u_2 \right]. \end{aligned} \quad (4.21)$$

Likewise, we expand  $w_\eta^N$  as

$$\nabla w_\eta^N = \nabla w_0^N + \eta \nabla u_1^N + \eta^2 \nabla u_2^N + o(\eta^2),$$

with

$$\begin{cases} -\operatorname{div} [C_0 (p + \nabla w_0^N)] = 0 & \text{in } Q_N, \\ w_0^N \text{ is } Q_N\text{-periodic,} \end{cases} \quad (4.22)$$

$$\begin{cases} -\operatorname{div} [C_0 \nabla u_1^N] = \operatorname{div} [\chi C_1 (p + \nabla w_0^N)] & \text{in } Q_N, \\ u_1^N \text{ is } Q_N\text{-periodic,} \end{cases} \quad (4.23)$$

and

$$\begin{cases} -\operatorname{div} [C_0 \nabla u_2^N] = \operatorname{div} [\chi C_1 \nabla u_1^N] & \text{in } Q_N, \\ u_2^N \text{ is } Q_N\text{-periodic.} \end{cases}$$

The homogenized matrix  $A_\eta^{*,N}(\omega)$  therefore satisfies

$$\left| A_\eta^{*,N}(\omega) - \left[ A_0^{*,N}(\omega) + \eta A_1^{*,N}(\omega) + \eta^2 A_2^{*,N}(\omega) \right] \right| \leq C \eta^3, \quad (4.24)$$

where  $C$  is independent of  $\eta$ ,  $N$  and  $\omega$ , and where the matrices  $A_0^{*,N}(\omega)$ ,  $A_1^{*,N}(\omega)$  and  $A_2^{*,N}(\omega)$  are defined by

$$\begin{aligned} A_0^{*,N}(\omega) p &= \frac{1}{|Q_N|} \int_{Q_N} C_0 (p + \nabla w_0^N), \\ A_1^{*,N}(\omega) p &= \frac{1}{|Q_N|} \int_{Q_N} \chi C_1 (p + \nabla w_0^N) + \frac{1}{|Q_N|} \int_{Q_N} C_0 \nabla u_1^N, \\ A_2^{*,N}(\omega) p &= \frac{1}{|Q_N|} \int_{Q_N} \chi C_1 \nabla u_1^N + \frac{1}{|Q_N|} \int_{Q_N} C_0 \nabla u_2^N. \end{aligned} \quad (4.25)$$

### SQS conditions

In line with the motivation we have mentioned above in Section 4.1.3, we are now in position to introduce the conditions that we use to *select* particular configurations of the environment within  $Q_N$  for which we compute the solution to (4.15), and, in turn, compute the approximation (4.16) of  $A_\eta^*$ . Our conditions are based upon the comparison of (4.21) and (4.25).

**Definition 4.2.** *For finite fixed  $N$ , we say that an environment  $\omega \in \Omega$  satisfies the SQS condition of*

- order 0 if  $A_0^{*,N}(\omega) = A_0^*$ , that is to say, for any  $p \in \mathbb{R}^d$ ,

$$\frac{1}{|Q_N|} \int_{Q_N} C_0(\cdot, \omega)(p + \nabla w_0^N(\cdot, \omega)) = \mathbb{E} \left[ \int_Q C_0(p + \nabla w_0) \right], \quad (4.26)$$

- order 1 if  $A_1^{*,N}(\omega) = A_1^*$ , that is to say, for any  $p \in \mathbb{R}^d$ ,

$$\begin{aligned} \frac{1}{|Q_N|} \int_{Q_N} \left[ \chi(\cdot, \omega) C_1(\cdot, \omega)(p + \nabla w_0^N(\cdot, \omega)) + C_0(\cdot, \omega) \nabla u_1^N(\cdot, \omega) \right] \\ = \mathbb{E} \left[ \int_Q \chi C_1(p + \nabla w_0) + C_0 \nabla u_1 \right], \end{aligned} \quad (4.27)$$

- order 2 if  $A_2^{*,N}(\omega) = A_2^*$ , that is to say, for any  $p \in \mathbb{R}^d$ ,

$$\begin{aligned} \frac{1}{|Q_N|} \int_{Q_N} \left[ \chi(\cdot, \omega) C_1(\cdot, \omega) \nabla u_1^N(\cdot, \omega) + C_0(\cdot, \omega) \nabla u_2^N(\cdot, \omega) \right] \\ = \mathbb{E} \left[ \int_Q \chi C_1 \nabla u_1 + C_0 \nabla u_2 \right]. \end{aligned} \quad (4.28)$$

**Remark 4.3.** *In full generality, we do not claim that there exist environments that satisfy these conditions. This might be the case that no such environment exists. One may for instance simply remark that a random variable that takes value  $-1$  and  $+1$  both with probability  $1/2$  never has value zero, which is its expectation! In some situations, we therefore have to relax the above conditions (see Section 4.2.4 below), but we temporarily leave these technicalities aside and assume that suitable environments exist.*

Consider now the two expansions (4.20) and (4.24). It is immediate to see, by subtraction, that

$$A_\eta^{*,N}(\omega) - A_\eta^* = (A_0^{*,N}(\omega) - A_0^*) + \eta(A_1^{*,N}(\omega) - A_1^*) + \eta^2(A_2^{*,N}(\omega) - A_2^*) + o(\eta^2).$$

Therefore it is readily seen that, if the configuration  $\omega$  satisfies the SQS conditions of Definition 4.2 up to the order  $k$  included ( $k = 0, 1, 2$  in our definition, but clearly one could consider higher order conditions derived likewise), then

$$A_\eta^{*,N}(\omega) - A_\eta^* = o(\eta^k), \quad (4.29)$$

where the constant in the right-hand side is independent of  $\eta$ ,  $N$  and  $\omega$ . Taking the expectation over such configurations therefore formally provides a more accurate approximation of  $A_\eta^*$ .

Now that we have derived the conditions (4.26)–(4.27)–(4.28) (which we henceforth call the *SQS conditions*) in the perturbative setting, we will actually use them in the non-perturbative setting, namely for a similar two-phase composite material, but with  $\eta$  *not* small. Of course, a property like (4.29) cannot be expected any longer since the homogenized matrix  $A^*$  is no longer a polynomial in a small coefficient that encodes a perturbation. Nevertheless, it can be expected that selecting the configurations using these conditions may improve the approximation, in particular by reducing the variance. We show in Sections 4.3 and 5.4 that it is indeed the case, theoretically and experimentally.

For the time being, we need to make a *practical* observation. The right-hand side of conditions (4.26)–(4.27)–(4.28) need to be evaluated in order to practically encode the SQS conditions. In principle, the computation of those right-hand sides are exact expectations, that can only be determined using an asymptotic limit, and are therefore almost as challenging to compute in practice as  $A^*$  itself.

We therefore need to restrict the generality of our setting (4.12) and consider cases where those right-hand sides are indeed amenable to a simple, inexpensive computation. This is the purpose of the next section.

### 4.2.3 Practical evaluation of the SQS conditions

In order to make our approach practical, we need, as mentioned above, to consider settings where the expectations present in the right-hand sides of (4.26)–(4.27)–(4.28) may be computed effectively.

#### Condition of order 0

We first consider (4.26) and its right-hand side

$$\mathbb{E} \left[ \int_Q C_0(x, \cdot) (p + \nabla w_0(x, \cdot)) dx \right]. \quad (4.30)$$

A natural assumption, which already covers a large portion of practically relevant situations, is

$$C_0(x, \omega) = C_0(x) \quad \text{is a deterministic, } \mathbb{Z}^d\text{-periodic matrix.} \quad (4.31)$$

The computation of (4.30) is then inexpensive since the solution  $w_0$  to (4.18) is in fact the deterministic solution to

$$-\operatorname{div} [C_0(p + \nabla w_0)] = 0 \quad \text{in } \mathbb{R}^d, \quad w_0 \text{ is } \mathbb{Z}^d\text{-periodic,}$$

which is unique up to the addition of a constant.

In addition, when  $N$  is an integer (and when the approximation chosen for (4.2) is the *periodic* approximation (4.3), as is indeed the case throughout this work), the solution to (4.22) is  $w_0^N \equiv w_0$  (up to an additive constant), and hence the condition (4.26) is systematically satisfied.

We henceforth assume that (4.31) holds, that  $N$  is an integer, and we proceed with the periodic approximation (4.3).

#### Condition of order 1

We next consider the SQS condition (4.27). One possible assumption to make that condition practical is

$$C_0(x, \omega) = C_0 \quad \text{is a deterministic, } \textit{constant} \text{ matrix.} \quad (4.32)$$

Since  $\nabla w_0 = 0$ , the right-hand side of (4.27) reads

$$\mathbb{E} \left[ \int_Q \chi C_1 (p + \nabla w_0) + C_0 \nabla u_1 \right] = \int_Q \mathbb{E} [\chi C_1] p + C_0 \mathbb{E} \int_Q \nabla u_1,$$

where the rightmost term vanishes in view of (4.19) and where the first term of the right-hand side may be computed using only characteristic properties of the environment considered. The condition (4.27) thus reads

$$\frac{1}{|Q_N|} \int_{Q_N} \chi(\cdot, \omega) C_1(\cdot, \omega) = \mathbb{E} \left[ \int_Q \chi C_1 \right]. \quad (4.33)$$

For instance, in a two-phase composite material mixing two *constant* and *deterministic* matrices  $C_0$  and  $C_1$ , we have

$$\mathbb{E} \left[ \int_Q \chi C_1 \right] = \mathbb{E} \left[ \int_Q \chi \right] C_1.$$

This quantity obviously only depends on the *volume fraction* of the two phases (recall (4.12)). Proceeding likewise with the left-hand side of the condition (4.27), we see that this condition reads

$$\frac{1}{|Q_N|} \int_{Q_N} \chi(x, \omega) dx = \mathbb{E} \left[ \int_Q \chi \right].$$

Interestingly (and not unexpectedly), we notice here that this condition on the volume fraction agrees with the condition we used to consider in the simple atomistic system of Section 4.2.1.

### Condition of order 2

We next proceed with condition (4.28). In addition to (4.32), we assume that

$$C_1(x, \omega) = C_1(x) \quad \text{is a deterministic, } \mathbb{Z}^d\text{-periodic matrix,} \quad (4.34)$$

and that

$$\chi(y, \omega) = \sum_{k \in \mathbb{Z}^d} X_k(\omega) \mathbb{1}_{Q+k}(y), \quad (4.35)$$

where  $X_k$  are identically distributed scalar random variables taking their values in  $[-1, 1]$ . We also assume that

$$\mathcal{C} = \sum_{k \in \mathbb{Z}^d} |\text{Cov}(X_0, X_k)| < \infty, \quad (4.36)$$

which is obviously satisfied if  $X_k$  are independent one from each other.

We then have the following result, which will be useful to make condition (4.28) practical. Its proof is postponed until Appendix 4.5.

**Lemma 4.4.** *Under the assumptions (4.32), (4.34), (4.35) and (4.36), the solution  $u_1$  to (4.19) satisfies*

$$\nabla u_1(y, \omega) = \mathbb{E}[X_0] \nabla \bar{u}_1(y) + \sum_{k \in \mathbb{Z}^d} \left( X_k(\omega) - \mathbb{E}[X_k] \right) \nabla \phi_1(y - k), \quad (4.37)$$

where  $\phi_1$  is the (unique up to the addition of a constant) solution in  $\{v \in L^2_{\text{loc}}(\mathbb{R}^d), \nabla v \in (L^2(\mathbb{R}^d))^d\}$  to

$$-\text{div} [C_0 \nabla \phi_1] = \text{div} [\mathbb{1}_Q C_1 p] \quad \text{in } \mathbb{R}^d \quad (4.38)$$

and  $\bar{u}_1$  is the (unique up to the addition of a constant) solution to

$$-\operatorname{div}[C_0 \nabla \bar{u}_1] = \operatorname{div}[C_1 p] \quad \text{in } \mathbb{R}^d, \quad \bar{u}_1 \text{ is } \mathbb{Z}^d\text{-periodic.} \quad (4.39)$$

The sum in (4.37) is a convergent series in  $L^2(Q \times \Omega)$ .

Using simpler arguments, we see that the solution  $u_1^N$  to (4.23) satisfies

$$\nabla u_1^N(y, \omega) = \mathbb{E}[X_0] \nabla \bar{u}_1(y) + \sum_{k \in \mathbb{Z}^d \cap Q_N} (X_k(\omega) - \mathbb{E}[X_k]) \nabla \phi_1^N(y - k), \quad (4.40)$$

where  $\bar{u}_1$  is defined by (4.39) and  $\phi_1^N$  is the (unique up to the addition of a constant) solution to

$$-\operatorname{div}[C_0 \nabla \phi_1^N] = \operatorname{div}[\mathbf{1}_Q C_1 p] \quad \text{in } Q_N, \quad \phi_1^N \text{ is } Q_N\text{-periodic.} \quad (4.41)$$

In practice, we can easily obtain an accurate approximation of  $\phi_1$  since the right-hand side of (4.38) has compact support. Truncating (4.38) over a sufficiently large bounded domain (with homogeneous Dirichlet boundary conditions) provides such an accurate approximation. Given (4.32), the right-hand side of Condition (4.28) rewrites  $\mathbb{E} \left[ \int_Q \chi C_1 \nabla u_1 \right]$  since  $\mathbb{E} \left[ \int_Q \nabla u_2 = 0 \right]$ . In view of (4.37), this quantity is in turn expanded as

$$\begin{aligned} & \mathbb{E} \left[ \int_Q \chi C_1 \nabla u_1 \right] \\ &= (\mathbb{E}[X_0])^2 \int_Q C_1 \nabla \bar{u}_1 + \sum_{k \in \mathbb{Z}^d} \mathbb{E} \left[ \int_Q X_0 (X_k - \mathbb{E}[X_k]) C_1(y) \nabla \phi_1(y - k) dy \right] \\ &= (\mathbb{E}[X_0])^2 \int_Q C_1 \nabla \bar{u}_1 \\ & \quad + \sum_{k \in \mathbb{Z}^d} \mathbb{E} \left[ \int_Q (X_0 - \mathbb{E}[X_0]) (X_k - \mathbb{E}[X_k]) C_1(y) \nabla \phi_1(y - k) dy \right], \end{aligned} \quad (4.42)$$

where, as mentioned above,  $\nabla \phi_1$  can be easily and accurately computed, while the series in  $k \in \mathbb{Z}^d$  may be truncated in an efficient manner because of the rapid decay at infinity of  $\nabla \phi_1$  (see [BCLBL12b, Lemma 3.1]).

We correspondingly expand the left-hand side of (4.28). The second term vanishes,



while the first term reads, in view of (4.40),

$$\begin{aligned}
& \frac{1}{|Q_N|} \int_{Q_N} \chi(y, \omega) C_1(y) \nabla u_1^N(y, \omega) dy \\
&= \sum_{j \in \mathbb{Z}^d \cap Q_N} \frac{1}{|Q_N|} \int_{Q_N} X_j(\omega) \mathbf{1}_{Q+j}(y) C_1(y) \mathbb{E}[X_0] \nabla \bar{u}_1(y) dy \\
&+ \sum_{k, j \in \mathbb{Z}^d \cap Q_N} \frac{1}{|Q_N|} \int_{Q_N} X_j(\omega) \mathbf{1}_{Q+j}(y) C_1(y) (X_k(\omega) - \mathbb{E}[X_k]) \nabla \phi_1^N(y - k) dy \\
&= (\mathbb{E}[X_0])^2 \int_Q C_1(y) \nabla \bar{u}_1(y) dy \\
&+ \mathbb{E}[X_0] \left( \frac{1}{|Q_N|} \sum_{j \in \mathbb{Z}^d \cap Q_N} (X_j(\omega) - \mathbb{E}[X_j]) \right) \int_Q C_1(y) \nabla \bar{u}_1(y) dy \\
&+ \mathbb{E}[X_0] \sum_{k \in \mathbb{Z}^d \cap Q_N} \frac{1}{|Q_N|} \int_{Q_N} C_1(y) (X_k(\omega) - \mathbb{E}[X_k]) \nabla \phi_1^N(y - k) dy \\
&+ \sum_{k, j \in \mathbb{Z}^d \cap Q_N} \frac{1}{|Q_N|} \int_{Q+j} (X_j(\omega) - \mathbb{E}[X_j]) C_1(y) (X_k(\omega) - \mathbb{E}[X_k]) \nabla \phi_1^N(y - k) dy. \tag{4.43}
\end{aligned}$$

In this particular (however still very generic) setting, we infer from (4.42) and (4.43) that Condition (4.28) reads as

$$\begin{aligned}
& \frac{1}{|Q_N|} \sum_{k, j \in Q_N \cap \mathbb{Z}^d} (X_k(\omega) - \mathbb{E}[X_k]) (X_j(\omega) - \mathbb{E}[X_j]) I_{k, j}^N \\
&+ \frac{1}{|Q_N|} \mathbb{E}[X_0] \sum_{k \in Q_N \cap \mathbb{Z}^d} (X_k(\omega) - \mathbb{E}[X_k]) I_k^N = \sum_{k \in \mathbb{Z}^d} \text{Cov}(X_0, X_k) I_k^\infty, \tag{4.44}
\end{aligned}$$

where

$$I_k^\infty = \int_{Q+k} C_1(y) \nabla \phi_1(y), \tag{4.45}$$

$$I_{k, j}^N = \int_{Q+j} C_1(y) \nabla \phi_1^N(y - k) dy, \tag{4.46}$$

$$I_k^N = \int_{Q_N} C_1(y) \nabla \phi_1^N(y - k) dy + \int_Q C_1(y) \nabla \bar{u}_1(y) dy. \tag{4.47}$$

### Summary

In the prototypical case where

$$A(x, \omega) = C_0 + \chi(x, \omega) C_1(x),$$

where  $C_0$  is constant,  $C_1$  is  $\mathbb{Z}^d$  periodic and  $\chi$  takes the form (4.35) (and where we consider the periodic approximation (4.3) of (4.2)), we have that:

- The condition (4.26) (SQS condition of order 0) is systematically fulfilled.
- In view of (4.33) and (4.35), the condition (4.27) (SQS condition of order 1) rewrites as

$$\frac{1}{|Q_N|} \sum_{k \in \mathbb{Z}^d \cap Q_N} X_k(\omega) = \mathbb{E}[X_0]. \tag{4.48}$$

- In view of (4.44), the condition (4.28) (SQS condition of order 2) writes as

$$\begin{aligned} \frac{1}{|Q_N|} \sum_{k,j \in Q_N \cap \mathbb{Z}^d} \bar{X}_k(\omega) \bar{X}_j(\omega) I_{k,j}^N \\ + \frac{1}{|Q_N|} \mathbb{E}[X_0] \sum_{k \in Q_N \cap \mathbb{Z}^d} \bar{X}_k(\omega) I_k^N = \sum_{k \in \mathbb{Z}^d} \mathbb{Cov}(X_0, X_k) I_k^\infty, \end{aligned} \quad (4.49)$$

where  $\bar{X}_k(\omega) = X_k(\omega) - \mathbb{E}[X_k]$ .

The conditions (4.48) and (4.49) are henceforth called the SQS 1 and SQS 2 conditions, respectively.

**Remark 4.5.** If (4.48) is satisfied, then the coefficient  $I_k^N$  in (4.49) can be replaced by

$$\bar{I}_k^N = \int_{Q_N} C_1(y) \nabla \phi_1^N(y - k) dy$$

and there is no need to compute  $\bar{u}_1$ .

#### 4.2.4 Selection Monte Carlo sampling

We are now in position to describe the selection Monte Carlo sampling we employ. We recall that the classical Monte Carlo sampling reads as follows:

**Algorithm 1 (Classical Monte Carlo).**

For  $m = 1, \dots, M$ ,

1. Generate a random environment  $\omega_m$ .
2. Solve the truncated corrector problem (4.3).
3. Compute  $A_N^*(\omega_m)$ .

Compute the estimator  $\mathcal{I}_{MC}^M = \frac{1}{M} \sum_{m=1}^M A_N^*(\omega_m)$  for  $A^*$ .

In contrast, our selection Monte Carlo sampling algorithm, in the particular case described in Section 4.2.3, reads as follows:

**Algorithm 2.**

The algorithm requires a tolerance  $\text{tol} > 0$ , fixed by the user.

##### 1. Offline stage

- (a) Solve the equation (4.38).
- (b) Compute  $(I_k^\infty)_{k \in \mathbb{Z}^d}$  defined by (4.45).
- (c) Compute the right-hand side of the SQS conditions (4.48) and (4.49).
- (d) Solve the equations (4.39) and (4.41).
- (e) Compute  $(I_{k,j}^N)_{k,j \in \mathbb{Z}^d \cap Q_N}$  and  $(I_k^N)_{k \in \mathbb{Z}^d \cap Q_N}$  defined by (4.46) and (4.47).

##### 2. Online stage

For  $m = 1, \dots, M$ ,

- (a) Generate a random environment  $\omega_m$ .
- (b) Using  $I_{k,j}^N$  and  $I_k^N$ , compute the left-hand sides of (4.48) and (4.49).
- (c) If the left-hand sides differ from the right-hand sides by more than  $\text{tol}$ , return to Step 2a.
- (d) Solve the truncated corrector problem (4.3).
- (e) Compute  $A_N^*(\omega_m)$ .

Compute the estimator  $\mathcal{I}_{SQS}^M = \frac{1}{M} \sum_{m=1}^M A_N^*(\omega_m)$  for  $A^*$ .

**Remark 4.6.** As pointed out above, the series in  $k \in \mathbb{Z}^d$  in the right-hand side of (4.49) may be truncated in an efficient manner because of the rapid decay at infinity of  $\nabla\phi_1$ . Therefore only a few factors  $I_k^\infty$  have to be computed at Step 1b.

**Remark 4.7.** When several SQS conditions (in practice SQS 1 and SQS 2) have to be simultaneously satisfied, we simply add them up using some weighting parameter. We have not observed any particular sensitivity of our numerical results (collected in Section 5.4 below) with respect to the adjustment of this parameter, provided it remains not too close to 0 and 1.

We have already mentioned that, in many situations, there might not be *any* random environments that satisfy some, or all, of the SQS conditions (4.26)–(4.27)–(4.28) we wish to enforce. Therefore, some adaptation is in order, and we have used in Algorithm 2 a tolerance parameter  $\text{tol} > 0$  for the SQS conditions to be satisfied.

However, if these conditions are enforced within some given tolerance as in Algorithm 2, the following issue arises. Since the motivation for precisely considering the SQS conditions is that they are fulfilled *asymptotically*, the larger the truncated computational domain we consider (that is, the larger  $N$ ), the less restrictive the conditions are, and therefore the less effective the variance reduction is likely to be. To circumvent this difficulty, a first possibility is to consider a tolerance that decreases when the size of  $Q_N$  increases. We consider this variant in our theoretical study of Section 4.3.2 below (see formula (4.69)). More precisely, we require in Proposition 4.14 that

$$\text{the SQS condition is satisfied with the tolerance } \frac{\lambda}{\sqrt{|Q_N|}}$$

for some  $\lambda$ . In practice, implementing such a threshold is not an easy matter, as the rate and the constants need to be adequately adjusted. In order to avoid such technicalities, we prefer to take a slightly different perspective, the purpose of which is to always select a *fixed proportion* of the original sample of the  $\mathcal{M}$  environments drawn. Practically, we pick the  $M$  configurations that best satisfy the SQS conditions, for some  $M$  fixed (given  $\mathcal{M}$ ).

The practical algorithm we employ is therefore as follows:

**Algorithm 3 (Selection Monte Carlo sampling).**

*The algorithm requires a number of trials  $\mathcal{M}$ , fixed by the user.*

1. **Offline stage 1:** same as the offline stage of Algorithm 2.

2. **Offline stage 2: selection step**

For  $m = 1, \dots, \mathcal{M}$ ,

- (a) Generate a random environment  $\omega_m$ .
- (b) Using  $I_{k,j}^N$  and  $I_k^N$ , compute the left-hand sides of (4.48) and (4.49).
- (c) Compute the error  $\text{error}_m$  between the left-hand sides and the right-hand sides of (4.48) and (4.49).

Sort the random environments  $(\omega_m)_{1 \leq m \leq M}$  according to  $\text{error}_m$ . Keep the  $M$  best realizations, and reject the others.

### 3. Online stage: resolution

For  $m = 1, \dots, M$ ,

- (a) Solve the truncated corrector problem (4.3).
- (b) Compute  $A_N^*(\omega_m)$ .

Compute the estimator  $\mathcal{I}_{SQS}^M = \frac{1}{M} \sum_{m=1}^M A_N^*(\omega_m)$  for  $A^*$ .

We wish to make a couple of comments about this selection Monte Carlo approach.

In full generality, the cost of Monte Carlo approaches is usually dominated by the cost of draws, and therefore selection algorithms are targeted to reject as few draws as possible.

In the present context, where boundary value problems such as (4.3) are to be solved repeatedly, the cost of draws for the environment is negligible in front of the cost of the solution procedure for such boundary value problems. Likewise, evaluating the quantities present in e.g. (4.49) is not expensive. Therefore, the purpose of the selection mechanism is to limit the number of boundary value problems to be solved, even though this comes at the (tiny) price of rejecting many environments. This also explains why we employ a simplistic rejection procedure for the selection, while in other situations of Monte Carlo samplings, one would invest in a more clever selection procedure.

A second observation is that, as potentially for any selection procedure, our selection introduces a bias (i.e. a modification of the systematic error in (4.9)). The point is to ensure that the gain in variance superseeds the bias introduced by the variance reduction approach.

Our next section addresses some theoretical aspects of our approach.

## 4.3 Elements of theoretical analysis

This section contains some elements of analysis that we are able to provide. We begin with a (somewhat) general result of convergence, and next, in some simplified cases, study our approach more thoroughly.

### 4.3.1 Proof of convergence of the approach

Formally, our approach consists in replacing an empirical average provided by the classical Monte Carlo approach to compute  $\mathbb{E}[A_N^*]$  by an empirical average *restricted* to some environments within  $Q_N$  satisfying some additional condition(s) (see Section 4.2.4). We work at a fixed size  $N$  of the truncation domain  $Q_N$  and recall that  $A_N^*(\omega)$  is defined by (4.4). Mathematically, our approach amounts to considering conditional expectations of the type  $\mathbb{E}[A_N^* \mid \text{SQS}]$ , where SQS encodes that one, or several, of the conditions summarized in (4.48)–(4.49) are satisfied.

The least we can expect from our approach is that it converges to the correct limit when  $N \rightarrow \infty$ , namely  $A^*$ , as in (4.8).

The theorem we now state establishes this fact. In order to prove it, we need to make some assumptions on our setting (see the details below), and also to make specific the SQS conditions we use. In Theorem 4.8 below, we specifically use the SQS 1 condition, in the form (4.48).

In order to state a result as general as possible, we therefore consider a condition that reads  $\frac{1}{|Q_N|} \sum_{k \in \mathbb{Z}^d \cap Q_N} f(X_k) = \mathbb{E}[f(X_0)]$  for some function  $f$ . In practice, our specific SQS 1 condition (4.48) corresponds to the choice  $f(x) = x$ .

**Theorem 4.8.** *Let  $(X_k)_{k \in \mathbb{Z}^d}$  be a sequence of independent and identically distributed scalar random variables following a common law  $\mu$ . We assume that  $\mu$  is absolutely continuous with respect to the Lebesgue measure on  $\mathbb{R}$ , and that, for any  $k \in \mathbb{Z}^d$ ,  $X_k(\omega) \in [-1, 1]$  almost surely. We consider the stationary random field*

$$A(y, \omega) = C_0 + \sum_{k \in \mathbb{Z}^d} X_k(\omega) \mathbb{1}_{Q+k}(y) C_1(y),$$

where  $C_0$  is constant and  $C_1$  is  $\mathbb{Z}^d$ -periodic and bounded. We also assume that  $C_0 + C_1(y)$  and  $C_0 - C_1(y)$  are uniformly coercive, and that  $C_0$  and  $C_1$  are symmetric.

Let  $f : \mathbb{R} \mapsto \mathbb{R}$  be a measurable function with compact level sets. We assume that  $f$  is not constant. Then we have

$$\mathbb{E} \left[ A_N^* \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \xrightarrow{N \rightarrow \infty} A^*, \quad (4.50)$$

where  $A_N^*(\omega)$  is defined by (4.4) and  $A^*$  is defined by (4.7).

Some remarks are in order.

**Remark 4.9.** *As is the case throughout this article, we have considered the periodic approximation (4.3) of (4.2). The proof of Theorem 4.8 actually carries over to the case of Neumann or Dirichlet boundary conditions, or any alternate truncation problem that provides some  $A^{*,N}(\omega)$  such that  $A_{\text{Neu}}^{*,N}(\omega) \leq A^{*,N}(\omega) \leq A_{\text{Dir}}^{*,N}(\omega)$  (see additional details in [Min15, Appendix]).*

**Remark 4.10.** *The assumptions regarding independence of the  $X_k$ , absolute continuity of their common law with respect to the Lebesgue measure and compactness of the level sets of  $f$  are necessary for technical reasons, since we need to apply a general result from [BO14]. See below for details.*

The proof of Theorem 4.8 is based on the following result, which is a particular case of a more general result due to C. Bernardin and S. Olla (see [BO14, Theorem B.2.2]):

**Theorem 4.11** (C. Bernardin and S. Olla, [BO14]). *Consider  $n$  scalar random variables  $X_1, \dots, X_n$ , that are independent and that all share the same probability distribution  $\mu(x) dx$  on  $\mathbb{R}$ . Consider a measurable function  $f : \mathbb{R} \mapsto \mathbb{R}$ , which is assumed to be not constant and to have compact level sets. Let  $f_0 = \mathbb{E}[f(X_1)] = \int_{\mathbb{R}} f(x) \mu(x) dx$ . Consider also a bounded and continuous function  $F : \mathbb{R}^k \mapsto \mathbb{R}$ . Then*

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[ F(X_1, \dots, X_k) \mid \frac{1}{n} \sum_{i=1}^n f(X_i) = f_0 \right] = \mathbb{E} [F(X_1, \dots, X_k)]. \quad (4.51)$$

Note that, when  $n \rightarrow \infty$ , the quantity  $\frac{1}{n} \sum_{i=1}^n f(X_i)$  almost surely converges to  $f_0$ .

Theorem 4.11 shows that conditioning on the manifold  $\frac{1}{n} \sum_{i=1}^n f(X_i) = f_0$  does not change the value (when  $n \rightarrow \infty$ ) of the expectation of a function  $F$  of a *finite* number  $k$  of random variables.

In our context, the variable  $X_i$  is the value of the field  $A$  on the cell  $Q + i$ . The conditioning in the left-hand side of (4.51) is identical to the conditioning in the left-hand side of (4.50).

The difference between Theorem 4.11 and our result lies in the quantity of which we compute the expectation. In our case, this quantity is  $A_N^*(\omega)$ , which is (asymptotically when  $N \rightarrow \infty$ ) a function of *all* the variables  $X_i$  and not only of a *finite* number of them. We hence cannot directly use Theorem 4.11. The proof of our result essentially amounts to introducing an upper bound and a lower bound on  $A_N^*(\omega)$  that both read as a sum of functions that depend on a *finite* number of random variables (see e.g. (4.54) below). We will then be in position to apply Theorem 4.11 on these functions.

*Proof of Theorem 4.8.* We fix some  $p \in \mathbb{R}^d$ . For the sake of clarity, the approximate homogenized matrix  $A_N^*(\omega)$  defined by (4.4) is here denoted  $A_{\text{per}}^{*,N}(\omega)$ , to emphasize that we have considered periodic boundary conditions. Since the matrix  $A$  is symmetric, we have

$$p^T A_{\text{per}}^{*,N}(\omega) p = \inf \{ \mathcal{J}_{Q_N}(v, \omega), \quad v \in H_{\text{per}}^1(Q_N) \},$$

where

$$\mathcal{J}_{Q_N}(v, \omega) = \frac{1}{|Q_N|} \int_{Q_N} (p + \nabla v)^T A(\cdot, \omega) (p + \nabla v).$$

We have considered in (4.3) periodic boundary conditions. As is well-known, other boundary conditions can be used, and these alternate approximations will be useful for the proof.

**Step 1: Upper bound.** We first introduce an approximation of  $A^*$  using a truncated corrector problem complemented with homogeneous Dirichlet boundary conditions. We consider the problem

$$\begin{cases} -\operatorname{div} \left( A(\cdot, \omega) \left( p + \nabla w_{p, \text{Dir}}^N(\cdot, \omega) \right) \right) = 0 & \text{in } Q_N, \\ w_{p, \text{Dir}}^N(\cdot, \omega) = 0 & \text{on } \partial Q_N, \end{cases}$$

which yields an approximation of  $A^*$  that we denote  $A_{\text{Dir}}^{*,N}(\omega)$  and which is defined by

$$\forall p \in \mathbb{R}^d, \quad A_{\text{Dir}}^{*,N}(\omega) p = \frac{1}{|Q_N|} \int_{Q_N} A(\cdot, \omega) (p + \nabla w_{p, \text{Dir}}^N(\cdot, \omega)).$$

As shown in [BP04], we know that

$$\lim_{N \rightarrow \infty} A_{\text{Dir}}^{*,N}(\omega) = A^* \quad \text{a.s.} \quad (4.52)$$

Since  $A$  is symmetric, we have

$$p^T A_{\text{Dir}}^{*,N}(\omega) p = \inf \{ \mathcal{J}_{Q_N}(v, \omega), \quad v \in H_0^1(Q_N) \}.$$

The matrix  $A_{\text{Dir}}^{*,N}(\omega)$  is always larger (in the sense of symmetric matrices) than  $A_{\text{per}}^{*,N}(\omega)$ . Indeed, let  $v \in H_0^1(Q_N)$ , and consider its  $Q_N$ -periodic extension  $\tilde{v}$ . Then this function belongs to  $H_{\text{per}}^1(Q_N)$ . We hence have that

$$p^T A_{\text{per}}^{*,N}(\omega) p \leq \mathcal{J}_{Q_N}(\tilde{v}, \omega) = \mathcal{J}_{Q_N}(v, \omega).$$

Minimizing over  $v \in H_0^1(Q_N)$ , we get that

$$p^T A_{\text{per}}^{*,N}(\omega) p \leq p^T A_{\text{Dir}}^{*,N}(\omega) p \quad \text{a.s.} \quad (4.53)$$

Just as  $A_{\text{per}}^{*,N}(\omega)$ , the matrix  $A_{\text{Dir}}^{*,N}(\omega)$  depends on all the random variables  $X_i(\omega)$ ,  $i \in Q_N \cap \mathbb{Z}^d$ . But, thanks to the use of homogeneous Dirichlet boundary conditions, it can be bounded from above by a sum of matrices that depend only on a finite number of random variables. To show this, we proceed as follows.

For any positive integers  $N$  and  $R$ , we introduce the integer part  $M$  of  $N/R$ . Then  $Q_N$  can be decomposed into a set of cubes of size  $R^d$ , up to some boundary layer  $B_{N,R}$ :

$$Q_N = \left( \bigcup_{j \in \mathbb{Z}^d, |j| \leq M} Rj + Q_R \right) \cup B_{N,R}.$$

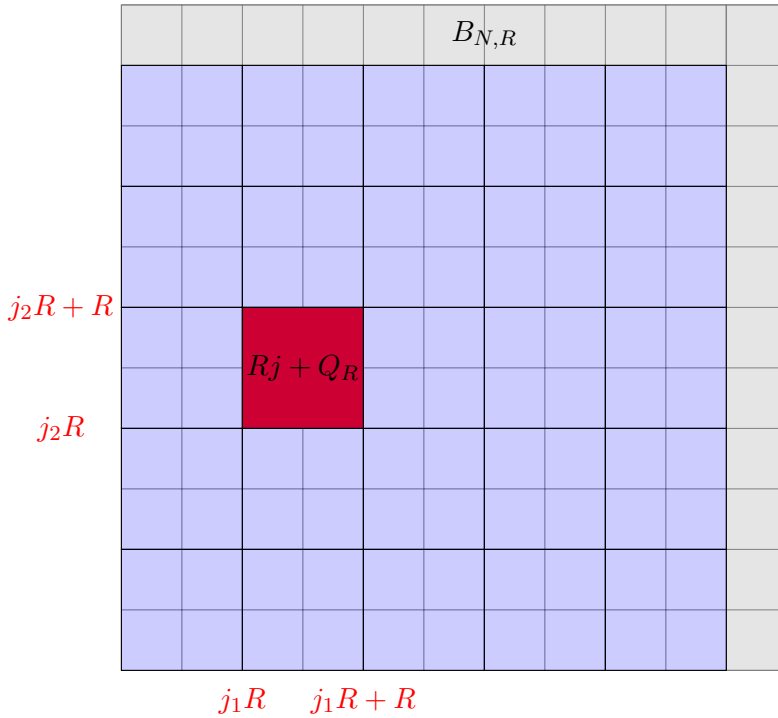


Figure 4.1 – The supercell  $Q_N$  (here represented for  $N = 11$ ) is split into cells of size  $R^d$  (here  $R = 2$ ), up to some boundary layer  $B_{N,R}$ .

For any  $j \in \mathbb{Z}^d$ ,  $|j| \leq M$ , consider a function  $v_j \in H_0^1(Rj + Q_R)$ . We now define the function  $v$  on  $Q_N$  as:

- for any  $x \in Rj + Q_R$ , we set  $v(x) = v_j(x)$ ;
- if  $x \in B_{N,R}$ , we set  $v(x) = 0$ .

The function  $v$  belongs to  $H_0^1(Q_N)$ . We hence write that

$$p^T A_{\text{Dir}}^{*,N}(\omega)p \leq \mathcal{J}_{Q_N}(v, \omega) = \frac{|Q_R|}{|Q_N|} \sum_{j \in \mathbb{Z}^d, |j| \leq M} \mathcal{J}_{Rj+Q_R}(v_j, \omega).$$

Minimizing over the functions  $v_j \in H_0^1(Rj + Q_R)$ , we hence get that

$$p^T A_{\text{Dir}}^{*,N}(\omega)p \leq \frac{|Q_R|}{|Q_N|} \sum_{j \in \mathbb{Z}^d, |j| \leq M} Y_j(\omega) \quad \text{a.s.} \quad (4.54)$$

where

$$Y_j(\omega) = \inf \left\{ \mathcal{J}_{Rj+Q_R}(v, \omega), \quad v \in H_0^1(Rj + Q_R) \right\}.$$

Since  $A$  is stationary, we note that all the random variables  $Y_j(\omega)$  share the same law. Moreover, we observe that  $Y_0(\omega) = p^T A_{\text{Dir}}^{*,R}(\omega)p$ , which is the approximation of the homogenized matrix using Dirichlet boundary conditions on  $Q_R$ .

We now take the conditional expectation of (4.54), and use the fact that the variables  $Y_j$  all share the same law:

$$\begin{aligned} \mathbb{E} \left[ p^T A_{\text{Dir}}^{*,N}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \\ \leq \frac{R^d M^d}{N^d} \mathbb{E} \left[ p^T A_{\text{Dir}}^{*,R}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right]. \end{aligned}$$

We now observe that  $p^T A_{\text{Dir}}^{*,R}(\omega)p$  only depends on a *finite* number of random variables, namely only on  $X_k(\omega)$  with  $k \in Q_R \cap \mathbb{Z}^d$ . We are thus in position to use Theorem 4.11, which yields the limit of the above right-hand side when  $N \rightarrow \infty$ . Hence, for any fixed  $R$ , we have

$$\limsup_{N \rightarrow \infty} \mathbb{E} \left[ p^T A_{\text{Dir}}^{*,N}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \leq \mathbb{E} \left[ p^T A_{\text{Dir}}^{*,R}(\omega)p \right].$$

Letting  $R$  go to  $\infty$  in the above bound and using (4.52), we obtain that

$$\limsup_{N \rightarrow \infty} \mathbb{E} \left[ p^T A_{\text{Dir}}^{*,N}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \leq p^T A^* p.$$

Using (4.53), we deduce that

$$\forall p \in \mathbb{R}^d, \quad \limsup_{N \rightarrow \infty} p^T U_N p \leq p^T A^* p, \quad (4.55)$$

where

$$U_N = \mathbb{E} \left[ A_{\text{per}}^{*,N}(\omega) \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right]. \quad (4.56)$$

**Step 2: Lower bound.** We now introduce an approximation of  $A^*$  using a truncated problem complemented with homogeneous Neumann boundary conditions. We consider the problem

$$\begin{cases} -\operatorname{div} \left( A(\cdot, \omega) (p + \nabla w_{p, \text{Neu}}^N(\cdot, \omega)) \right) = 0 & \text{in } Q_N, \\ n^T A(\cdot, \omega) (p + \nabla w_{p, \text{Neu}}^N(\cdot, \omega)) = n^T p & \text{on } \partial Q_N, \end{cases} \quad (4.57)$$



which yields an approximation of  $A^*$  that we denote  $A_{\text{Neu}}^{*,N}(\omega)$  and which is defined by

$$A_{\text{Neu}}^{*,N}(\omega) = \left( S_{\text{Neu}}^{*,N}(\omega) \right)^{-1}, \quad (4.58)$$

where  $S_{\text{Neu}}^{*,N}(\omega)$  is defined by

$$\forall p \in \mathbb{R}^d, \quad S_{\text{Neu}}^{*,N}(\omega)p = \frac{1}{|Q_N|} \int_{Q_N} p + \nabla w_{p,\text{Neu}}^N(\cdot, \omega). \quad (4.59)$$

See Remark 4.12 below for some heuristic justification of (4.58)–(4.59).

As recalled in [Min15, Appendix], we have that

$$\lim_{N \rightarrow \infty} A_{\text{Neu}}^{*,N}(\omega) = A^* \quad \text{a.s.} \quad (4.60)$$

and

$$p^T A_{\text{Neu}}^{*,N}(\omega)p \leq p^T A_{\text{per}}^{*,N}(\omega)p \quad \text{a.s.} \quad (4.61)$$

In addition, we have the following variational characterization:

$$p^T S_{\text{Neu}}^{*,N}(\omega)p = \inf \{ \mathcal{E}_{Q_N}(\sigma, \omega), \quad \sigma \in V(Q_N) \}, \quad (4.62)$$

where

$$\mathcal{E}_{Q_N}(\sigma, \omega) = \frac{1}{|Q_N|} \int_{Q_N} (p + \sigma)^T A^{-1}(\cdot, \omega)(p + \sigma)$$

and

$$V(Q_N) = \left\{ \sigma \in (L^2(Q_N))^d, \quad \text{div } \sigma = 0 \text{ in } Q_N, \quad n^T \sigma = 0 \text{ on } \partial Q_N \right\}.$$

The matrix  $S_{\text{Neu}}^{*,N}(\omega)$  (and hence the matrix  $A_{\text{Neu}}^{*,N}(\omega)$ ) depends on all the variables  $X_i(\omega)$ ,  $i \in Q_N \cap \mathbb{Z}^d$ . However, thanks to the characterization (4.62), it can be bounded from above by a sum of matrices that depend only on a finite number of random variables.

To show this, we proceed as in Step 1 of the proof. For any positive integers  $N$  and  $R$ , we introduce the integer part  $M$  of  $N/R$ , and decompose  $Q_N$  into a set of cubes of size  $R^d$ , up to some boundary layer  $B_{N,R}$  (see Figure 4.1):

$$Q_N = \left( \cup_{j \in \mathbb{Z}^d, |j| \leq M} Rj + Q_R \right) \cup B_{N,R}.$$

For any  $j \in \mathbb{Z}^d$ ,  $|j| \leq M$ , consider a function  $\sigma_j \in V(Rj + Q_R)$ . We now define the function  $\sigma$  on  $Q_N$  as:

- for any  $x \in Rj + Q_R$ , we set  $\sigma(x) = \sigma_j(x)$ ;
- if  $x \in B_{N,R}$ , we set  $\sigma(x) = 0$ .

We claim that  $\sigma \in V(Q_N)$ . We indeed first have that  $\sigma \in (L^2(Q_N))^d$ . We next consider  $\varphi \in C_0^\infty(Q_N)$  and compute that

$$\begin{aligned} \langle \text{div } \sigma, \varphi \rangle &= -\langle \sigma, \nabla \varphi \rangle \\ &= - \sum_{j \in \mathbb{Z}^d, |j| \leq M} \int_{Rj + Q_R} \sigma_j \cdot \nabla \varphi \\ &= - \sum_{j \in \mathbb{Z}^d, |j| \leq M} \int_{\partial(Rj + Q_R)} n_j^T \sigma_j \varphi \\ &= 0, \end{aligned}$$

where  $n_j$  is the outward normal to the domain  $Rj + Q_R$ . We hence have checked that  $\sigma \in V(Q_N)$ .

We next write that

$$p^T S_{\text{Neu}}^{\star, N}(\omega)p \leq \mathcal{E}_{Q_N}(\sigma, \omega) = \frac{|Q_R|}{|Q_N|} \sum_{j \in \mathbb{Z}^d, |j| \leq M} \mathcal{E}_{Rj+Q_R}(\sigma_j, \omega).$$

Minimizing over the functions  $\sigma_j \in V(Rj + Q_R)$ , we hence get that

$$p^T S_{\text{Neu}}^{\star, N}(\omega)p \leq \frac{|Q_R|}{|Q_N|} \sum_{j \in \mathbb{Z}^d, |j| \leq M} Z_j(\omega) \quad \text{a.s.} \quad (4.63)$$

where

$$Z_j(\omega) = \inf \{ \mathcal{E}_{Rj+Q_R}(\sigma, \omega), \quad \sigma \in V(Rj + Q_R) \}.$$

Since  $A$  is stationary, we note that all the random variables  $Z_j(\omega)$  share the same law. Moreover, we observe that  $Z_0(\omega) = p^T S_{\text{Neu}}^{\star, R}(\omega)p$ .

We now take the conditional expectation of (4.63), and use the fact that the variables  $Z_j$  all share the same law:

$$\begin{aligned} \mathbb{E} \left[ p^T S_{\text{Neu}}^{\star, N}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \\ \leq \frac{R^d M^d}{N^d} \mathbb{E} \left[ p^T S_{\text{Neu}}^{\star, R}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right]. \end{aligned}$$

We observe that  $p^T S_{\text{Neu}}^{\star, R}(\omega)p$  only depends on a *finite* number of random variables, namely only on  $X_k$  with  $k \in Q_R \cap \mathbb{Z}^d$ . We are thus in position to use Theorem 4.11, which yields the limit of the above right-hand side when  $N \rightarrow \infty$ . Hence, for any fixed  $R$ , we have

$$\limsup_{N \rightarrow \infty} \mathbb{E} \left[ p^T S_{\text{Neu}}^{\star, N}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \leq \mathbb{E} \left[ p^T S_{\text{Neu}}^{\star, R}(\omega)p \right].$$

Letting  $R$  go to  $\infty$  in the above bound and using (4.58) and (4.60), we obtain that

$$\limsup_{N \rightarrow \infty} \mathbb{E} \left[ p^T S_{\text{Neu}}^{\star, N}(\omega)p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \leq p^T (A^\star)^{-1} p.$$

Using (4.58) and (4.61), we deduce that

$$\limsup_{N \rightarrow \infty} \mathbb{E} \left[ p^T (A_{\text{per}}^{\star, N}(\omega))^{-1} p \mid \frac{1}{|Q_N|} \sum_{k \in Q_N \cap \mathbb{Z}^d} f(X_k) = \mathbb{E}[f(X_0)] \right] \leq p^T (A^\star)^{-1} p.$$

Using Jensen inequality, we infer from the above bound that

$$\forall p \in \mathbb{R}^d, \quad \limsup_{N \rightarrow \infty} p^T (U_N)^{-1} p \leq p^T (A^\star)^{-1} p, \quad (4.64)$$

where the matrix  $U_N$  is defined by (4.56).

**Step 3: Conclusion.** We eventually show that (4.55) and (4.64) imply that  $U_N$  converges to  $A^*$  when  $N \rightarrow \infty$ .

From the assumptions on  $A$ , we know that there exists  $0 < a_- \leq a_+ < \infty$  such that, for any  $N$  and almost surely,  $a_- \leq A_{\text{per}}^{*,N}(\omega) \leq a_+$ . Hence, for any  $N$ , the symmetric matrix  $U_N$  satisfies  $a_- \leq U_N \leq a_+$ . We can thus extract a subsequence  $U_{\varphi(N)}$  that converges to some symmetric matrix  $B$ . Let us show that  $B = A^*$ .

Let  $p \in \mathbb{R}^d$ . We first observe that, by definition,

$$\limsup_{k \rightarrow \infty} p^T U_k p \geq \lim_{k \rightarrow \infty} p^T U_{\varphi(k)} p = p^T B p.$$

We thus infer from (4.55) that

$$\forall p \in \mathbb{R}^d, \quad p^T B p \leq p^T A^* p. \quad (4.65)$$

We now proceed likewise with  $U_k^{-1}$ . We observe that,

$$\limsup_{k \rightarrow \infty} p^T U_k^{-1} p \geq \lim_{k \rightarrow \infty} p^T U_{\varphi(k)}^{-1} p = p^T B^{-1} p.$$

We thus infer from (4.64) that

$$\forall p \in \mathbb{R}^d, \quad p^T B^{-1} p \leq p^T (A^*)^{-1} p. \quad (4.66)$$

Collecting (4.65) and (4.66), we deduce that  $B = A^*$ .

The sequence  $U_N$  is bounded, and we have shown that any converging subsequence converges to  $A^*$ . This implies that  $U_N$  converges to  $A^*$  when  $N \rightarrow \infty$ , which is exactly the result (4.50). This concludes the proof of Theorem 4.8.  $\square$

**Remark 4.12.** In view of (4.57), we can check that

$$\frac{1}{|Q_N|} \int_{Q_N} A(\cdot, \omega) (p + \nabla w_{p, \text{Neu}}^N(\cdot, \omega)) = p.$$

The definition (4.58)–(4.59) can hence be understood as

$$\left\langle A(\cdot, \omega) (p + \nabla w_{p, \text{Neu}}^N(\cdot, \omega)) \right\rangle = A_{\text{Neu}}^{*,N}(\omega) \left\langle p + \nabla w_{p, \text{Neu}}^N(\cdot, \omega) \right\rangle,$$

where  $\langle \cdot \rangle = |Q_N|^{-1} \int_{Q_N} \cdot$  is the average on  $Q_N$ .

### 4.3.2 Complete analysis in some simple cases

In this section, we aim at improving the convergence result (4.50) of the previous section by quantifying both the statistical and systematic errors, in order to assess the efficiency of our approach. We are only able to proceed in simple situations where all the quantities are indeed accessible using analytic calculations. These two situations are examined in Sections 4.3.2 and 4.3.2 respectively. For the sake of brevity, and because the proofs are not very enlightening and are not likely to carry over to more general cases, we do not provide the proofs of our claims here. We refer to [Min15] where they are presented in details.

We establish below that our approach preserves the *rate* of decay of the standard Monte Carlo sampling both for the systematic and the statistical error (and thus, in particular, the systematic error remains, in rate, smaller than the statistical error). Furthermore, the *prefactor* in the statistical error is significantly reduced by our approach.

### “Zero-dimensional” homogenization

As simplest possible situation, we consider a function  $f : \mathbb{R} \mapsto \mathbb{R}$  and the random variables  $(X_i)_{1 \leq i \leq n}$ . We assume that these random variables are independent and that they are all centered Gaussian random variables with unit variance. We also assume that  $f \in C^1(\mathbb{R})$  and that  $\mathbb{E}[|f(X_1)| + |f'(X_1)|] < \infty$ . Note that it is not surprising to make some smoothness assumptions on  $f$  as we are here after *rates* of convergence, and not only a convergence result as in Section 4.3.1.

We set

$$\xi : x \mapsto \frac{1}{n} \sum_{i=1}^n x_i.$$

Assume we want to compute  $\mathbb{E}[f(X_1)]$ . A classical Monte Carlo approach would approximate this by the limit of the empirical mean  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n f(X_i(\omega))$ . In this particular instance, the simplest version of our variance reduction approach instead considers  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n f(X_i(\omega))$  for realizations  $X(\omega)$  that satisfy  $\xi(X(\omega)) = 0$ .

In this simple case, the bias of the classical approach is actually identically zero: of course,  $\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f(X_i)\right]$  does not depend on  $n$ . The statistical error is controlled by the

Central Limit Theorem and is asymptotically of order  $\sqrt{\frac{\text{Var}[f(X_1)]}{n}}$ .

**Proposition 4.13.** *Under the assumptions of this section, the bias of the selection method is of order  $1/n$ . More specifically,*

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f(X_i) \mid \xi(X) = 0\right] - \mathbb{E}[f(X_1)] = -\frac{1}{2n} \mathbb{E}[f'(X_1)] + O\left(\frac{1}{n^2}\right). \quad (4.67)$$

*The variance of the selection method is reduced by a factor independent of  $n$ . More specifically,*

$$\frac{\text{Var}\left[\frac{1}{n} \sum_{i=1}^n f(X_i) \mid \xi(X) = 0\right]}{\text{Var}\left[\frac{1}{n} \sum_{i=1}^n f(X_i)\right]} = 1 - \frac{(\mathbb{E}[f'(X_1)])^2}{\text{Var}[f(X_1)]} + O\left(\frac{1}{n}\right). \quad (4.68)$$

In view of (4.67)–(4.68), we observe that, in spite of introducing a bias of order  $O(1/n)$ , our approach reduces the statistical error from  $\frac{\lambda_{\text{MC}}}{\sqrt{n}}$  to  $\frac{\lambda_{\text{SQS}}}{\sqrt{n}}$  (with  $\lambda_{\text{SQS}} < \lambda_{\text{MC}}$ ), and therefore, for sufficiently large  $n$ , reduces the total error.

The following result covers the case where we insert a non-zero tolerance in Algorithm 2.

**Proposition 4.14.** *Under the assumptions of this section, consider the selection method where we condition on the realizations such that  $\frac{z_0}{\sqrt{n}} \leq \xi(X(\omega)) \leq \frac{z_1}{\sqrt{n}}$ , for some  $z_0$  and  $z_1 > z_0$  in  $\mathbb{R}$ . Then, for any choice of  $z_0$  and  $z_1 > z_0$ , the variance of the selection method is reduced by a factor independent of  $n$ :*

$$\frac{\text{Var}\left[\frac{1}{n} \sum_{i=1}^n f(X_i) \mid \frac{z_0}{\sqrt{n}} \leq \xi(X) \leq \frac{z_1}{\sqrt{n}}\right]}{\text{Var}\left[\frac{1}{n} \sum_{i=1}^n f(X_i)\right]} = 1 - (1 - C) \frac{(\mathbb{E}[f'(X_1)])^2}{\text{Var}[f(X_1)]} + O\left(\frac{1}{n}\right), \quad (4.69)$$

where  $C = \text{Var} \left[ X_1 \mid z_0 \leq X_1 \leq z_1 \right]$ .

The conditioning  $z_0/\sqrt{n} \leq \xi(X) \leq z_1/\sqrt{n}$  is deliberately chosen in order to match the rate of the Central Limit Theorem. It corresponds to the selection of a fixed *proportion* of samples (as in Algorithm 3 when  $\mathcal{M}$  is proportional to  $M$ ). Note that  $C > 0$ , hence the variance is less reduced than when conditioning at  $\xi(X) = 0$  (which is the case considered in Proposition 4.13). Note also that the variance is reduced (with respect to the standard Monte Carlo sampling) if, and only if,  $1 - C \geq 0$ . We are yet unable to conclude that this is the case in general. We simply note that, when  $z_1 = -z_0 > 0$ , then  $C = 1$ , yielding no gain.

### One-dimensional homogenization

In the one-dimensional case, the homogenization of a random field

$a : (y, \omega) \mapsto \sum_{i \in \mathbb{Z}} g(X_i(\omega)) \mathbb{1}_{(i, i+1)}(y)$  (where  $g$  is valued, say, in  $[a_-, a_+]$  with  $a_- > 0$ ) is a simple harmonic average. It is readily seen that

$$a_N^*(\omega) = \left( \frac{1}{N} \sum_{i=1}^N \frac{1}{g(X_k)} \right)^{-1} = \varphi \left( \frac{1}{N} \sum_{i=1}^N \frac{1}{g(X_k)} \right) \quad \text{with } \varphi(x) = 1/x.$$

Formally, the problem is thus analogous to that of the previous section, for a certain  $\varphi : \mathbb{R} \mapsto \mathbb{R}$  instead of  $\varphi = \text{Id}$ . Therefore, it is sufficient to prove consistency and variance reduction for quantities of the form  $\varphi \left( \frac{1}{N} \sum_{i=1}^N f(X_i) \right)$ .

**Proposition 4.15.** *Consider a smooth function  $\varphi : \mathbb{R} \mapsto \mathbb{R}$ . Under the assumptions of this section, the bias of the standard method and that of the selection method respectively are*

$$\mathbb{E} \left[ \varphi \left( \frac{1}{N} \sum_{i=1}^N f(X_i) \right) \right] - \varphi(f_0) = \frac{\varphi''(f_0)}{2N} \text{Var}[f(X_1)] + O \left( \frac{1}{N^2} \right) \quad (4.70)$$

and

$$\begin{aligned} & \mathbb{E} \left[ \varphi \left( \frac{1}{N} \sum_{i=1}^N f(X_i) \right) \mid \xi(X) = 0 \right] - \varphi(f_0) \\ &= \frac{\varphi''(f_0)}{2N} (\text{Var}[f(X_1)] - (\mathbb{E}[f'(X_1)])^2) - \frac{\varphi'(f_0)}{2N} \mathbb{E}[X_1 f'(X_1)] + o \left( \frac{1}{N} \right), \end{aligned} \quad (4.71)$$

with  $f_0 = \mathbb{E}[f(X_1)]$ .

The variance of the selection method is reduced by a factor independent of  $N$ :

$$\frac{\text{Var} \left[ \varphi \left( \frac{1}{N} \sum_{i=1}^N f(X_i) \right) \mid \xi(X) = 0 \right]}{\text{Var} \left[ \varphi \left( \frac{1}{N} \sum_{i=1}^N f(X_i) \right) \right]} = 1 - \frac{(\mathbb{E}[f'(X_1)])^2}{\text{Var}[f(X_1)]} + o(1). \quad (4.72)$$

To keep things simple, we do not investigate whether a more general result, accounting for some tolerance in the manner our condition is fulfilled (in the spirit of Proposition 4.14), holds here.

Proposition 4.15 shows that the bias is unchanged in rate, while the prefactor for the variance is reduced. Since the variance only decays at the rate  $1/\sqrt{N}$  while the bias decays

at the rate  $1/N$ , we see that our approach indeed reduces the total error for sufficiently large  $N$ .

In the numerical practice (mimicking in this one-dimensional setting what is actually performed for higher dimensional settings – although it is in some sense unnecessary here), we generate several, independent realizations of the  $N$ -tuples  $(X_i)_{1 \leq i \leq N}$  corresponding to as many draws of environments within the “cube”  $Q_N$ . In the classical Monte Carlo approach, we keep all such  $N$ -tuples. In our approach, we only consider those that satisfy an additional criterion.

An empirical mean (aimed at approximating  $A^*$ ) is then computed. The systematic error and the statistical error of the latter approximation are precisely related to the errors estimated in (4.70)–(4.71)–(4.72) respectively. Thus a theoretical assessment of our practical approach.

## 4.4 Numerical experiments

We first present in this section some numerical experiments that show the robustness of our variance reduction approach with respect to the tolerance with which we enforce the SQS conditions (see Section 4.4.1). We next turn to studying the performance of our approach in Section 4.4.2.

We consider the test-case when  $A$  reads as in (4.12), that is

$$A_\eta(x, \omega) = C_0(x, \omega) + \eta \chi(x, \omega) C_1(x, \omega),$$

with  $\eta = 0.5$ ,  $C_0 = C_1 = \text{Id}$ , and  $\chi$  is of the form (4.35), that is

$$\chi(y, \omega) = \sum_{k \in \mathbb{Z}^d} X_k(\omega) \mathbb{1}_{Q+k}(y).$$

The random variables  $X_k$  are i.i.d. and follow a Bernoulli law of parameter  $1/2$  valued in  $\{-1, +1\}$ . The contrast (i.e. the ratio of the largest value of  $A$  divided by its minimum value) is equal to 3. The influence of the contrast on the efficiency of our approach is investigated at the end of Section 4.4.2 (see Table 4.1). We consider there much larger values of the contrast (however all smaller than 20).

In what follows, we only consider Algorithm 3, where we take  $M = 100$  and  $\mathcal{M} = 2000$  (thus an acceptance rate of 5%).

In this setting, the SQS 1 condition as stated in (4.48) is satisfied if and only if the numbers of cells within which  $X_k(\omega) = 1$  is equal to the number of cells within which  $X_k(\omega) = -1$ . We enforce this by randomly selecting  $|Q_N|/2$  cells within the  $|Q_N|$  cells that are in  $Q_N$ , and setting  $X_k = 1$  on these cells and  $X_k = -1$  on the others.

In all our tests, we have kept the computational time fixed, or almost fixed, since the additional time needed by the selection step (namely Steps 1 and 2 of Algorithm 3) is roughly 5% of the total original computational time.

### 4.4.1 Robustness of the approach

As pointed out above, the SQS 2 condition as stated in (4.49) is only enforced in Algorithm 3 up to some tolerance. In this section, we experimentally investigate how this tolerance affects the quality of the approximation and the efficiency of the approach. To mimic the difficulty associated with the SQS 2 condition, we have also performed some

tests where we only enforce the SQS 2 condition up to some tolerance, and not exactly. The results of our numerical tests are displayed in Figures 4.2 through 4.5.

Figures 4.2 and 4.3 show the sensitivity of the variance reduction ratio upon the first order condition (4.48). Following Algorithm 3, we have sorted the realizations with respect to the error in (4.48). On Figure 4.2, the left-most dot displays the ratio  $V_{\text{SQS } 1}/V_{\text{MC}}$  between the empirical variance  $V_{\text{SQS } 1}$  among the best  $M = 100$  realizations and the reference Monte Carlo variance  $V_{\text{MC}}$ . The second dot shows the ratio between the empirical variance among the next best  $M = 100$  realizations and the reference Monte Carlo variance. We next proceed with the subsequent groups of  $M = 100$  realizations. We work with  $\mathcal{M} = 2000$ , hence there are 20 groups of 100 realizations, and hence 20 dots on Figure 4.2. On Figure 4.3, we display the same ratio of variances, but the  $x$  axis provides (for each group of  $M = 100$  realizations) the maximum error with which the first order condition (4.48) is fulfilled (rather than the index of that group, as on Figure 4.2). Hence, the first group (left-most dot) corresponds to a vanishing error, the second group corresponds to an error between 0 and  $\text{tol}$ , the third group corresponds to an error between  $\text{tol}$  and  $2\text{tol}$ , and so on and so forth.

Figures 4.4 and 4.5 show the sensitivity upon the second order condition (4.49). Here, the first order condition (4.48) is directly embedded into the random environment generator, so that every realization that is considered actually satisfies (4.48). Following Algorithm 3, we have sorted the realizations with respect to the error in (4.49). We present the results on Figures 4.4 and 4.5 following the same procedure as for Figures 4.2 and 4.3. We plot the ratio  $\frac{V_{\text{SQS } 2}}{V_{\text{exact SQS } 1}}$  between the variance  $V_{\text{SQS } 2}$  among the  $M = 100$  realizations that exactly satisfy the SQS 1 condition and best satisfy the SQS 2 condition on the one hand, and, on the other hand, the variance  $V_{\text{exact SQS } 1}$  of the realizations that exactly satisfy the SQS 1 condition.

We conclude that our approach is robust in this respect. Even if the SQS conditions (4.48)–(4.49) are not *exactly* satisfied, but only with some small tolerance, we obtain a significant variance reduction.

We next investigate how much the two SQS conditions (4.48) and (4.49) are independent (in a probabilistic sense) one from the other. Results are shown on Figure 4.6. The rightmost histogram in black is the distribution of the criterion SQS 1 (namely, the left-hand side of (4.48)) among all realizations (we have considered  $M = 100$  samples). The left histogram in black is the distribution of the criterion SQS 2 (namely, the left-hand side of (4.49)) among all realizations (we have used the same samples for both black histograms). Finally, the left histogram in blue is the conditional distribution of the criterion SQS 2 among  $M = 100$  realizations that exactly satisfy the criterion SQS 1. These two histograms on the left of the figure are sufficiently close to each other to state that conditioning with respect to SQS 1 does not change the distribution of the SQS 2 criterion. Therefore, enforcing first the condition (4.48) does not make challenging the subsequent selection of realizations that best satisfy (4.49).

#### 4.4.2 Efficiency of the approach

In this section, we investigate how the efficiency of our approach depends (i) on the size of the truncated domain  $Q_N$  and (ii) on the contrast in  $A$ .

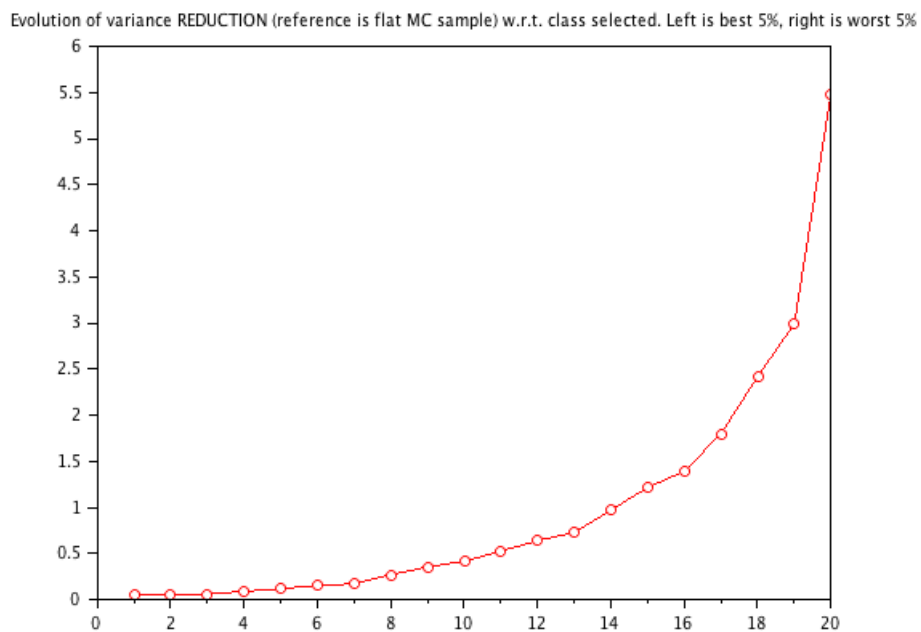


Figure 4.2 – Variance ratio  $V_{\text{SQS } 1}/V_{\text{MC}}$  for the 20 groups of realizations (sorted according to their SQS 1 error)

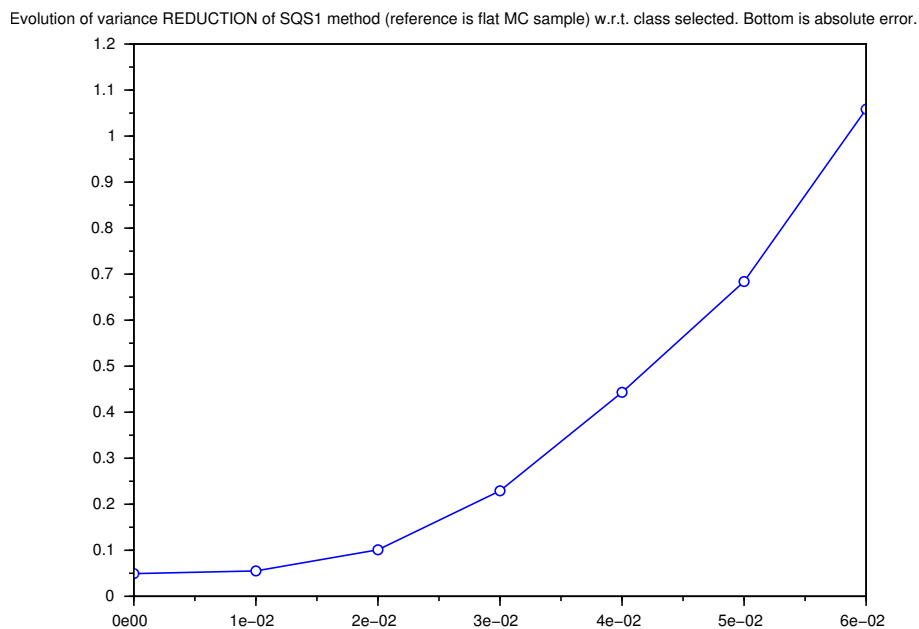


Figure 4.3 – Variance ratio  $V_{\text{SQS } 1}/V_{\text{MC}}$  as a function of the error in (4.48)



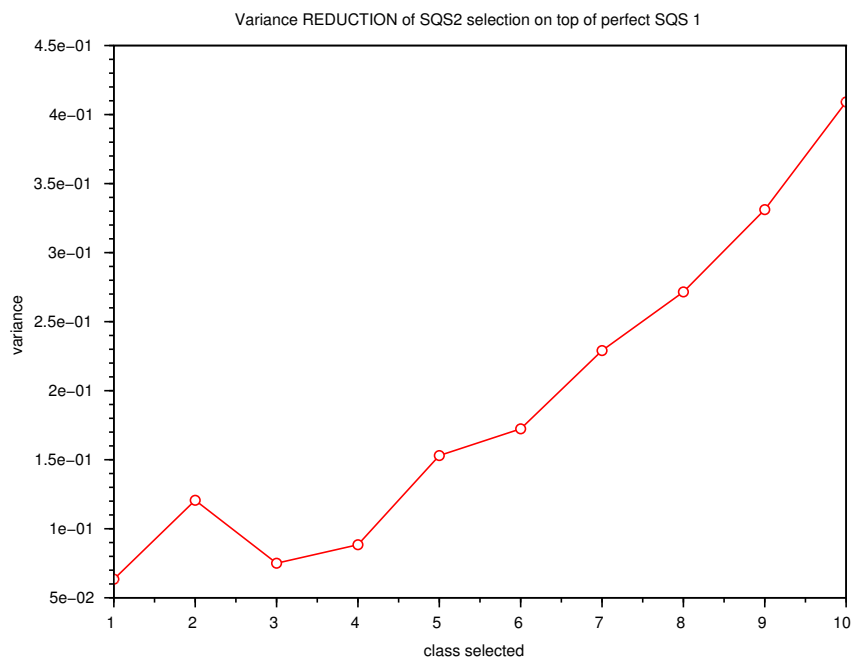


Figure 4.4 – Variance ratio  $\frac{V_{\text{SQS 2}}}{V_{\text{exact SQS 1}}}$  for the different groups of realizations (sorted according to their SQS 2 error; the SQS 1 condition is exactly satisfied). Only the 10 best groups are shown.

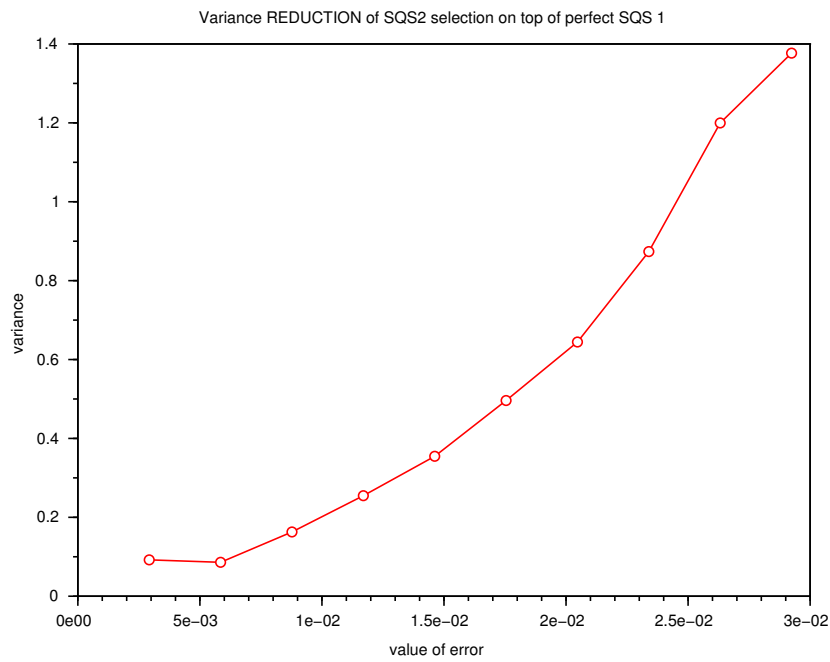


Figure 4.5 – Variance ratio  $\frac{V_{\text{SQS 2}}}{V_{\text{exact SQS 1}}}$  as a function of the error in (4.49) (the condition (4.48) is exactly satisfied). Only the 10 best groups are shown.

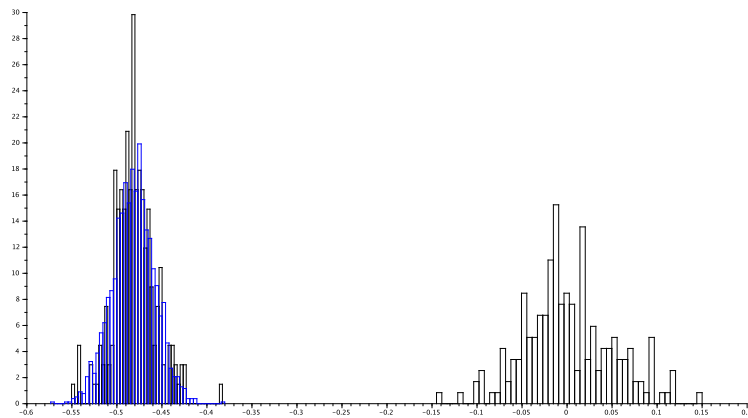


Figure 4.6 – Left: Empirical probability distribution function of the SQS 2 criterion (black histogram: no conditioning; blue histogram: the samples exactly satisfy the SQS 1 criterion). Right: empirical probability distribution of the SQS 1 criterion.

### Efficiency with respect to $N$

Figure 4.7 shows the different estimators of the first entry  $[A^*]_{11}$  of the homogenized matrix and their respective confidence intervals. The black curve is the standard Monte Carlo estimator defined by (4.10). The variance is large. In red, we display the estimator obtained by selecting realizations that exactly satisfy the SQS 1 condition. The variance is much smaller, leading in turn to a narrower confidence interval. In blue we display the estimator obtained with realizations satisfying exactly the SQS 1 condition and selected according to the SQS 2 condition (see Algorithm 3). The variance is much smaller than when using the SQS 1 approach, even when the supercell size  $N$  is small.

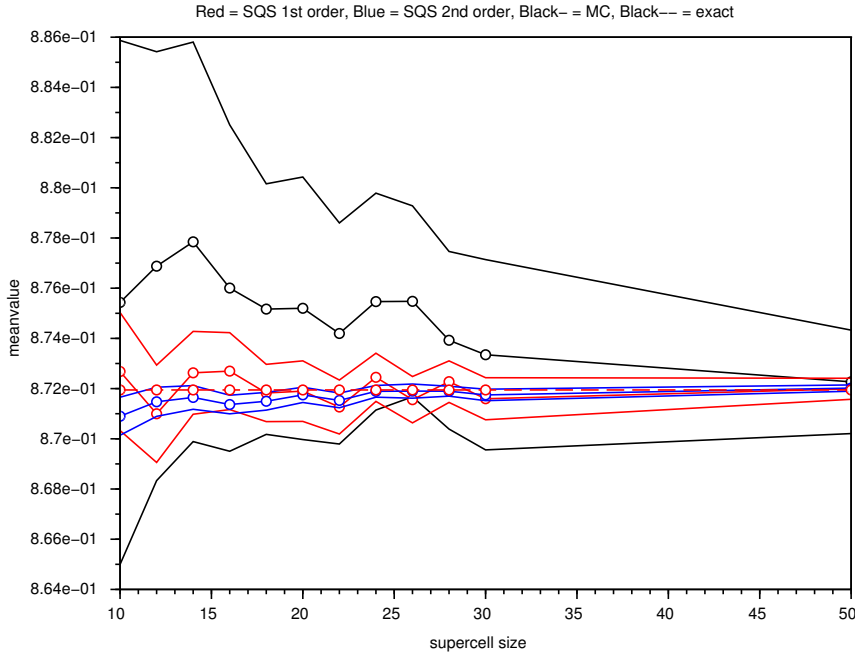


Figure 4.7 – Estimators of  $[A^*]_{11}$  (along with confidence intervals) as a function of the supercell side length  $N$ . Black curve: Monte Carlo method. Red curve: SQS 1 method. Blue curve: SQS 2 method (see text).

Figure 4.8 shows a representation of the total error as a function of the supercell size. The reference result is defined as follows. In addition to the fact that we work on a finite domain  $Q_N$  and with a finite number of samples, there is of course a slight finite element error due to the finiteness of the meshsize used to solve (4.3). This is why we do not take as reference the exact homogenized matrix, which is here known and equal to  $A^* = \sqrt{(1 + \eta)(1 - \eta)} \text{Id}$ . We rather take as reference the empirical expectation of  $A_{N_{\text{ref}}}^*(\omega)$  over  $\mathcal{M}_{\text{ref}} = 2000$  random realizations exactly satisfying the SQS 1 condition, and for the largest supercell size we have considered, that is  $N_{\text{ref}} = 50$ . The reference value is thus defined as

$$A_{\text{ref}}^* = \mathbb{E} [A_{N_{\text{ref}}}^* \mid \text{SQS 1 condition is exactly satisfied}],$$

which is in practice computed as an empirical mean over  $\mathcal{M}_{\text{ref}} = 2000$  realizations.

On the black curve, we see the total error of the standard Monte Carlo method, defined as

$$\text{total error} = \left| \frac{1}{M} \sum_{m=1}^M A_N^*(\omega_m) - A_{\text{ref}}^* \right|.$$

The two other curves show the same quantity, where the  $M$  environments considered now satisfy exactly the SQS 1 condition (red curve) and additionally the SQS 2 condition (blue curve). We observe that our two methods (SQS 1 and SQS 2) yield a total error roughly 7 times smaller than the Monte Carlo error. Note that we are not able to distinguish between our two methods because the reference value  $A_{\text{ref}}^*$  is computed with parameters  $N_{\text{ref}}$  and  $\mathcal{M}_{\text{ref}}$  not large enough for that purpose.

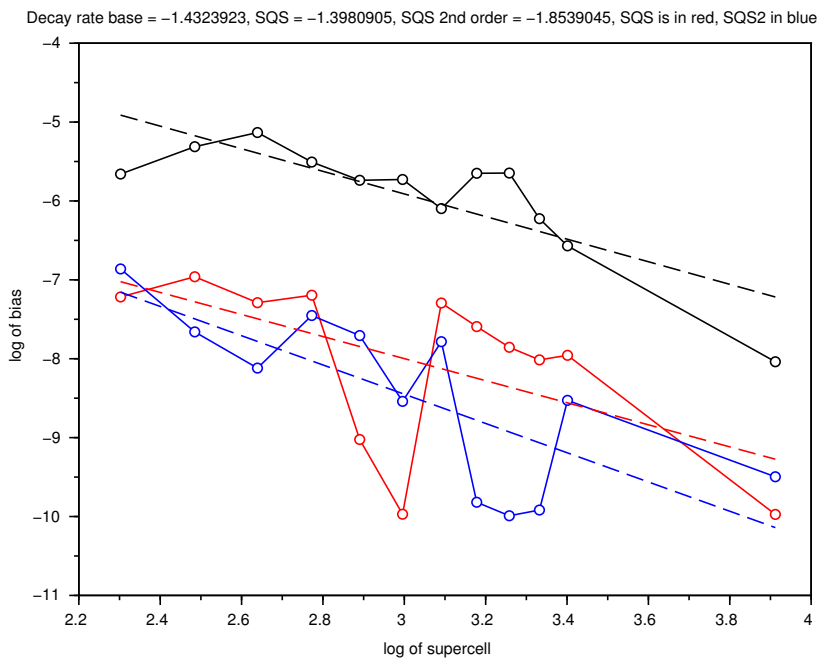


Figure 4.8 – log-log plot of the total error as a function of the supercell side length  $N$  (natural logarithm). Black curve: Monte Carlo method. Red curve: SQS 1 method. Blue curve: SQS 2 method (see text).

Figure 4.9 shows the empirical variance of the different estimators of  $[A^*]_{11}$  as a function of the supercell size. The black curve is the standard Monte Carlo estimator defined by (4.10). In red, we display the estimator obtained by selecting realizations that exactly satisfy the SQS 1 condition. In blue, we display the estimator obtained with realizations exactly satisfying the SQS 1 condition and selected according to the SQS 2 condition (see Algorithm 3). We observe that, each time we consider an additional SQS condition, the empirical variance of the estimator is significantly reduced (even if this SQS condition is not exactly enforced; recall that we consider here only the 5 % best samples in terms of the SQS 2 condition, but that we are unable to enforce it exactly). On our test-case, enforcing the SQS 1 condition leads to a variance 20 times smaller than that of the standard Monte Carlo approach, while additionally enforcing the SQS 2 condition leads to an additional variance reduction of a factor of 10.

We also observe on Figure 4.9 that all variances decay as  $\lambda/|Q_N|$ , where

$$\lambda_{\text{SQS 2}} < \lambda_{\text{exact SQS 1}} < \lambda_{\text{MC}}.$$

This corroborates in higher dimension the behaviour predicted in Section 4.3.2. In particular, the gain in variance does not decrease when the supercell becomes larger.

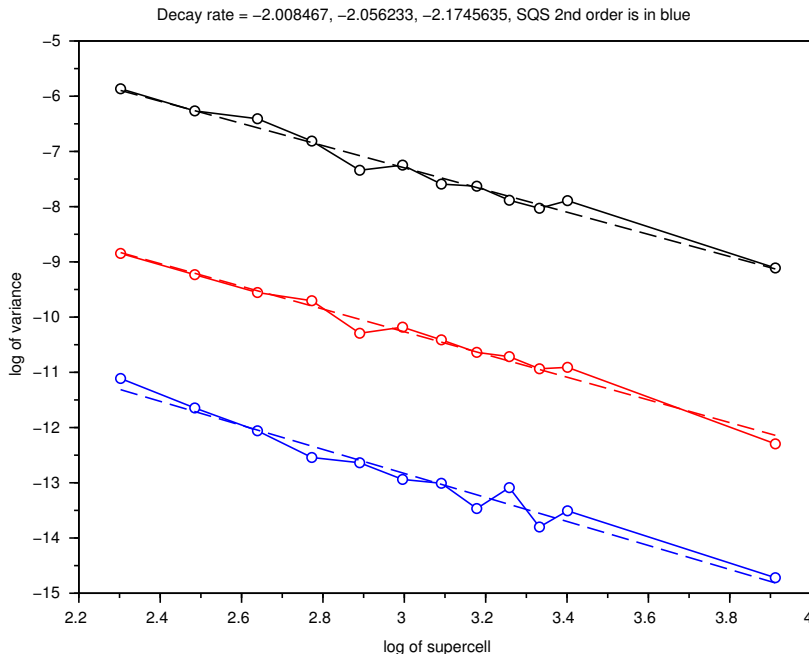


Figure 4.9 – log-log plot of the variance as a function of the supercell side length  $N$  (natural logarithm). Black curve: Monte Carlo method. Red curve: SQS 1 method. Blue curve: SQS 2 method.

### Efficiency with respect to the contrast

We eventually investigate how the contrast in the field  $A$  affects the gain in variance. Results are shown on Table 4.1. We observe that the gain decreases when the contrast increases. Note that this is also the case with the antithetic variable and the control variate techniques that we have previously studied (see [BCLBL12a, LM15b, LM15a]).

However, our SQS 2 approach still yields a significant gain of a factor of 10 when the contrast is equal to 20.

### Conclusion

We may summarize the numerical results by saying that, on the set of two-dimensional test-cases we have considered, the variance is significantly reduced by the approach presented in this article.

Although definite conclusions are yet to be obtained for more challenging, possibly three-dimensional, test-cases, and possibly for more elaborate equations, the results obtained are promising. The robustness and versatility of the approach give us hope for its general, efficient applicability.

Contrast	$V_{\text{MC}}$	$V_{\text{exact SQS 1}}$	$V_{\text{SQS 2}}$	$\frac{V_{\text{MC}}}{V_{\text{exact SQS 1}}}$	$\frac{V_{\text{MC}}}{V_{\text{SQS 2}}}$
1.22	0.0000273	5.801e-08	6.858e-10	470	39821
1.50	0.0001097	0.0000009	1.585e-08	118	6921
1.86	0.0002488	0.0000047	0.0000001	52.6	1996
2.33	0.0004478	0.0000151	0.0000006	29.5	720
3.00	0.0007118	0.0000379	0.0000024	18.8	296
4.00	0.0010496	0.0000814	0.0000080	12.8	131
5.67	0.0014769	0.0001600	0.0000244	9.23	60.5
9.00	0.0020289	0.0003021	0.0000739	6.71	27.4
19.0	0.0028330	0.0006061	0.0002554	4.67	11.1

Table 4.1 – For different values of the contrast, we show the Monte Carlo variance (column #2), the variance of the SQS 1 method (column #3) and the variance of the SQS 2 method (column #4). We next show the variance ratio  $V_{\text{MC}}/V_{\text{exact SQS 1}}$  for the SQS 1 approach (column #5) and the variance ratio  $V_{\text{MC}}/V_{\text{SQS 2}}$  for the SQS 2 approach (column #6). The supercell size is fixed at  $N = 20$ .

## Acknowledgements

The first two authors would like to thank E. Cancès for introducing them to the SQS approach in the context of solid state physics, pointing out to them references [vPDFN10, WFBZ90, ZWFB90], as well as for several stimulating discussions in the early stages of this work. The authors also thank X. Blanc for his suggestions on a previous version of this article.

This work was partially supported by ONR under Grant N00014-12-1-0383 and by EOARD under Grant FA8655-13-1-3061. This work has benefited from a French government grant managed by ANR within the frame of the national program Investments for the Future ANR-11-LABX-022-01.

## 4.5 Proof of Lemma 4.4

We follow the arguments of the proof of [BCLBL12b, Lemma 3.2].

The existence and uniqueness (up to the addition of a constant) of  $\phi_1$  solution to (4.38) is established in [BCLBL12b, Lemma 3.1]. We next point out that (4.39) admits a unique (up to the addition of a constant) solution in  $H_{\text{per}}^1(Q)$ . It is a simple consequence of the Lax-Milgram lemma.

We now prove that the sum in (4.37) is a convergent series in  $L^2(Q \times \Omega)$ . For this purpose, we compute the norm of the remainder of the series, using the notation  $\overline{X}_k(\omega) =$

$X_k(\omega) - \mathbb{E}[X_k]$ :

$$\begin{aligned}
& \left\| \sum_{|k| \geq N+1} \bar{X}_k \nabla \phi_1(\cdot - k) \right\|_{L^2(Q \times \Omega)}^2 \\
&= \sum_{|k| \geq N+1} \sum_{|\ell| \geq N+1} \mathbb{E} [\bar{X}_k \bar{X}_\ell] \int_Q \nabla \phi_1(y - k) \cdot \nabla \phi_1(y - \ell) dy \\
&\leq \sum_{|k| \geq N+1} \sum_{|\ell| \geq N+1} |\text{Cov}(X_k, X_\ell)| \|\nabla \phi_1\|_{L^2(Q-k)} \|\nabla \phi_1\|_{L^2(Q-\ell)} \\
&\leq \sum_{|k| \geq N+1} \sum_{|\ell| \geq N+1} |\text{Cov}(X_k, X_\ell)| \|\nabla \phi_1\|_{L^2(Q-k)}^2,
\end{aligned}$$

where we have used at the last line the discrete Cauchy-Schwarz inequality between  $|\text{Cov}(X_k, X_\ell)|^{1/2} \|\nabla \phi_1\|_{L^2(Q-k)}$  and  $|\text{Cov}(X_k, X_\ell)|^{1/2} \|\nabla \phi_1\|_{L^2(Q-\ell)}$ . We next write, using the stationarity of  $X_k$  and (4.36), that

$$\begin{aligned}
& \left\| \sum_{|k| \geq N+1} \bar{X}_k \nabla \phi_1(\cdot - k) \right\|_{L^2(Q \times \Omega)}^2 \\
&\leq \sum_{|k| \geq N+1} \|\nabla \phi_1\|_{L^2(Q-k)}^2 \sum_{|\ell| \geq N+1} |\text{Cov}(X_k, X_\ell)| \\
&\leq \mathcal{C} \sum_{|k| \geq N+1} \|\nabla \phi_1\|_{L^2(Q-k)}^2.
\end{aligned}$$

The above right-hand side converges to 0 as  $N \rightarrow \infty$  since  $\nabla \phi_1 \in (L^2(\mathbb{R}^d))^d$ .

Hence, the right-hand side of (4.37) defines a function  $T \in (L^2(Q \times \Omega))^d$ . As  $\partial_i T_j = \partial_j T_i$ , there exists a function  $\tilde{u}_1$  such that

$$\nabla \tilde{u}_1 = T = \mathbb{E}[X_0] \nabla \bar{u}_1 + \sum_{k \in \mathbb{Z}^d} (X_k(\omega) - \mathbb{E}[X_k]) \nabla \phi_1(\cdot - k).$$

As  $\bar{u}_1$  is  $\mathbb{Z}^d$ -periodic, we infer from the above equality that

$$\nabla \tilde{u}_1 \text{ is stationary and } \int_Q \mathbb{E}(\nabla \tilde{u}_1) = 0. \quad (4.73)$$

Next, we compute

$$C_0 \nabla \tilde{u}_1 = \mathbb{E}[X_0] C_0 \nabla \bar{u}_1 + \sum_{k \in \mathbb{Z}^d} (X_k(\omega) - \mathbb{E}[X_k]) C_0 \nabla \phi_1(\cdot - k).$$

Taking the divergence of this equation and using (4.32) and (4.34), we thus find that, in the distribution sense,

$$\begin{aligned}
-\text{div} [C_0 \nabla \tilde{u}_1] &= \sum_{k \in \mathbb{Z}^d} - (X_k(\omega) - \mathbb{E}[X_k]) \text{div} [C_0 \nabla \phi_1(\cdot - k)] \\
&\quad - \mathbb{E}[X_0] \text{div} [C_0 \nabla \bar{u}_1] \\
&= \sum_{k \in \mathbb{Z}^d} (X_k(\omega) - \mathbb{E}[X_k]) \text{div} [\mathbb{1}_{Q+k} C_1 p] \\
&\quad + \mathbb{E}[X_0] \text{div} [C_1 p] \\
&= \text{div} [\chi(\cdot, \omega) C_1 p]. \quad (4.74)
\end{aligned}$$

---

Collecting (4.73) and (4.74), we see that  $\tilde{u}_1$  solves (4.19). As the solution to this equation is unique up to the addition of a (possibly random) constant  $C(\omega)$ , we obtain that  $\tilde{u}_1 = u_1 + C(\omega)$ , hence proving (4.37).





## Chapter 5

# A parameter identification problem in stochastic homogenization

Ce **Chapitre** reprend l'intégralité d'un article écrit en collaboration avec Frédéric Legoll, Marielle Simon et Amaël Obliger et accepté dans ESAIM ProcS [[LMOS15](#)].

Dans un cadre d'homogénéisation stochastique discrète, nous cherchons à identifier des paramètres de la loi de probabilité du milieu microscopique, à l'aide de quantités macroscopiques. Nous formulons ce problème à la manière d'un problème d'optimisation aux moindres carrés, et nous démontrons en dimension 1 d'espace qu'il est asymptotiquement bien posé. Ce résultat théorique est complété par des expérimentations numériques.

## A parameter identification problem in stochastic homogenization

Frédéric Legoll<sup>1,2</sup>, William Minvielle<sup>3,2</sup>, Amaël Obliger<sup>4,5</sup> and Marielle Simon<sup>6</sup>

legoll@lami.enpc.fr, minvielw@cermics.enpc.fr, amael.obliger@upmc.fr, marielle.simon@ens-lyon.fr

<sup>1</sup> Laboratoire Navier, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal,  
Pascal,  
77455 Marne-La-Vallée Cedex 2, France ;

<sup>2</sup> INRIA Rocquencourt, MATHERIALS research-team, Domaine de Voluceau, B.P. 105, 78153 Le  
Chesnay Cedex, France ;

<sup>3</sup> CERMICS, École Nationale des Ponts et Chaussées, Université Paris-Est, 6 et 8 avenue Blaise Pascal,  
77455 Marne-La-Vallée Cedex 2, France ;

<sup>4</sup> Sorbonne Universités, UPMC Univ. Paris 06, UMR 8234 PHENIX, 75005 Paris, France ;

<sup>5</sup> ANDRA, Parc de la Croix-Blanche, 1-7, rue Jean-Monnet, 92298 Châtenay-Malabry, France

<sup>6</sup> UMPA, UMR-CNRS 5669, ENS Lyon, 46 allée d'Italie, 69007 Lyon, France.

**Abstract.** *In porous media physics, calibrating model parameters through experiments is a challenge. This process is plagued with errors that come from modelling, measurement and computation of the macroscopic observables through random homogenization – the forward problem – as well as errors coming from the parameters fitting procedure – the inverse problem. In this work, we address these issues by considering a least-square formulation to identify parameters of the microscopic model on the basis on macroscopic observables, including homogenized coefficients. In particular, we discuss the selection of the macroscopic observables which we need to know in order to uniquely determine these parameters. To gain a better intuition and explore the problem without a too high computational load, we mostly focus on the one-dimensional case. We show that the Newton algorithm can be efficiently used to robustly determine optimal parameters, even if some small statistical noise is present in the system.*

**Résumé.** *En physique des milieux poreux, calibrer certains paramètres d'un modèle microscopique sur la base d'expériences donnant accès à des grandeurs macroscopiques est un enjeu majeur. Cette démarche est entachée d'erreurs de modèle, de mesure et de calculs dans la procédure d'homogénéisation: le problème direct est biaisé. La résolution du problème inverse, lorsqu'il s'agit d'estimer les paramètres à partir des observations, engendre aussi des erreurs. Nous considérons ici une formulation "moindres carrés" du problème, cherchant à minimiser l'erreur entre les quantités macroscopiques observées et celles calculées via l'homogénéisation aléatoire. Nous discutons en particulier de la nature des informations macroscopiques nécessaires pour déterminer de manière univoque les paramètres de la densité de probabilité des propriétés microscopiques. Afin d'explorer plus facilement cette question, nous nous intéressons ici essentiellement au cas unidimensionnel. Nous montrons que le problème peut être résolu de manière efficace par l'algorithme de Newton, même en présence d'un petit bruit statistique.*

## 5.1 Introduction

Modelling porous media is a challenge, in particular because the geometry of such materials can be extremely complex. Rock samples are often described as a pile of layers of solid phase which do not permit flows, creating voids in-between layers that are connected by channels, the size and shape of which is difficult to describe (and to observe experimentally, although, in rare cases, imaging methods such as micro-tomography can be used). Besides these issues related to the description of the geometry of the media, another difficulty is to properly model the physical phenomena occurring in the flow. To circumvent these difficulties, a possible approach consists in completely forgetting the exact geometry of the system except for a few parameters (e.g. the size of the channels), and consider that the channels form a simple network, often taken to be  $\mathbb{Z}^d$ . This results in the so-called pore-network models (PNM), initially introduced by Fatt in the 1950s [Fat56] and which have been widely used since then. The void space of a rock (its porosity) is described by a pore network connected by channels. In this framework, the geometry of pores and channels is idealized. Some microscopic properties are assigned to network elements (e.g. the conductance of the channels) and rules are defined to compute the upscaled (homogenized) properties on the basis of this microscopic description. In turn, these upscaled properties can be compared to the available experimental data. The aim is to construct a microscopic network with the same effective properties as those of a real representative sample of rock.

In this work, we follow this approach, and assume that pores are located at the vertices of a simple lattice. Physical properties are described by some random field at the microscopic scale. We focus on monophasic transport phenomena in porous media, where the sample of rock is mainly characterized by its permeability. These phenomena are described by the Darcy law, where the local flux of water is assumed to be proportional to the local pressure gradient, and the microscopic properties of interest are the conductances of the channels. In the pore network model, conductances are solely assigned to channels, and it is assumed that pores do not contribute to the flow. Following Darcy's equation, the microscopic pressure field is computed in the network by ensuring mass conservation at each pore. The equation to solve is therefore a *discrete* linear elliptic equation in divergence form, with random coefficients (see (5.13)–(5.14) below for a more detailed physical description, and (5.4) for a more mathematical description).

The conductances of the channels (i.e. their microscopic permeabilities) depend on their size. Therefore the construction of the network starts by randomly attributing a size to each channel. In practice, this channel size distribution can be inferred from experiments such as mercury porosimetry: we denote it by  $\mathcal{L}_{\text{exp}}$ . Several issues of different nature arise in this procedure. As a consequence, it turns out that the effective properties (e.g. macroscopic permeability) that are computed for a pore network with channel sizes distributed according to  $\mathcal{L}_{\text{exp}}$  are different from the experimental effective properties. The extraction procedure, which provides a channel size distribution, is thus somewhat slightly inconsistent. The main goal of this work consists in improving that distribution, when starting from the experimental initial guess, in order to eventually achieve a better agreement between measured and computed effective properties.

From a more mathematical standpoint, the question can be phrased in the following terms. Consider a second-order divergence-form operator whose coefficients are random. If the distribution of the coefficients is stationary and ergodic, then (under some additional technical assumptions) this random operator can be replaced, over large scales, by an effective operator with constant homogenized coefficients (see Theorem 5.5 below). Random homogenization theory actually provides formulas to compute the homogenized quanti-

ties. We thus have at our disposal a procedure to compute macroscopic quantities if we know the microscopic quantities, and to solve the so-called *forward problem*. However, in practice, given a heterogeneous materials, it is a difficult question to decide on the law of the microscopic physical properties. On the other hand, macroscopic quantities are more easily accessible. It is thus of interest to consider the *inverse problem*, and try to extract some information on the properties of the materials at the microscopic scale on the basis of macroscopic quantities.

In the same spirit, if one makes assumptions on the microscopic law, then macroscopic quantities can be computed, and for instance compared to experimental values. In view of the possible discrepancy between the two, one could question or revisit the assumptions made at the microscopic scale.

Of course, homogenization is an averaging process, which filters out many features of the microscopic coefficients. There is thus no hope to recover a full information about the microstructure (in our case, the *probability distribution* of the conductances) from the only knowledge of macroscopic quantities. We adopt here a more restricted objective. We assume a functional form for the distribution of the microscopic conductances (namely, a Weibull distribution). Our aim is to recover the *parameters* (denoted here  $\theta$ ) of that microscopic law of the basis of macroscopic quantities.

We point out that our approach is not specific to Weibull laws, and that it could be used for other distribution laws with parameters  $\theta$ . What we need is that the random field  $A(x, \omega)$  used at the microscopic scale can be written as

$$A(x, \omega) = \mathcal{F}(u(x, \omega), \theta)$$

where  $u(x, \omega)$  is a field of random variables that are uniformly distributed and  $\mathcal{F}$  smoothly depends on the parameters  $\theta$  (see (5.16) in our particular case). Computing the derivatives of the microscopic random field  $A(x, \omega)$  (and next of the macroscopic, homogenized quantities) with respect to  $\theta$  is then easy. Our motivation for choosing Weibull laws comes from physical reasons: based on experimental results, it appears to be a reasonable choice.

Likewise, our approach is not specific to *discrete* elliptic equations. It could also be applied for problems modelled by *continuous* elliptic partial differential equations (PDEs) with random, highly oscillatory coefficients. Here, we consider discrete equations because the pore network model, which is naturally written in terms of discrete equations, is commonly used for such materials.

The question of recovering the unknown parameters  $\theta$  of the microscopic distribution from homogenized (and more generally macroscopic) quantities belongs to the wide family of *inverse problems*. In this work, a major point of interest is the selection of the macroscopic quantities which we need to know in order to uniquely determine the parameters  $\theta$ . This point is discussed in Section 5.3.2. We refer to [NPS12] for a review article on inverse problems in a multiscale context.

The article is organised as follows. In Section 5.2, we recall some elements of stochastic homogenization and describe the physical problem that motivates this work (including the choice of Weibull laws). We conclude that section with results specific to the one-dimensional case. In particular, random variables distributed according to a Weibull law are not isolated from 0 or  $+\infty$ , and thus the microstructure does not satisfy the classical assumption of ellipticity, namely (5.3) below. We show in Section 5.2.4 that, in the one-dimensional case, homogenization still holds under a weaker assumption, that in turn is satisfied by Weibull random variables.

Next, in Section 5.3, we introduce our parameter fitting problem, formulated as least-square optimization. We first consider the general (multi-dimensional) case before turning to the one-dimensional case. In that latter case, we discuss the macroscopic quantities that are needed to uniquely determine the parameters  $\theta$ . More precisely, Weibull laws have two parameters, and the knowledge of a single homogenized quantity (namely the macroscopic permeability) is, as expected, insufficient to determine the two unknown parameters. We show there (in the one-dimensional case) that, if we additionally specify the relative variance of the effective macroscopic permeability, then we are in position to uniquely determine the two parameters of the microscopic Weibull law.

Section 5.4 is dedicated to numerical results, again in the one-dimensional case. We show that the Newton algorithm can be efficiently used in the current least-square optimization setting. In particular, in practice, the exact homogenized coefficients cannot be computed, and only a random approximation of them is available. We monitor here how this randomness propagates to the optimal parameters. The extension of this work to the two-dimensional case will be addressed in a future work [LMOS].

## 5.2 Discrete homogenization theory

For the sake of completeness, we recall first, in Sections 5.2.1 and 5.2.2, some elements of homogenization for discrete elliptic equations with random coefficients. We refer to [Kün83, Koz87] for seminal contributions on this topic. For homogenization of elliptic partial differential equations (PDEs), we refer to the seminal work [PV81], to the textbooks [BLP78, CD99, JKO94], to [ES08] for a general, numerically oriented presentation, and to the review article [ACLB<sup>+</sup>12].

Next, in Section 5.2.3, we describe the physical background that motivates this work. We eventually turn in Section 5.2.4 to the one-dimensional case, where explicit formulae can be obtained.

### 5.2.1 Homogenization result

We first recall some definitions useful for stochastic homogenization, before turning to the specific case of discrete elliptic equations.

Throughout this article,  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space and we denote by  $\mathbb{E}(X) = \int_{\Omega} X(\omega) d\mathbb{P}(\omega)$  the expectation value of any random variable  $X \in L^1(\Omega, d\mathbb{P})$ . We next fix  $d \in \mathbb{N}^*$  (the ambient physical dimension), and assume that the group  $(\mathbb{Z}^d, +)$  acts on  $\Omega$ . We denote by  $(\tau_k)_{k \in \mathbb{Z}^d}$  this action, and assume that it preserves the measure  $\mathbb{P}$ , that is, for all  $k \in \mathbb{Z}^d$  and all  $B \in \mathcal{F}$ ,  $\mathbb{P}(\tau_k B) = \mathbb{P}(B)$ . We assume that the action  $\tau$  is *ergodic*, that is, if  $B \in \mathcal{F}$  is such that  $\tau_k B = B$  for any  $k \in \mathbb{Z}^d$ , then  $\mathbb{P}(B) = 0$  or 1. In addition, we introduce the following notion of stationarity:

**Definition 5.1.** *We say that a function  $\psi : \mathbb{Z}^d \times \Omega \rightarrow \mathbb{R}$  is stationary if*

$$\forall x, z \in \mathbb{Z}^d, \quad \psi(x + z, \omega) = \psi(x, \tau_z \omega) \quad a.s. \quad (5.1)$$

We now focus on the case of discrete elliptic equations. We view  $\mathbb{Z}^d$  as a lattice, whose unit vectors are denoted by  $e_i$ ,  $i \in \{1, \dots, d\}$ . Each vertex  $x \in \mathbb{Z}^d$  of the lattice is connected to  $2d$  other vertices:  $x \pm e_i$ ,  $i \in \{1, \dots, d\}$ . We write  $x \sim y$  if  $x$  and  $y$  are neighbours (i.e. connected), and  $e = (x, y)$  the corresponding (non-oriented) edge. For any vertex  $x \in \mathbb{Z}^d$  and any direction  $1 \leq i \leq d$ , we denote by  $a_i(x, \omega) \in (0, \infty)$  the random conductance of the

edge  $(x, x + e_i)$ . We next introduce the diagonal matrix  $A$  defined for any vertex  $x \in \mathbb{Z}^d$  by

$$A(x, \omega) = \text{diag}\left(a_1(x, \omega), \dots, a_d(x, \omega)\right). \quad (5.2)$$

We assume that, for any direction  $i$ , the conductances  $\{a_i(x, \cdot)\}_{x \in \mathbb{Z}^d}$  form an i.i.d. sequence of random variables. The matrix  $A$  is therefore stationary.

We introduce the following assumption:

**ASSUMPTION 5.2** (Ellipticity – boundedness condition). *There exist two positive deterministic constants  $c$  and  $C$  such that the matrix  $A$  defined by (5.2) satisfies*

$$\forall \xi \in \mathbb{R}^d, \quad \forall x \in \mathbb{Z}^d, \quad c|\xi|^2 \leq \xi \cdot A(x, \omega)\xi \leq C|\xi|^2 \quad a.s. \quad (5.3)$$

In view of (5.2), note that this simply means that  $0 < c \leq a_j(x, \omega) \leq C$  almost surely, for any  $1 \leq j \leq d$  and any  $x \in \mathbb{Z}^d$ .

We next introduce discrete differential operators on the lattice  $\mathbb{Z}^d$ .

**Definition 5.3.** *For a function  $g : \mathbb{Z}^d \rightarrow \mathbb{R}$ , the gradient  $\nabla g : \mathbb{Z}^d \rightarrow \mathbb{R}^d$  is defined by*

$$(\nabla g)(x) = \begin{pmatrix} g(x + e_1) - g(x) \\ \vdots \\ g(x + e_d) - g(x) \end{pmatrix}.$$

*For a function  $G = (G_1, \dots, G_d) : \mathbb{Z}^d \rightarrow \mathbb{R}^d$ , the function  $\nabla^* G : \mathbb{Z}^d \rightarrow \mathbb{R}$  is defined by*

$$-(\nabla^* G)(x) = \sum_{i=1}^d (G_i(x) - G_i(x - e_i)).$$

We think of  $\nabla^* G$  as the negative divergence of  $G$ . The operator  $\nabla^*$  is the  $\ell^2$  transpose of  $\nabla$  in the following sense: for any compactly supported functions  $g : \mathbb{Z}^d \rightarrow \mathbb{R}$  and  $G : \mathbb{Z}^d \rightarrow \mathbb{R}^d$ ,

$$\sum_{x \in \mathbb{Z}^d} g(x) \nabla^* G(x) = \sum_{x \in \mathbb{Z}^d} \nabla g(x) \cdot G(x).$$

Hereafter, the notation  $a \cdot b$  stands for the usual scalar product in  $\mathbb{R}^d$ .

We additionally define rescaled discrete differential operators as follows:

**Definition 5.4.** *For a function  $g : \varepsilon\mathbb{Z}^d \rightarrow \mathbb{R}$ , the gradient  $\nabla_\varepsilon g : \varepsilon\mathbb{Z}^d \rightarrow \mathbb{R}^d$  is defined by*

$$(\nabla_\varepsilon g)(x) = \frac{1}{\varepsilon} \begin{pmatrix} g(x + \varepsilon e_1) - g(x) \\ \vdots \\ g(x + \varepsilon e_d) - g(x) \end{pmatrix}.$$

*For a function  $G = (G_1, \dots, G_d) : \varepsilon\mathbb{Z}^d \rightarrow \mathbb{R}^d$ , the function  $\nabla_\varepsilon^* G : \varepsilon\mathbb{Z}^d \rightarrow \mathbb{R}$  is defined by*

$$-(\nabla_\varepsilon^* G)(x) = \sum_{i=1}^d \frac{G_i(x) - G_i(x - \varepsilon e_i)}{\varepsilon}.$$

The following homogenization result holds (we refer to [Kün83, Theorems 3 and 4] for a proof):

**Theorem 5.5.** *Let  $\mathcal{D}$  be a bounded domain of  $\mathbb{R}^d$  and  $f \in C^0(\overline{\mathcal{D}})$ . Let  $A$  be the random stationary matrix field given by (5.2). We assume that (5.3) holds. Let  $u_\varepsilon \in \ell^2(\varepsilon\mathbb{Z}^d; \mathbb{R})$  be the unique solution to*

$$\nabla_\varepsilon^* [A(x/\varepsilon, \omega) \nabla_\varepsilon u_\varepsilon(x, \omega)] = f(x) \quad \text{in } \mathcal{D} \cap \varepsilon\mathbb{Z}^d, \quad u_\varepsilon(x, \omega) = 0 \quad \text{in } (\mathbb{R}^d \setminus \mathcal{D}) \cap \varepsilon\mathbb{Z}^d. \quad (5.4)$$

When  $\varepsilon$  goes to 0,  $u_\varepsilon(\cdot, \omega)$  converges to some homogenized function  $u^*$ .

For any  $\xi \in \mathbb{R}^d$ , introduce the corrector  $\varphi_\xi$  in the direction  $\xi$  as the unique solution (defined on  $\mathbb{Z}^d \times \Omega$ ) to

$$\begin{cases} -\nabla^* [A(\cdot, \omega)(\xi + \nabla\varphi_\xi(\cdot, \omega))] = 0 \quad \text{in } \mathbb{Z}^d, \text{ a.s.}, \\ \nabla\varphi_\xi \text{ is stationary in the sense of (5.1),} \\ \forall x \in \mathbb{Z}^d, \quad \mathbb{E}[\nabla\varphi_\xi(x, \cdot)] = 0, \\ \varphi_\xi(0, \omega) = 0 \text{ a.s.} \end{cases} \quad (5.5)$$

Introduce next the constant matrix  $A^*$  defined by

$$\forall \xi \in \mathbb{R}^d, \quad A^*\xi = \mathbb{E}[A(x, \cdot)(\xi + \nabla\varphi_\xi(x, \cdot))] \quad (5.6)$$

and the unique solution  $u^* \in H_0^1(\mathcal{D})$  to the (continuous) partial differential equation

$$-\operatorname{div}[A^*\widehat{\nabla}u^*] = f \quad \text{in } \mathcal{D},$$

where  $\widehat{\nabla}$  and  $\operatorname{div}$  are the usual (continuous) gradient and divergence differential operators.

Then, we have the (strong) convergence  $u_\varepsilon \xrightarrow[\varepsilon \rightarrow 0]{} u^*$ , in the sense that

$$\varepsilon^d \sum_{x \in \mathcal{D} \cap \varepsilon\mathbb{Z}^d} |u_\varepsilon(x, \omega) - u^*(x)|^2 \xrightarrow[\varepsilon \rightarrow 0]{} 0 \quad \text{almost surely.} \quad (5.7)$$

Note that, in the right-hand side of (5.6), the vector  $A(\xi + \nabla\varphi_\xi)$  is stationary, and therefore the expectation may be evaluated at any  $x \in \mathbb{Z}^d$ . Note also that, in general,  $\varphi_\xi$  itself is not stationary, as the one-dimensional case shows. Only its gradient is.

**Remark 5.6.** *We can define, on  $\mathcal{D}$ , the function*

$$\tilde{u}_\varepsilon(x, \omega) = \sum_{k \in \varepsilon\mathbb{Z}^d \cap \mathcal{D}} u_\varepsilon(k, \omega) \mathbb{1}_{k+\varepsilon Q}(x), \quad \text{where } Q = (0, 1)^d.$$

Then (5.7) implies that  $\tilde{u}_\varepsilon(\cdot, \omega) \xrightarrow[\varepsilon \rightarrow 0]{} u^*$  in  $L^2(\mathcal{D})$  almost surely.

### 5.2.2 Approximation on finite boxes

The corrector problem (5.5) is untractable in practice, since it is posed in the entire lattice  $\mathbb{Z}^d$ . Approximations are therefore in order. The standard procedure amounts to considering finite boxes (see e.g. [BP04]). For a positive integer  $N$ , we denote by  $\Lambda_N$  the finite box  $\{0, \dots, N\}^d$  and by  $\mathcal{E}_N$  the set of edges in  $\Lambda_N$  (see Figure 5.1).

The truncated corrector  $\varphi_\xi^N$  defined on  $\Lambda_N \times \Omega$  is the unique solution to

$$\begin{cases} -\nabla^* [A(\cdot, \omega)(\xi + \nabla\varphi_\xi^N(\cdot, \omega))] = 0 \quad \text{in } \Lambda_N, \text{ a.s.}, \\ \varphi_\xi^N(\cdot, \omega) \text{ is } \Lambda_N\text{-periodic,} \\ \varphi_\xi^N(0, \omega) = 0 \text{ a.s.} \end{cases} \quad (5.8)$$



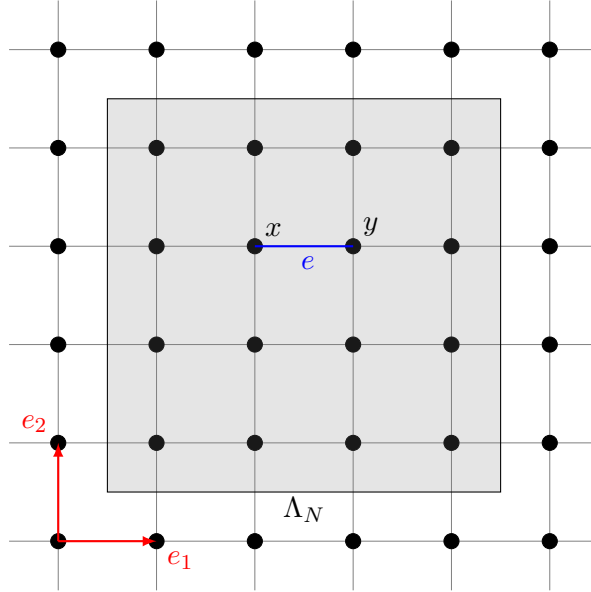


Figure 5.1 – Finite box  $\Lambda_N$  in  $\mathbb{Z}^2$

The homogenized matrix  $A^*$ , which is deterministic, is then approximated by the matrix  $A_N^*$  defined by

$$\forall \xi \in \mathbb{R}^d, \quad A_N^*(\omega)\xi = \frac{1}{|\Lambda_N|} \sum_{x \in \Lambda_N} A(x, \omega)(\xi + \nabla \varphi_\xi^N(x, \omega)). \quad (5.9)$$

Because of truncation, the practical approximation  $A_N^*$  is random. In the large  $N$  limit, the deterministic value is attained, thanks to ergodicity. More precisely,  $A_N^*(\omega)$  converges almost surely towards  $A^*$  as  $N$  goes to infinity, thanks to the ergodic theorem.

**Remark 5.7.** In (5.8), we have complemented the elliptic equation in  $\Lambda_N$  with periodic boundary conditions. Other choices could be made, such as imposing homogeneous Dirichlet boundary conditions:  $\varphi_\xi^N(\cdot, \omega) = 0$  on  $\partial\Lambda_N$  (see e.g. [BP04] for a similar discussion in the case of continuous PDEs). In the numerical experiments of Section 5.4, we only use periodic boundary conditions, following (5.8).

In practice, we work on a finite box  $\Lambda_N$ , on which the apparent homogenized matrix  $A_N^*$  is random. It is therefore natural to introduce  $M$  i.i.d. realizations of the random field  $A(x, \omega)$  and solve (5.8)–(5.9) for each of them, thereby obtaining i.i.d. realizations  $A_N^{*,m}(\omega)$ ,  $1 \leq m \leq M$ . We next introduce the empirical mean

$$\overline{A}_{N,M}^*(\omega) = \frac{1}{M} \sum_{m=1}^M A_N^{*,m}(\omega) \quad (5.10)$$

which is, according to the Central Limit Theorem, a converging approximation of  $\mathbb{E}[A_N^*]$ . We have that

$$\overline{A}_{N,M}^*(\omega) \xrightarrow{M \rightarrow \infty} \mathbb{E}[A_N^*] \quad \text{a.s.}$$

In addition, for any entry  $1 \leq i, j \leq d$  of the matrix, we have that, with a probability of 95 %,

$$\left| \left( \overline{A}_{N,M}^*(\omega) \right)_{ij} - \mathbb{E} \left[ (A_N^*)_{ij} \right] \right| \leq 1.96 \sqrt{\frac{\text{Var} (A_N^*)_{ij}}{M}}.$$

The error when approximating  $A^*$  by  $\overline{A}_{N,M}^*$  can be written as the sum of two contributions,

$$A^* - \overline{A}_{N,M}^* = \left( A^* - \mathbb{E}[A_N^*] \right) + \left( \mathbb{E}[A_N^*] - \overline{A}_{N,M}^* \right). \quad (5.11)$$

The second term in the right-hand side of (5.11) is the *statistical* error. The first term is the *systematic* error, due to the fact that, for any finite  $N$ ,  $\mathbb{E}[A_N^*] \neq A^*$ . The dominated convergence theorem ensures that this error vanishes as  $N \rightarrow \infty$ . Many studies have been recently devoted to proving sharp estimates on the rate of this convergence, following the seminal works [Yur86, BP04]. In [GNO14, Theorem 2], the authors show (for any dimension  $d \geq 2$ ) that the systematic error is of order  $N^{-d} \ln^d(N)$  when using periodic boundary conditions (namely, solving (5.8)), and that  $\text{Var}(A_N^*)$  scales as  $N^{-d}$ . Likewise, in the case of *continuous* PDEs, when  $d \geq 3$ , optimal estimates on  $\text{Var}(A_N^*)$  have been established in [Nol14, Theorem 1.3 and Proposition 1.4].

**Remark 5.8.** *The estimator  $\overline{A}_{N,M}^*$  only agrees with  $\mathbb{E}[A_N^*]$  in the limit of an asymptotically large number  $M$  of realizations. Note that variance reduction approaches have been introduced in this context (see e.g. [BCLBL12a, BCLBL12b, CLBL10], and [LM15b] for the extension to a nonlinear setting) to obtain approximations of  $\mathbb{E}[A_N^*]$  in a more efficient manner than by using  $\overline{A}_{N,M}^*$ .*

In the sequel, we will identify the parameters of the microscopic probability distribution on the basis of two types of macroscopic quantities:

1. the homogenized permeability, which is in practice approximated by  $\overline{A}_{N,M}^*$ ;
2. the relative variance of any entry  $(A_N^*(\omega))_{ij}$ , defined by

$$\text{VarR} \left[ (A_N^*)_{ij} \right] := \frac{\text{Var} \left[ (A_N^*)_{ij} \right]}{\left( \mathbb{E} \left[ (A_N^*)_{ij} \right] \right)^2},$$

which is in practice approximated by

$$S_{N,M} = \frac{1}{(\overline{A}_{N,M}^*)_{ij}^2} \left( \frac{1}{M} \sum_{m=1}^M \left( (A_N^{*,m}(\omega))_{ij} - (\overline{A}_{N,M}^*)_{ij} \right)^2 \right). \quad (5.12)$$

### 5.2.3 Physical problem

We describe here the physical background which inspires this work. As pointed out above, from a physical viewpoint, understanding the microscopic properties of charged porous media is of great importance. Such materials have elaborate geometries that make direct computations very challenging. To circumvent this issue, we use here the Pore Network Model (PNM), which involves a simplified model of the geometry. In the PNM model, pores are located at the vertices of the lattice  $\mathbb{Z}^d$ . Neighbouring pores are connected by channels, which allow water to flow. Each channel  $(x, x + e_i)$  is endowed with its random conductance  $a_i(x, \omega) > 0$ , the probability distribution of which is discussed below.

Experiments provide measures on the *macroscopic permeability*  $K_{\text{obs}}$ , which is modelled as a homogenized coefficient  $K^*$ . In practice, as explained in Section 5.2.2, the homogenized coefficient can only be approximated through a computation on a large box. Assuming that the conductance field  $a_i(x, \omega)$  is given for any direction  $1 \leq i \leq d$  and any vertex  $x$  on the finite lattice  $\Lambda_N$ , the PNM model consists in computing the pressure field  $P(x, \omega)$

by solving the conservation equations (i.e., Darcy law) in the network. This leads to the following linear system:

$$\forall x \in \Lambda_N, \quad \sum_{y \sim x} \tilde{a}(x, y, \omega) (P(y, \omega) - P(x, \omega)) = 0, \quad (5.13)$$

where  $\tilde{a}(x, y, \omega)$  is the conductance of the non-oriented edge  $(x, y)$ . Some boundary conditions need to be imposed to make this problem well-posed, they are discussed below. We next see, by definition of  $a_i$ , that

$$\begin{aligned} & \sum_{y \sim x} \tilde{a}(x, y, \omega) (P(y, \omega) - P(x, \omega)) \\ = & \sum_{i=1}^d a_i(x, \omega) (P(x + e_i, \omega) - P(x, \omega)) + \sum_{i=1}^d a_i(x - e_i, \omega) (P(x - e_i, \omega) - P(x, \omega)) \\ = & \nabla^* [A(\cdot, \omega) \nabla P(\cdot, \omega)](x), \end{aligned} \quad (5.14)$$

where the matrix  $A$  is defined in terms of  $\{a_i\}_{i=1}^d$  by (5.2).

We now describe (in the two-dimensional case, for the sake of simplicity) the boundary conditions imposed on (5.13). They are designed to mimic experimental conditions. We first recall that the large box reads  $\Lambda_N = \{0, \dots, N\}^2$ . The pressure field is assumed to be periodic in the vertical direction, whereas a macroscopic gradient is imposed in the horizontal direction as follows. Imagine that all vertices with coordinates  $(0, \cdot)$  are connected to one fixed vertex denoted by  $O$ , representing a pressure reservoir at pressure  $P_O$ . Likewise, all vertices with coordinates  $(N, \cdot)$  are connected to one fixed vertex denoted by  $I$  at pressure  $P_I$  (see Figure 5.2). Then, the boundary conditions write

$$\text{for all } j \in \{0, \dots, N\}, \quad P(0, j) = P_O \quad \text{and} \quad P(N, j) = P_I.$$

Once (5.13) is solved with the above boundary conditions, the *macroscopic permeability*  $K_N^*$  is defined by

$$K_N^*(\omega) := \frac{N}{P_O - P_I} \frac{1}{|\Lambda_N|} \sum_{x \in \Lambda_N} \tilde{a}(x, x + e_1, \omega) (P(x, \omega) - P(x + e_1, \omega)). \quad (5.15)$$

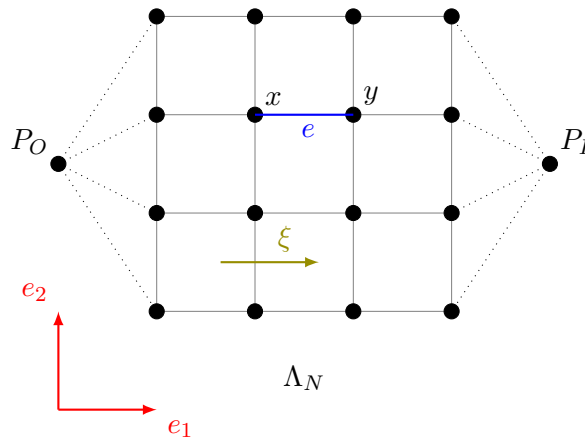


Figure 5.2 – Finite lattice with boundary conditions

Let us now show that Equations (5.13)–(5.15) are actually the same as Equations (5.8)–(5.9) written above. By linearity of (5.13)–(5.15), we can assume that  $P_O = 0$  and  $P_I = N$  without loss of generality. Let  $P$  be a solution to (5.13). We introduce  $\varphi_{e_1}$  such that

$$P(x, \omega) = x \cdot e_1 + \varphi_{e_1}(x, \omega).$$

In view of (5.13) and (5.14), we see that  $\varphi_{e_1}(\cdot, \omega)$  is solution to

$$\forall x \in \Lambda_N, \quad \nabla^* [A(\cdot, \omega)(e_1 + \nabla \varphi_{e_1}(\cdot, \omega))](x) = 0$$

with  $\varphi_{e_1}((0, j), \omega) = \varphi_{e_1}((N, j), \omega) = 0$  for any  $j$  and  $\varphi_{e_1}(\cdot, \omega)$  is periodic in the vertical direction. Up to the choice of boundary conditions, we thus recognize (5.8) for  $\xi = e_1$ . We also infer from (5.15) that

$$\begin{aligned} K_N^*(\omega) &= \frac{N}{P_O - P_I} \frac{1}{|\Lambda_N|} \sum_{x \in \Lambda_N} \tilde{a}(x, x + e_1, \omega) (P(x, \omega) - P(x + e_1, \omega)) \\ &= \frac{-1}{|\Lambda_N|} \sum_{x \in \Lambda_N} a_1(x, \omega) (\varphi_{e_1}(x, \omega) - \varphi_{e_1}(x + e_1, \omega) - e_1 \cdot e_1) \\ &= \frac{1}{|\Lambda_N|} \sum_{x \in \Lambda_N} a_1(x, \omega) e_1^T (e_1 + \nabla \varphi_{e_1}(x, \omega)) \\ &= e_1^T A_N^*(\omega) e_1, \end{aligned}$$

where  $A_N^*(\omega)$  is defined by (5.9), and where we have used (5.2) in the last line. Thus, up to the choice of boundary conditions in the corrector problem, the formulation (5.13)–(5.15) is identical to the formulation (5.8)–(5.9).

We eventually discuss the choice of the probability distribution for the conductances. Based on experimental results, it is reasonable to assume the following:

**ASSUMPTION 5.9.** *We assume that the radius  $r$  of the channels are i.i.d. random variables distributed according to a Weibull law of parameter  $\theta := (\lambda, k) \in (\mathbb{R}_+^*)^2$ , that we denote  $\mathcal{W}(\lambda, k)$ . We recall that such random variables are positive, with a probability density that reads (see Figure 5.3)*

$$\forall r > 0, \quad f(r; k, \lambda) = \frac{k}{\lambda} \left(\frac{r}{\lambda}\right)^{k-1} \exp\left(-\left(\frac{r}{\lambda}\right)^k\right),$$

corresponding to the cumulative distribution function

$$F(r; k, \lambda) = \int_0^r f(s; k, \lambda) ds = 1 - \exp\left(-\left(\frac{r}{\lambda}\right)^k\right).$$

Note that the radius of all channels (independently of their direction  $1 \leq i \leq d$ ) share the same probability distribution.

In practice, a Weibull distribution is generated as follows. Let  $u(\omega)$  be a random variable uniformly distributed in  $[0, 1]$ . Then

$$r(\omega) = \lambda \left[ -\ln(1 - u(\omega)) \right]^{1/k}$$

is distributed according to the Weibull law of parameter  $(\lambda, k)$ .

Physical arguments lead to the fact that the conductance  $a_i(x, \omega)$  of any channel  $(x, x + e_i)$  is directly related to its radius  $r(x, x + e_i, \omega)$ . Hereafter, we assume that

$$a_i(x, \omega) = C_0 r^4(x, x + e_i, \omega) = C_0 \lambda^4 \left[ -\ln(1 - u_i(x, \omega)) \right]^{4/k}, \quad (5.16)$$

where  $C_0$  is a constant (for instance, for a Poiseuille flow,  $C_0 = \pi/(8\eta)$  where  $\eta$  is the fluid viscosity and  $u_i(x, \omega)$  is uniformly distributed in  $[0, 1]$ ). For the sake of simplicity, we take  $C_0 = 1$  in the sequel. Therefore, we assume that

*The conductances  $\{a_i(x, \omega)\}_{x \in \mathbb{Z}^d, 1 \leq i \leq d}$  form an i.i.d. sequence of random variables that are distributed according to the Weibull law of parameter  $(\lambda^4, k/4)$ .* (5.17)

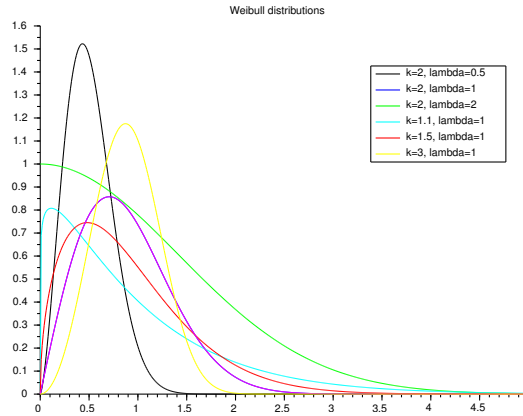


Figure 5.3 – Examples of Weibull distributions.

**Remark 5.10.** *Note that the Weibull distribution is isolated neither from 0 nor from  $\infty$ . The above model therefore does not satisfy the ellipticity condition (5.3). First, we show in Section 5.2.4 below that, in the one-dimensional case, the assumption (5.3) is not necessary, and that homogenization holds under a weaker assumption. Second, we refer to [Bis11] for similar studies (again under assumptions weaker than (5.3)) in higher-dimensional cases.*

**Remark 5.11.** *The numerical tests of Section 5.4 are performed with the above model, and thus aim at identifying the two parameters  $\lambda$  and  $k$ . We however note that nothing in our approach is specific to this particular model using Weibull laws. This choice is only motivated by physical reasons.*

Since the conductances  $a_i(x, \omega)$  are all i.i.d. (for any  $1 \leq i \leq d$  and any  $x \in \mathbb{Z}^d$ ), the problem is invariant by any rotation of angle  $\pi/2$ . The homogenized matrix  $A^*$  is therefore proportional to the identity matrix  $\text{Id}_d$ , and reads

$$A^* = K^* \text{Id}_d$$

where  $K^* \in (0, \infty)$  is the homogenized permeability. We can also write that

$$K^* = e_1 \cdot A^* e_1.$$

In practice, we only have access to  $A_N^*(\omega)$ , which is a symmetric matrix (but is a priori not proportionnal to the identity matrix). All directions are statistically identical, hence we only focus in the sequel on

$$K_N^*(\omega) := e_1 \cdot A_N^*(\omega) e_1. \quad (5.18)$$

### 5.2.4 The one dimensional case

The purpose of this section is two-fold. First, we recall explicit formulas for the homogenized quantities in terms of the microscopic field  $A(x, \omega)$ . We derive these formulas assuming that (5.3) holds. Second, we show that we can relax Assumption (5.3) and still state a homogenization result.

#### Explicit formulas in the elliptic case (5.3)

In the one-dimensional case, the problem (5.8)–(5.9) can be analytically solved. We have

$$A_N^*(\omega) = \left( \frac{1}{N} \sum_{x \in \Lambda_N} \frac{1}{A(x, \omega)} \right)^{-1} \quad \text{for almost all } \omega. \quad (5.19)$$

Likewise, the problem (5.5)–(5.6) can also be solved, yielding the formula

$$A^* = \left( \mathbb{E} \left[ \frac{1}{A(x, \cdot)} \right] \right)^{-1}, \quad (5.20)$$

which can be evaluated at any  $x \in \mathbb{Z}$  due to the stationarity of  $A$ .

We note that, as soon as  $A(x, \omega) > 0$  a.s. for any  $x \in \mathbb{Z}$  and  $A^{-1}(x, \cdot) \in L^1(\Omega)$  (this latter condition being independent of  $x$ ), formulas (5.19) and (5.20) are well-defined. The aim of the next section is to recall that, in the one-dimensional case, these assumptions are enough for homogenization to hold.

#### Relaxing Assumption (5.3)

In this section, we show that the following assumption is enough for homogenization to hold:

**ASSUMPTION 5.12.** *We assume that the coefficient  $A$  is almost surely positive and finite and satisfies*

$$A^{-1}(x, \cdot) \in L^1(\Omega). \quad (5.21)$$

Of course, by stationarity, if (5.21) is satisfied for some  $x \in \mathbb{Z}$ , then it is satisfied for all  $x \in \mathbb{Z}$ .

**Theorem 5.13.** *Let  $\mathcal{D}$  be a bounded domain of  $\mathbb{R}$ ,  $f \in C^0(\overline{\mathcal{D}})$  and  $A$  be a random stationary scalar field (defined on  $\mathbb{Z} \times \Omega$ ) that satisfies (5.21). Let  $u_\varepsilon \in \ell^2(\varepsilon\mathbb{Z}; \mathbb{R})$  be the unique solution to*

$$\nabla_\varepsilon^* [A(x/\varepsilon, \omega) \nabla_\varepsilon u_\varepsilon(x, \omega)] = f(x) \quad \text{in } \mathcal{D} \cap \varepsilon\mathbb{Z}, \quad u_\varepsilon(x, \omega) = 0 \quad \text{in } (\mathbb{R} \setminus \mathcal{D}) \cap \varepsilon\mathbb{Z}, \quad (5.22)$$

and let  $u^* \in H_0^1(\mathcal{D})$  be the unique solution to the (continuous) boundary value problem

$$- [A^*(u^*)]' = f \quad \text{in } \mathcal{D}, \quad (5.23)$$

where  $A^*$  is defined by (5.20).

Then, when  $\varepsilon \rightarrow 0$ ,  $u_\varepsilon(\cdot, \omega)$  converges to the homogenized solution  $u^*$ , in the sense that

$$\varepsilon \sum_{x \in \mathcal{D} \cap \varepsilon\mathbb{Z}} |u_\varepsilon(x, \omega) - u^*(x)|^2 \xrightarrow{\varepsilon \rightarrow 0} 0 \quad \text{almost surely.} \quad (5.24)$$

Note that (5.22) is almost surely well-posed. Indeed, since  $A$  is stationary and  $0 < A(0, \omega) < \infty$  almost surely, we have that, almost surely,  $0 < A(x, \omega) < \infty$  for any  $x \in \mathbb{Z}$ . For those  $\omega$ , problem (5.22) is well-posed. Likewise, since  $A$  is almost surely finite (resp.  $A^{-1}(0, \cdot) \in L^1(\Omega)$ ), we have that  $A^* < \infty$  (resp.  $A^* > 0$ ) and hence (5.23) is well-posed.

*Proof.* The proof proceeds by truncation of the coefficient  $A$  in the neighbourhood of 0 and  $+\infty$ . For the sake of simplicity, we take  $\mathcal{D} = (0, 1)$ . For any  $m \in \mathbb{N}^*$ , we introduce the coefficient  $A_m$  defined on  $\mathbb{Z} \times \Omega$  by

$$A_m(x, \omega) := \begin{cases} \frac{1}{m} & \text{if } 0 < A(x, \omega) < \frac{1}{m}, \\ A(x, \omega) & \text{if } \frac{1}{m} \leq A(x, \omega) \leq m, \\ m & \text{if } A(x, \omega) > m. \end{cases}$$

We set

$$A_m^* = \left( \mathbb{E} \left[ \frac{1}{A_m(0, \cdot)} \right] \right)^{-1}.$$

For almost all  $\omega$  (i.e. those such that  $A(0, \omega) > 0$ ), we have

$$\begin{aligned} \lim_{m \rightarrow \infty} \frac{1}{A_m(0, \omega)} &= \frac{1}{A(0, \omega)}, \\ \forall m \in \mathbb{N}^*, \quad 0 < \frac{1}{A_m(0, \omega)} &\leq 1 + \frac{1}{A(0, \omega)}, \end{aligned}$$

where the right-hand side of the above second line belongs to  $L^1(\Omega)$ , in view of the assumption (5.21). Therefore, the dominated convergence theorem implies that

$$\lim_{m \rightarrow \infty} A_m^* = A^*. \quad (5.25)$$

Let  $u_\varepsilon^m \in \ell^2(\varepsilon\mathbb{Z}; \mathbb{R})$  be the unique solution to

$$\nabla_\varepsilon^* [A_m(x/\varepsilon, \omega) \nabla_\varepsilon u_\varepsilon^m(x, \omega)] = f(x) \quad \text{in } (0, 1) \cap \varepsilon\mathbb{Z}, \quad u_\varepsilon^m(x, \omega) = 0 \quad \text{in } (\mathbb{R} \setminus (0, 1)) \cap \varepsilon\mathbb{Z}, \quad (5.26)$$

and let  $u_m^* \in H_0^1(0, 1)$  be the unique solution to the (continuous) boundary value problem

$$- [A_m^*(u_m^*)]' = f \quad \text{in } (0, 1).$$

We write

$$\|u_\varepsilon(\cdot, \omega) - u^*\|_{\ell_\varepsilon^2} \leq \|u_\varepsilon(\cdot, \omega) - u_\varepsilon^m(\cdot, \omega)\|_{\ell_\varepsilon^2} + \|u_\varepsilon^m(\cdot, \omega) - u_m^*\|_{\ell_\varepsilon^2} + \|u_m^* - u^*\|_{\ell_\varepsilon^2} \quad (5.27)$$

where, for any function  $v$ ,

$$\|v\|_{\ell_\varepsilon^2} := \sqrt{\varepsilon \sum_{x \in (0, 1) \cap \varepsilon\mathbb{Z}} v^2(x)}.$$

We successively study the three terms of the right-hand side of (5.27).

First, we have

$$\lim_{\varepsilon \rightarrow 0} \|u_m^* - u^*\|_{\ell_\varepsilon^2} = \|u_m^* - u^*\|_{L^2(0, 1)},$$

and the convergence (5.25) implies that

$$\lim_{m \rightarrow \infty} \lim_{\varepsilon \rightarrow 0} \|u_m^* - u^*\|_{\ell_\varepsilon^2} = 0. \quad (5.28)$$

Second, the coefficient  $A_m$  satisfies the ellipticity condition (5.3), so we infer from Theorem 5.5 that, for any  $m \in \mathbb{N}^*$ ,

$$\lim_{\varepsilon \rightarrow 0} \|u_\varepsilon^m(\cdot, \omega) - u_m^*\|_{\ell_\varepsilon^2} = 0 \quad \text{a.s.} \quad (5.29)$$

We eventually turn to the first term of the right-hand side of (5.27). Let

$$F_\varepsilon(x) = \varepsilon \sum_{y \in (0, x] \cap \varepsilon \mathbb{Z}} f(y),$$

which satisfies, for any  $x$ ,  $|F_\varepsilon(x)| \leq \|f\|_{L^\infty}$ . Integrating once the equations (5.22) and (5.26), we can show that there exist two random variables  $C_\varepsilon(\omega)$  and  $C_\varepsilon^m(\omega)$ , independent of  $x$ , such that

$$A_m\left(\frac{x}{\varepsilon}, \omega\right) \nabla_\varepsilon u_\varepsilon^m(x, \omega) = -F_\varepsilon(x) + C_\varepsilon^m(\omega), \quad (5.30)$$

$$A\left(\frac{x}{\varepsilon}, \omega\right) \nabla_\varepsilon u_\varepsilon(x, \omega) = -F_\varepsilon(x) + C_\varepsilon(\omega). \quad (5.31)$$

Using the boundary conditions on  $u_\varepsilon^m$  and  $u_\varepsilon$ , we get

$$C_\varepsilon(\omega) = \frac{\mathcal{N}_\varepsilon(\omega)}{\mathcal{D}_\varepsilon(\omega)} \quad \text{and} \quad C_\varepsilon^m(\omega) = \frac{\mathcal{N}_\varepsilon^m(\omega)}{\mathcal{D}_\varepsilon^m(\omega)}$$

where

$$\begin{aligned} \mathcal{D}_\varepsilon(\omega) &= \varepsilon \sum_{x \in (0, 1) \cap \varepsilon \mathbb{Z}} A\left(\frac{x}{\varepsilon}, \omega\right)^{-1}, & \mathcal{N}_\varepsilon(\omega) &= \varepsilon \sum_{x \in (0, 1) \cap \varepsilon \mathbb{Z}} A\left(\frac{x}{\varepsilon}, \omega\right)^{-1} F_\varepsilon(x), \\ \mathcal{D}_\varepsilon^m(\omega) &= \varepsilon \sum_{x \in (0, 1) \cap \varepsilon \mathbb{Z}} A_m\left(\frac{x}{\varepsilon}, \omega\right)^{-1}, & \mathcal{N}_\varepsilon^m(\omega) &= \varepsilon \sum_{x \in (0, 1) \cap \varepsilon \mathbb{Z}} A_m\left(\frac{x}{\varepsilon}, \omega\right)^{-1} F_\varepsilon(x). \end{aligned}$$

All these quantities are well-defined for almost all  $\omega$ . We claim that

$$\lim_{m \rightarrow \infty} \limsup_{\varepsilon \rightarrow 0} |C_\varepsilon(\omega) - C_\varepsilon^m(\omega)| = 0 \quad \text{a.s.} \quad (5.32)$$

To prove this claim, we start by writing that

$$C_\varepsilon(\omega) - C_\varepsilon^m(\omega) = \frac{\mathcal{N}_\varepsilon(\omega) - \mathcal{N}_\varepsilon^m(\omega)}{\mathcal{D}_\varepsilon(\omega)} + \frac{\mathcal{N}_\varepsilon^m(\omega)}{\mathcal{D}_\varepsilon^m(\omega)\mathcal{D}_\varepsilon(\omega)} (\mathcal{D}_\varepsilon^m(\omega) - \mathcal{D}_\varepsilon(\omega)). \quad (5.33)$$

Introduce

$$b_m(x, \omega) = \left| \frac{1}{A_m(x, \omega)} - \frac{1}{A(x, \omega)} \right| \quad \text{and} \quad \mathcal{B}_\varepsilon^m(\omega) = \varepsilon \sum_{x \in (0, 1) \cap \varepsilon \mathbb{Z}} b_m\left(\frac{x}{\varepsilon}, \omega\right).$$

For any  $m \in \mathbb{N}^*$ , we get

$$|\mathcal{N}_\varepsilon^m(\omega) - \mathcal{N}_\varepsilon(\omega)| \leq \|f\|_{L^\infty} \varepsilon \sum_{x \in (0, 1) \cap \varepsilon \mathbb{Z}} b_m\left(\frac{x}{\varepsilon}, \omega\right) = \|f\|_{L^\infty} \mathcal{B}_\varepsilon^m(\omega), \quad (5.34)$$

$$|\mathcal{D}_\varepsilon^m(\omega) - \mathcal{D}_\varepsilon(\omega)| \leq \varepsilon \sum_{x \in (0, 1) \cap \varepsilon \mathbb{Z}} b_m\left(\frac{x}{\varepsilon}, \omega\right) = \mathcal{B}_\varepsilon^m(\omega), \quad (5.35)$$

$$|\mathcal{N}_\varepsilon^m(\omega)| \leq \|f\|_{L^\infty} \mathcal{D}_\varepsilon^m(\omega). \quad (5.36)$$



Using the ergodic theorem for the stationary functions  $A^{-1}$ ,  $A_m^{-1}$  and  $b_m$ , we have that, for any  $m \in \mathbb{N}^*$ , almost surely,

$$\lim_{\varepsilon \rightarrow 0} \mathcal{D}_\varepsilon(\omega) = \frac{1}{A^\star}, \quad \lim_{\varepsilon \rightarrow 0} \mathcal{D}_\varepsilon^m(\omega) = \frac{1}{A_m^\star}, \quad \lim_{\varepsilon \rightarrow 0} \mathcal{B}_\varepsilon^m(\omega) = \mathcal{B}_\star^m := \mathbb{E} \left[ \left| \frac{1}{A_m(0, \cdot)} - \frac{1}{A(0, \cdot)} \right| \right]. \quad (5.37)$$

We introduce

$$\Omega_{\text{conv}} = \left\{ \omega \in \Omega; \lim_{\varepsilon \rightarrow 0} \mathcal{D}_\varepsilon(\omega) = \frac{1}{A^\star} \text{ and } \forall m \in \mathbb{N}^*, \lim_{\varepsilon \rightarrow 0} \mathcal{D}_\varepsilon^m(\omega) = \frac{1}{A_m^\star}, \lim_{\varepsilon \rightarrow 0} \mathcal{B}_\varepsilon^m(\omega) = \mathcal{B}_\star^m \right\}$$

and we deduce that  $\mathbb{P}(\Omega_{\text{conv}}) = 1$ .

Let  $\omega \in \Omega_{\text{conv}}$ . In view of (5.37), we know that there exists  $\varepsilon_0^m(\omega)$  such that, for any  $\varepsilon < \varepsilon_0^m(\omega)$ , we have

$$\frac{1}{2A^\star} \leq \mathcal{D}_\varepsilon(\omega), \quad \frac{1}{2A_m^\star} \leq \mathcal{D}_\varepsilon^m(\omega) \leq \frac{3}{2A_m^\star}. \quad (5.38)$$

We thus infer from (5.33), (5.38), (5.34), (5.36) and (5.35) that, for any  $\omega \in \Omega_{\text{conv}}$ , any  $m \in \mathbb{N}^*$  and any  $\varepsilon < \varepsilon_0^m(\omega)$ , we have

$$\begin{aligned} |C_\varepsilon(\omega) - C_\varepsilon^m(\omega)| &\leq 2A^\star |\mathcal{N}_\varepsilon(\omega) - \mathcal{N}_\varepsilon^m(\omega)| + 4A^\star A_m^\star |\mathcal{N}_\varepsilon^m(\omega)| |\mathcal{D}_\varepsilon^m(\omega) - \mathcal{D}_\varepsilon(\omega)| \\ &\leq 2A^\star \|f\|_{L^\infty} \mathcal{B}_\varepsilon^m(\omega) + 4A^\star A_m^\star \|f\|_{L^\infty} \mathcal{D}_\varepsilon^m(\omega) \mathcal{B}_\varepsilon^m(\omega) \\ &\leq 2A^\star \|f\|_{L^\infty} \mathcal{B}_\varepsilon^m(\omega) + 6A^\star \|f\|_{L^\infty} \mathcal{B}_\varepsilon^m(\omega). \end{aligned} \quad (5.39)$$

Hence, for any  $\omega \in \Omega_{\text{conv}}$  and any  $m \in \mathbb{N}^*$ , we have

$$\limsup_{\varepsilon \rightarrow 0} |C_\varepsilon(\omega) - C_\varepsilon^m(\omega)| \leq 8A^\star \|f\|_{L^\infty} \mathcal{B}_m^\star.$$

The dominated convergence theorem implies that  $\lim_{m \rightarrow \infty} \mathcal{B}_m^\star = 0$ , hence, for any  $\omega \in \Omega_{\text{conv}}$ , we have

$$\lim_{m \rightarrow \infty} \limsup_{\varepsilon \rightarrow 0} |C_\varepsilon(\omega) - C_\varepsilon^m(\omega)| = 0.$$

Since  $\mathbb{P}(\Omega_{\text{conv}}) = 1$ , we have proved the claim (5.32).

We now proceed and deduce from (5.30) and (5.31) that

$$\begin{aligned} u_\varepsilon^m(z, \omega) &= \varepsilon \sum_{x \in (0, z) \cap \varepsilon \mathbb{Z}} A_m \left( \frac{x}{\varepsilon}, \omega \right)^{-1} (C_\varepsilon^m(\omega) - F_\varepsilon(x)), \\ u_\varepsilon(z, \omega) &= \varepsilon \sum_{x \in (0, z) \cap \varepsilon \mathbb{Z}} A \left( \frac{x}{\varepsilon}, \omega \right)^{-1} (C_\varepsilon(\omega) - F_\varepsilon(x)), \end{aligned}$$

hence

$$|u_\varepsilon^m(z, \omega) - u_\varepsilon(z, \omega)| \leq |C_\varepsilon^m(\omega) - C_\varepsilon(\omega)| \mathcal{D}_\varepsilon^m(\omega) + (|C_\varepsilon(\omega)| + \|f\|_{L^\infty}) \mathcal{B}_\varepsilon^m(\omega).$$

Using that  $|C_\varepsilon(\omega)| \leq \|f\|_{L^\infty}$ , we deduce that

$$\|u_\varepsilon^m(\cdot, \omega) - u_\varepsilon(\cdot, \omega)\|_{\ell_\varepsilon^2} \leq |C_\varepsilon^m(\omega) - C_\varepsilon(\omega)| \mathcal{D}_\varepsilon^m(\omega) + 2\|f\|_{L^\infty} \mathcal{B}_\varepsilon^m(\omega).$$

For any  $\omega \in \Omega_{\text{conv}}$ , any  $m \in \mathbb{N}^*$  and any  $\varepsilon < \varepsilon_0^m(\omega)$ , using (5.39) and (5.38), we obtain that

$$\|u_\varepsilon^m(\cdot, \omega) - u_\varepsilon(\cdot, \omega)\|_{\ell_\varepsilon^2} \leq 8A^\star \|f\|_{L^\infty} \mathcal{B}_\varepsilon^m(\omega) \frac{3}{2A_m^\star} + 2\|f\|_{L^\infty} \mathcal{B}_\varepsilon^m(\omega),$$

hence, for any  $\omega \in \Omega_{\text{conv}}$  and any  $m \in \mathbb{N}^*$ ,

$$\limsup_{\varepsilon \rightarrow 0} \|u_\varepsilon^m(\cdot, \omega) - u_\varepsilon(\cdot, \omega)\|_{\ell_\varepsilon^2} \leq 8A^* \|f\|_{L^\infty} \mathcal{B}_*^m \frac{3}{2A_m^*} + 2\|f\|_{L^\infty} \mathcal{B}_*^m,$$

and thus, almost surely,

$$\lim_{m \rightarrow \infty} \limsup_{\varepsilon \rightarrow 0} \|u_\varepsilon^m(\cdot, \omega) - u_\varepsilon(\cdot, \omega)\|_{\ell_\varepsilon^2} = 0. \quad (5.40)$$

Collecting (5.27), (5.28), (5.29) and (5.40), we obtain that

$$\limsup_{\varepsilon \rightarrow 0} \|u_\varepsilon(\cdot, \omega) - u^*\|_{\ell_\varepsilon^2} = 0 \quad \text{a.s.},$$

which is the convergence (5.24).  $\square$

### The case of Weibull laws

Following Section 5.2.3, assume that the conductances are given by (5.17), i.e. are distributed according to the Weibull law of parameter  $(\lambda^4, k/4)$ . For any  $k > 0$ , Assumption (5.3) is not satisfied. However, when  $k > 4$ , Assumption (5.21) is satisfied: in view of Theorem 5.13, homogenization holds and the homogenized coefficient is given by

$$A^* = \frac{\lambda^4}{\Gamma(1 - 4/k)}, \quad (5.41)$$

where  $\Gamma$  is the Euler Gamma function defined for any  $z > 0$  by

$$\Gamma(z) = \int_0^\infty t^{z-1} \exp(-t) dt. \quad (5.42)$$

The variance of  $A_N^*$  is finite if and only if  $k > 8$ . In the sequel, we work in the range  $k > 8$ .

## 5.3 A parameter fitting problem

We now describe the problem we consider, first in the general case (Section 5.3.1), next in the one-dimensional case (Section 5.3.2). In that latter section, we also motivate our choice of macroscopic quantities from which we fit the parameters of the Weibull law.

### 5.3.1 General case

We assume that we are given two observed quantities, the first coefficient of the macroscopic permeability matrix (see (5.18))

$$K_N^{*,\text{obs}}(\omega) = e_1 \cdot A_N^{*,\text{obs}}(\omega) e_1$$

and its relative variance  $S_N^{\text{obs}}$  for some parameter  $\theta_{\text{obs}} = (\lambda_{\text{obs}}, k_{\text{obs}})$  of the Weibull law. Note that the relative variance crucially depends on the size  $N^d$  of the finite box on which it is measured (in contrast to the apparent permeability, which converges to a finite value when  $N \rightarrow \infty$ ). We assume here that we know this size. In practice, these three quantities,  $N$ ,  $K_N^{*,\text{obs}}$  and  $S_N^{\text{obs}}$ , can be obtained by physical experiments. We therefore assume that there exists  $\theta_{\text{obs}}$  and  $N$  such that

$$\mathbb{E}[K_N^*(\cdot, \theta_{\text{obs}})] = K_N^{*,\text{obs}}, \quad \text{VarR}[K_N^*(\cdot, \theta_{\text{obs}})] = S_N^{\text{obs}}, \quad (5.43)$$

where, we recall,

$$\text{VarR}[K_N^*(\cdot, \theta_{\text{obs}})] = \frac{\text{Var}[K_N^*(\cdot, \theta_{\text{obs}})]}{(\mathbb{E}[K_N^*(\cdot, \theta_{\text{obs}})])^2}.$$

Given  $N$ ,  $K_N^{*,\text{obs}}$  and  $S_N^{\text{obs}}$ , our aim is to recover (an approximation of)  $\theta_{\text{obs}}$ . To that aim, we consider the function

$$F_{N,M} : \begin{cases} (\mathbb{R}_+^*)^2 & \rightarrow \mathbb{R}_+ \\ \theta & \mapsto \left( \frac{\overline{K}_{N,M}^*(\theta)}{K_N^{*,\text{obs}}} - 1 \right)^2 + \left( \frac{S_{N,M}(\theta)}{S_N^{\text{obs}}} - 1 \right)^2 \end{cases} \quad (5.44)$$

which penalizes the sum of the (relative) errors between

- on the one hand,  $\overline{K}_{N,M}^*(\theta)$  (which is an empirical estimator of  $\mathbb{E}[K_N^*(\cdot, \theta)]$  when  $M$  is large, see (5.10) and (5.18)) and  $K_N^{*,\text{obs}}$
- and, on the other hand,  $S_{N,M}(\theta)$  (which is an empirical estimator of the relative variance of  $K_N^*(\omega, \theta)$  when  $M$  is large, see (5.12) and (5.18)) and  $S_N^{\text{obs}}$ .

Of course, different weights could be assigned to the error on the permeability and the error on its relative variance. We eventually cast our parameter fitting problem in the form of the optimization problem

$$\inf_{\theta=(\lambda,k) \in (0,\infty) \times \mathcal{K}} F_{N,M}(\theta),$$

where  $\mathcal{K} \subset (0, \infty)$  is the admissible set of parameters  $k$  such that homogenization holds (even if Assumption (5.3) is not satisfied for any  $k > 0$ ) and the variance of  $K_N^*$  is also well-defined. In the one-dimensional case we focus on in this article,  $\mathcal{K} = (8, \infty)$ .

Note that  $F_{N,M}(\theta)$  is random, as it depends on the realizations used to evaluate  $\overline{K}_{N,M}^*(\theta)$  and  $S_{N,M}(\theta)$  (see (5.10) and (5.12)). For any  $\theta$ , in the limit when  $M \rightarrow \infty$ ,  $F_{N,M}(\theta)$  converges almost surely to the deterministic limit

$$F_N(\theta) = \left( \frac{\mathbb{E}[K_N^*(\cdot, \theta)]}{K_N^{*,\text{obs}}} - 1 \right)^2 + \left( \frac{\text{VarR}[K_N^*(\cdot, \theta)]}{S_N^{\text{obs}}} - 1 \right)^2.$$

Under Assumption (5.43), we have

$$F_N(\theta) = \left( \frac{\mathbb{E}[K_N^*(\cdot, \theta)]}{\mathbb{E}[K_N^*(\cdot, \theta_{\text{obs}})]} - 1 \right)^2 + \left( \frac{\text{VarR}[K_N^*(\cdot, \theta)]}{\text{VarR}[K_N^*(\cdot, \theta_{\text{obs}})]} - 1 \right)^2.$$

### 5.3.2 The one-dimensional case

Recall that, in that case,  $K_N^* = A_N^*$  (see (5.18)).

#### Theoretical result

We first identify the limit when  $N \rightarrow \infty$  of  $F_N$ , which we recall reads

$$F_N(\theta) = \left( \frac{\mathbb{E}[A_N^*(\cdot, \theta)]}{\mathbb{E}[A_N^*(\cdot, \theta_{\text{obs}})]} - 1 \right)^2 + \left( \frac{\text{VarR}[A_N^*(\cdot, \theta)]}{\text{VarR}[A_N^*(\cdot, \theta_{\text{obs}})]} - 1 \right)^2. \quad (5.45)$$

Obviously, the first term above converges to

$$\left( \frac{A^*(\theta)}{A^*(\theta_{\text{obs}})} - 1 \right)^2,$$

with, following (5.41),

$$A^* = \frac{\lambda^4}{\Gamma(1 - 4/k)}$$

where  $\Gamma$  is the Euler Gamma function.

For the second term of (5.45), it is clear that  $\text{VarR}[A_N^*(\cdot, \theta)]$  vanishes in the limit  $N \rightarrow \infty$ , since  $A_N^*(\cdot, \theta)$  converges almost surely to a deterministic limit. Furthermore, equation (5.19) implies that

$$\text{Var}[A_N^*] = \frac{(A^*)^4}{N} \text{Var} \left[ \frac{1}{A(0, \cdot)} \right] + o\left(\frac{1}{N}\right). \quad (5.46)$$

The conductances are distributed according to a Weibull law (see (5.17)), therefore

$$\text{Var}[A_N^*] = \frac{\lambda^{16}}{N \Gamma(1 - 4/k)^4} \left( \frac{\Gamma(1 - 8/k)}{\lambda^8} - \frac{\Gamma(1 - 4/k)^2}{\lambda^8} \right) + o\left(\frac{1}{N}\right),$$

hence the relative variance reads

$$\text{VarR}[A_N^*] = \frac{1}{N} \left( \frac{\Gamma(1 - 8/k)}{\Gamma(1 - 4/k)^2} - 1 \right) + o\left(\frac{1}{N}\right), \quad (5.47)$$

which implies that

$$\lim_{N \rightarrow \infty} \frac{\text{VarR}[A_N^*(\cdot, \theta)]}{\text{VarR}[A_N^*(\cdot, \theta_{\text{obs}})]} = \frac{\frac{\Gamma(1-8/k)}{\Gamma(1-4/k)^2} - 1}{\frac{\Gamma(1-8/k_{\text{obs}})}{\Gamma(1-4/k_{\text{obs}})^2} - 1}.$$

In the one-dimensional case, we are thus able to identify the limit as  $N \rightarrow \infty$  of  $F_N(\theta)$ , which reads

$$F_\infty^{1D}(\theta) := \lim_{N \rightarrow \infty} F_N(\theta) = \left( \frac{\lambda^4}{\lambda_{\text{obs}}^4} \frac{\Gamma(1 - 4/k_{\text{obs}})}{\Gamma(1 - 4/k)} - 1 \right)^2 + \left( \frac{\frac{\Gamma(1-8/k)}{\Gamma(1-4/k)^2} - 1}{\frac{\Gamma(1-8/k_{\text{obs}})}{\Gamma(1-4/k_{\text{obs}})^2} - 1} - 1 \right)^2. \quad (5.48)$$

Obviously, this function is minimal (and vanishes) when  $\theta = \theta_{\text{obs}}$ . It turns out that this minimizer is the unique minimizer, as shown below.

**Lemma 5.14.** *The function  $F_\infty^{1D}$  defined by (5.48) has a unique minimizer, which is  $\theta_{\text{obs}}$ .*

Homogenization is an averaging process, which filters out many features of the microscopic coefficient  $A$ . These features cannot be recovered from the knowledge of macroscopic quantities. The above lemma shows (in the one-dimensional case) that, if one assumes a given form for the probability distribution of  $A$  (here, a Weibull distribution), then one is able to recover the two parameters of that law on the basis of two macroscopic quantities, the permeability and its relative variance.

It is also obvious from (5.41) that knowing the macroscopic permeability is not enough to uniquely determine the two parameters  $\lambda$  and  $k$  of the Weibull law. Additional information is needed. Our choice of considering the relative variance of the permeability is motivated by the following observation. This quantity, in the one-dimensional case, only depends (at first order in  $N$ ) on  $k$  and does not depend on  $\lambda$ , as can be seen on (5.47).

Knowing this quantity is therefore very useful to estimate the parameter  $k$ . Once  $k$  has been identified, knowing the macroscopic permeability yields, using (5.41), an estimation of the parameter  $\lambda$ .

Of course, it is likely that the knowledge of quantities of interest alternate to the relative variance of the permeability may also prove useful to determine the unknown parameters. Note also that such alternate relevant quantities should be “different enough” from the homogenized permeability to indeed bring new information. We do not pursue in that direction.

We plot on Figure 5.4 the function  $\theta \mapsto F_{\infty}^{1D}(\theta)$  for  $\lambda_{\text{obs}} = 1$  and  $k_{\text{obs}} = 15$ . We observe that the function is not degenerated at its minimum, in the sense that its Hessian matrix at  $\theta_{\text{obs}}$  is positive definite, with eigenvalues equal to 16 and 0.04. We thus expect that a standard algorithm (such as the Newton algorithm) will be able to converge to the minimizer of  $F_{\infty}^{1D}$ . This is indeed the case, as shown in Section 5.4.

*Proof of Lemma 5.14.* The proof consists of three steps: in Step 1, we recall (and prove for the sake of completeness) that  $\ln \Gamma$  is a convex function. In Step 2, we prove that the function

$$\zeta : k \mapsto \frac{\Gamma(1 - 8/k)}{\Gamma(1 - 4/k)^2}$$

is monotone (hence injective). We conclude in Step 3.

**Step 1.** From (5.42), we compute that, for any  $z > 0$ ,

$$\begin{aligned} \Gamma'(z) &= \int_0^{\infty} (\ln t) t^{z-1} \exp(-t) dt, \\ \Gamma''(z) &= \int_0^{\infty} (\ln t)^2 t^{z-1} \exp(-t) dt, \end{aligned}$$

therefore  $\Gamma''(z) > 0$  and  $\Gamma$  is positive and convex on  $(0, \infty)$ . In addition, we have

$$(\ln \Gamma)''(z) = \frac{\Gamma(z)\Gamma''(z) - (\Gamma'(z))^2}{\Gamma^2(z)}$$

which is positive, in view of the Cauchy-Schwartz inequality:

$$(\Gamma'(z))^2 = \left( \int_0^{\infty} (\ln t) \sqrt{t^{z-1} \exp(-t)} dt \right)^2 < \Gamma''(z)\Gamma(z).$$

Therefore,  $\ln \Gamma$  is a strictly convex function.

**Step 2.** We define the function

$$\zeta : k \mapsto \frac{\Gamma(1 - 8/k)}{\Gamma(1 - 4/k)^2},$$

the derivative of which reads

$$\zeta'(k) = \frac{\Gamma(1 - 8/k)}{\Gamma(1 - 4/k)^2} \left( \frac{8}{k^2} \frac{\Gamma'(1 - 8/k)}{\Gamma(1 - 8/k)} - \frac{8}{k^2} \frac{\Gamma'(1 - 4/k)}{\Gamma(1 - 4/k)} \right).$$

For any  $k > 0$ , we have that  $1 - 8/k < 1 - 4/k$ . As a consequence of  $\ln \Gamma$  being strictly convex, we have that its derivative is increasing, therefore

$$\frac{\Gamma'(1 - 8/k)}{\Gamma(1 - 8/k)} < \frac{\Gamma'(1 - 4/k)}{\Gamma(1 - 4/k)}.$$

We can now conclude that  $\zeta'(k) < 0$ , hence  $\zeta$  is decreasing.

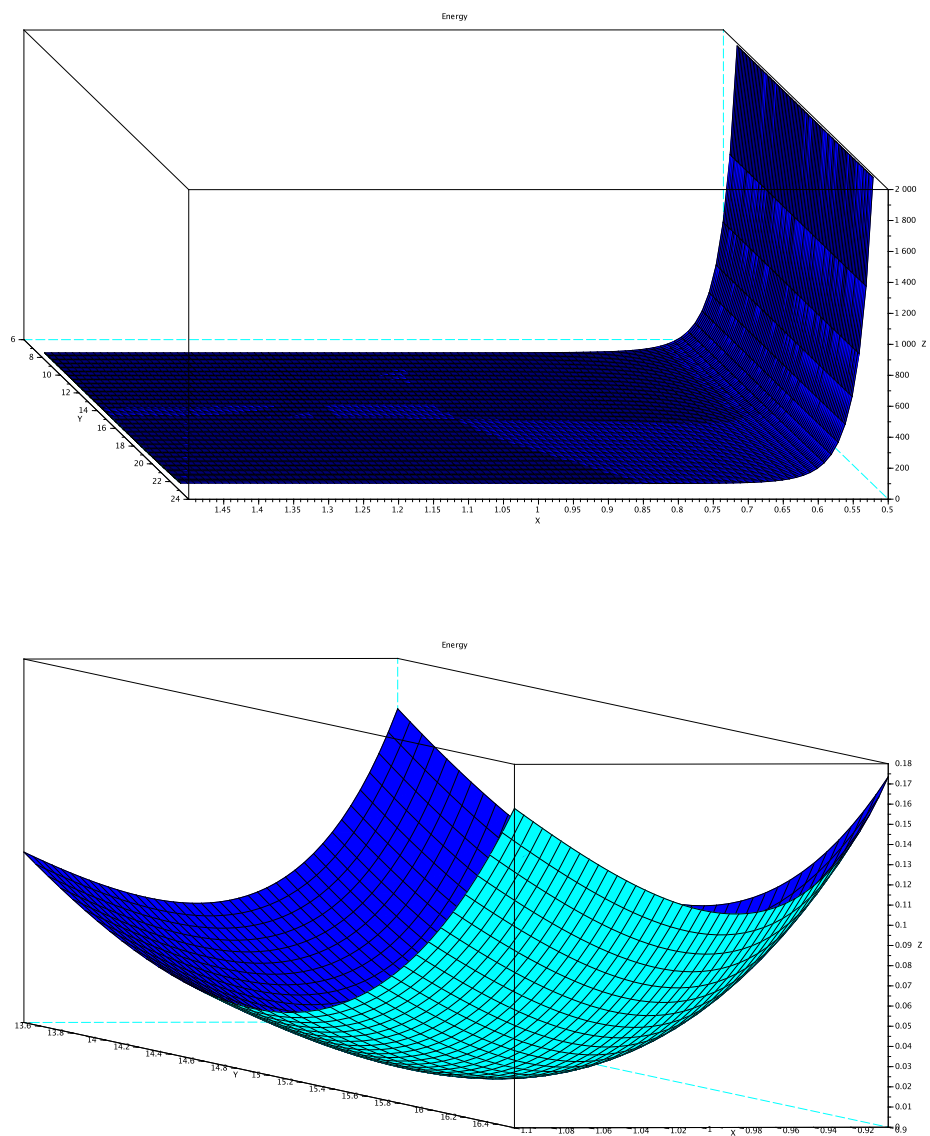


Figure 5.4 – Plot of  $\theta \mapsto F_{\infty}^{1D}(\theta)$  for  $\lambda_{\text{obs}} = 1$  and  $k_{\text{obs}} = 15$ . The bottom plot is a zoom of the top plot around the minimizer.

**Step 3.** By definition of  $\zeta$ , we have that

$$F_\infty^{1D}(\theta) = \left( \frac{\lambda^4}{\lambda_{\text{obs}}^4} \frac{\Gamma(1 - 4/k_{\text{obs}})}{\Gamma(1 - 4/k)} - 1 \right)^2 + \left( \frac{\zeta(k) - 1}{\zeta(k_{\text{obs}}) - 1} - 1 \right)^2.$$

We obviously have that  $\theta_{\text{obs}}$  is a minimizer of  $F_\infty^{1D}$ , with  $F_\infty^{1D}(\theta_{\text{obs}}) = 0$ . Conversely, let  $\theta$  be a minimizer of  $F_\infty^{1D}$ . We thus have  $F_\infty^{1D}(\theta) = 0$ , which implies that  $\zeta(k) = \zeta(k_{\text{obs}})$ . The function  $\zeta$  being monotone, this implies that  $k = k_{\text{obs}}$ . Since the first term in  $F_\infty^{1D}(\theta)$  also vanishes, we obtain that  $\lambda = \lambda_{\text{obs}}$  as well. This concludes the proof.  $\square$

### Practical situation

In the general (i.e. multi-dimensional) case, we have introduced in (5.44) the function  $F_{N,M}$  that we wish to minimize. Turning next to the one-dimensional case, we have theoretically identified its limit when  $M \rightarrow \infty$  and  $N \rightarrow \infty$ . In practice, we cannot take any of these limits, and have thus to work with  $F_{N,M}$  defined by

$$F_{N,M}(\theta) = \left( \frac{\overline{K}_{N,M}^*(\theta)}{K_N^{\text{obs}}} - 1 \right)^2 + \left( \frac{S_{N,M}(\theta)}{S_N^{\text{obs}}} - 1 \right)^2.$$

In view of (5.46), we see that, when  $M \rightarrow \infty$  and  $N \rightarrow \infty$ , the relative variance  $S_{N,M}(\theta)$  is close to

$$S_{N,M}(\theta) \approx \frac{(A^*)^2}{N} \text{Var} \left[ \frac{1}{A(0, \cdot)} \right] = \frac{\mathbb{E}[W^{-8}] - (\mathbb{E}[W^{-4}])^2}{N (\mathbb{E}[W^{-4}])^2} = \frac{\mathbb{E}[W^{-8}]}{N (\mathbb{E}[W^{-4}])^2} - \frac{1}{N},$$

where  $W$  is a random variable distributed according to the Weibull law  $\mathcal{W}(1, k)$ . Likewise,

$$\overline{K}_{N,M}^*(\theta) \approx A^* = \frac{\lambda^4}{\mathbb{E}[W^{-4}]}.$$

Let  $\{u_i(\omega)\}_{i=1}^N$  be a sequence of i.i.d. random variables uniformly distributed in  $[0, 1]$ . We define

$$w_i(k, \omega) := (-\ln(1 - u_i(\omega)))^{-1/k}, \quad (5.49)$$

so that  $\{1/w_i(k, \omega)\}_{i=1}^N$  are i.i.d. random variables distributed according to  $\mathcal{W}(1, k)$ . In the sequel, we approximate the function to minimize by

$$\tilde{F}_N^{1D}(\theta, \omega) = \left( \frac{\lambda^4}{K_N^{\text{obs}}} \left[ \frac{1}{N} \sum_{i=1}^N w_i^4(k, \omega) \right]^{-1} - 1 \right)^2 + \left( \frac{1}{S_N^{\text{obs}}} \left[ \frac{\sum_{i=1}^N w_i^8(k, \omega)}{(\sum_{i=1}^N w_i^4(k, \omega))^2} - \frac{1}{N} \right] - 1 \right)^2. \quad (5.50)$$

This function is consistent in the sense that it almost surely converges, when  $N \rightarrow \infty$ , to the exact function (5.48). On the other hand,  $\tilde{F}_N^{1D}(\theta, \omega)$  is random, and thus somewhat mimics the difficulties that one would encounter in the multi-dimensional case when working with  $F_{N,M}(\theta)$ .

## 5.4 Numerical results

We briefly explain in Section 5.4.1 how in practice we minimize the function (5.50), before turning in Section 5.4.2 to our numerical results. As pointed out in the introduction, we only consider here the one-dimensional case, and postpone the study of two-dimensional examples to the future work [LMOS].

### 5.4.1 Optimization algorithm

We provide in Appendix 5.5 expressions for the first and second derivatives of the function  $\tilde{F}_N^{1D}(\theta, \omega)$  defined by (5.50) with respect to  $\theta = (\lambda, k)$ . We are thus in position to use the Newton algorithm, and compute a sequence  $\theta_j$  according to

$$\theta_{j+1} = \theta_j - \mu_j \left[ \mathcal{H} \left( \tilde{F}_N^{1D} \right) (\theta_j) \right]^{-1} \nabla \tilde{F}_N^{1D}(\theta_j), \quad (5.51)$$

where  $\mathcal{H} \left( \tilde{F}_N^{1D} \right) \in \mathbb{R}^{2 \times 2}$  is the Hessian matrix of  $\tilde{F}_N^{1D}$  and  $\nabla \tilde{F}_N^{1D} \in \mathbb{R}^2$  is the gradient of  $\tilde{F}_N^{1D}$  (for the sake of simplicity, we keep implicit the dependence with respect to  $\omega$ ). In turn,  $\mu_j > 0$  is the step-size by which we move. To choose  $\mu_j$ , we have used a line-search algorithm (along the descent direction prescribed by the Newton algorithm) using Goldstein (respectively Armijo) rule to increase (respectively decrease) the step-size.

We note that the function  $\theta \mapsto F_\infty^{1D}(\theta)$  is not convex. It is possible to find some  $\theta$  such that the Hessian matrix  $\mathcal{H} \left( F_\infty^{1D} \right) (\theta)$  is not positive definite, but rather has (at least) one negative eigenvalue. We thus cannot expect the function  $\theta \mapsto \tilde{F}_N^{1D}(\theta)$  to be convex (even for large values of  $N$ ), and the Newton algorithm to be globally convergent. We are therefore careful to start the Newton iterations from an initial guess  $\theta_0$  (given by physical experiments) that we hope to be close enough to the minimizer of  $\tilde{F}_N^{1D}$ .

### 5.4.2 Numerical results

In all what follows, we set  $N = 10^5$ .

#### Robustness of the algorithm with respect to the initial guess

Our first numerical test is a simple one, to check whether the Newton algorithm (5.51) is indeed able to minimize the function  $\theta \mapsto \tilde{F}_N^{1D}(\theta, \omega)$ . We pick once for all one realization of the i.i.d. random variables  $\{u_i(\omega)\}_{1 \leq i \leq N}$  (which, we recall, are uniformly distributed in  $[0, 1]$ ). We then build  $\{w_i(k, \omega)\}_{1 \leq i \leq N}$  according to (5.49) and consider the function  $\theta \mapsto \tilde{F}_N^{1D}(\theta, \omega)$  defined by (5.50), where the observed quantities are defined by

$$K_N^{*,\text{obs}} = \lambda_{\text{obs}}^4 \left[ \frac{1}{N} \sum_{i=1}^N w_i^4(k_{\text{obs}}, \omega) \right]^{-1}, \quad S_N^{\text{obs}} = \frac{\sum_{i=1}^N w_i^8(k_{\text{obs}}, \omega)}{\left( \sum_{i=1}^N w_i^4(k_{\text{obs}}, \omega) \right)^2} - \frac{1}{N},$$

with  $\lambda_{\text{obs}} = 1$  and  $k_{\text{obs}} = 15$ . The function  $\theta \mapsto \tilde{F}_N^{1D}(\theta, \omega)$  obviously vanishes at  $\theta_{\text{obs}} = (\lambda_{\text{obs}}, k_{\text{obs}})$ .

We run the Newton algorithm (5.51) starting from several initial guesses  $\theta_0$ , and check that it indeed always converges to  $\theta_{\text{obs}}$  in a limited number of iterations. We also observe that, for some initial guesses, using an adaptive step-size  $\mu_j$  as in (5.51) is critical: in contrast, if one uses the step-size  $\mu_j = 1$ , then the algorithm may not converge, or converges after a much larger number of iterations.

#### Robustness with respect to statistical noise

For our second test, we proceed as follows. We first set  $\theta_{\text{ref}} = (\lambda_{\text{ref}}, k_{\text{ref}}) = (1, 15)$  and pick one realization of the i.i.d. random variables  $\{u_i(\bar{\omega})\}_{1 \leq i \leq N}$  (which, we recall, are uniformly



distributed in  $[0, 1]$ ). We then build  $\{w_i(k_{\text{ref}}, \bar{\omega})\}_{1 \leq i \leq N}$  according to (5.49) and define the macroscopic observed quantities as

$$K_N^{*,\text{obs}} = \lambda_{\text{ref}}^4 \left[ \frac{1}{N} \sum_{i=1}^N w_i^4(k_{\text{ref}}, \bar{\omega}) \right]^{-1}, \quad S_N^{\text{obs}} = \frac{\sum_{i=1}^N w_i^8(k_{\text{ref}}, \bar{\omega})}{\left( \sum_{i=1}^N w_i^4(k_{\text{ref}}, \bar{\omega}) \right)^2} - \frac{1}{N}. \quad (5.52)$$

We now fix the initial guess  $\theta_0 = (1.1, 16.5)$  (10% off the reference value  $\theta_{\text{ref}}$ ) and set  $M = 500$ . For any  $1 \leq m \leq M$ , we perform the following procedure:

- we draw a realization of  $N$  i.i.d. random variables  $\{u_i(\omega_m)\}_{1 \leq i \leq N}$  which is independent of the realization  $\{u_i(\omega_{m'})\}_{1 \leq i \leq N}$  for any  $m' \neq m$ , and independent of the realization  $\{u_i(\bar{\omega})\}_{1 \leq i \leq N}$  used to compute  $K_N^{*,\text{obs}}$  and  $S_N^{\text{obs}}$  in (5.52);
- using  $\{u_i(\omega_m)\}_{1 \leq i \leq N}$ , we build  $w_i(k, \omega_m)$  according to (5.49) and we consider the function  $\theta \mapsto \tilde{F}_N^{1\text{D}}(\theta, \omega_m)$  defined by (5.50), i.e.

$$\tilde{F}_N^{1\text{D}}(\theta, \omega_m) = \left( \frac{\lambda^4}{K_N^{*,\text{obs}}} \left[ \frac{1}{N} \sum_{i=1}^N w_i^4(k, \omega_m) \right]^{-1} - 1 \right)^2 + \left( \frac{1}{S_N^{\text{obs}}} \left[ \frac{\sum_{i=1}^N w_i^8(k, \omega_m)}{\left( \sum_{i=1}^N w_i^4(k, \omega_m) \right)^2} - \frac{1}{N} \right] - 1 \right)^2.$$

Recall that the macroscopic observed quantities are independent of  $\omega_m$ .

- we run the Newton algorithm (5.51) to minimize the function  $\theta \mapsto \tilde{F}_N^{1\text{D}}(\theta, \omega_m)$ . The optimal parameter found by the algorithm depends on  $\omega_m$  and is denoted  $\theta_{\text{opt}}(\omega_m)$ . Since the realization  $\omega_m$  is different from the reference realization  $\bar{\omega}$ , we have in general  $\theta_{\text{opt}}(\omega_m) \neq \theta_{\text{ref}}$ .

We show on Figure 5.5 the histogram of the optimal parameters  $\theta_{\text{opt}}(\omega_m)$  for  $1 \leq m \leq M$ . We see that these histograms are centered close to the reference value ( $k_{\text{ref}}$ , resp.  $\lambda_{\text{ref}}$ ). There is however a small bias, i.e.  $\mathbb{E}(\theta_{\text{opt}}) \neq \theta_{\text{ref}}$ . We also observe that the width of these histograms (related to the variance of  $k_{\text{opt}}$  and  $\lambda_{\text{opt}}$ ) is quite small.

**Remark 5.15.** *Of course, the variance of  $k_{\text{opt}}$  and  $\lambda_{\text{opt}}$  is related to  $N$ . In the limit  $N \rightarrow \infty$ , the function  $\tilde{F}_N^{1\text{D}}(\theta, \omega)$  almost surely converges to the deterministic limit  $F_\infty^{1\text{D}}(\theta)$  defined by (5.48), and we thus expect  $k_{\text{opt}}$  and  $\lambda_{\text{opt}}$  to almost surely converge to a deterministic limit. But this is not the regime we are interested in, since in practice (in the two-dimensional case), we have to work with the random function  $F_{N,M}$ .*

We next compare the variance of  $\theta_{\text{opt}}$  with the amount of randomness introduced in the function  $\tilde{F}_N^{1\text{D}}(\cdot, \omega)$  defined by (5.50). By construction,

$$\tilde{F}_N^{1\text{D}}(\theta, \omega) = \left( \frac{K_N^*(\theta, \omega)}{K_N^{*,\text{obs}}} - 1 \right)^2 + \left( \frac{S_N(k, \omega)}{S_N^{\text{obs}}} - 1 \right)^2$$

with

$$S_N(k, \omega) = \left[ \frac{\sum_{i=1}^N w_i^8(k, \omega)}{\left( \sum_{i=1}^N w_i^4(k, \omega) \right)^2} - \frac{1}{N} \right],$$

which is an approximation of the relative variance of  $K_N^*(\theta, \omega)$ . We show on Figure 5.6 the histograms, for  $1 \leq m \leq M$ , of  $K_N^*(\theta_0, \omega_m)$  and of  $S_N(k_0, \omega_m)$ , for the initial guess parameter  $\theta_0 = (1.1, 16.5)$ .

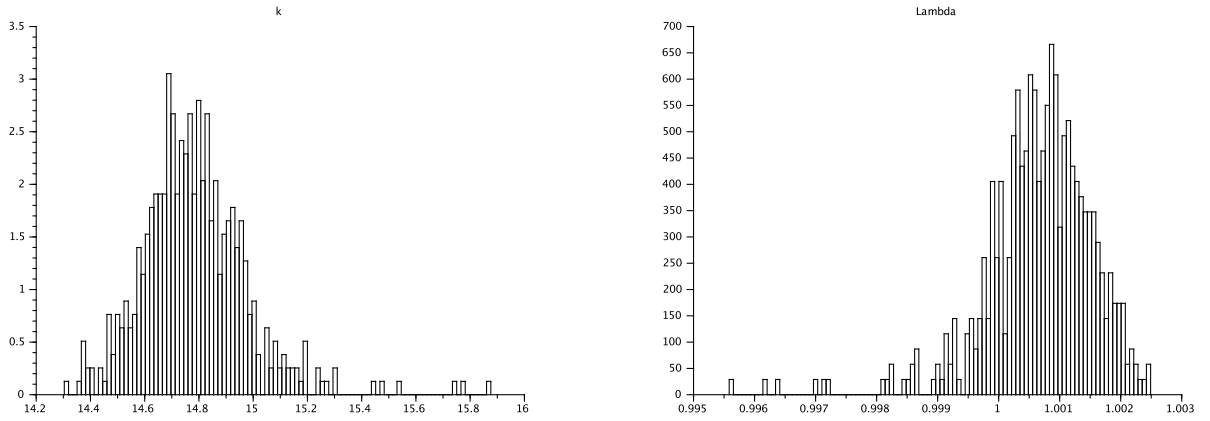


Figure 5.5 – Top: distribution of  $k_{\text{opt}}(\omega)$ . Bottom: distribution of  $\lambda_{\text{opt}}(\omega)$ .

On this test-case, we compute that  $\text{Var}[\lambda_{\text{opt}}] \approx 7.9 \cdot 10^{-7}$  and  $\text{Var}[k_{\text{opt}}] \approx 3.8 \cdot 10^{-2}$ , thus

$$\text{VarR}[\lambda_{\text{opt}}] \approx 7.9 \cdot 10^{-7} \quad \text{and} \quad \text{VarR}[k_{\text{opt}}] \approx 1.7 \cdot 10^{-4}.$$

On the other hand,  $A^*(\theta_0) \approx 1.2$ ,  $\text{Var}[A_N^*(\theta_0)] \approx 2.0 \cdot 10^{-6}$  and  $\text{Var}[S_N(k_0)] \approx 4.5 \cdot 10^{-15}$ , thus

$$\text{VarR}[K_N^*(\theta_0)] \approx 1.4 \cdot 10^{-6} \quad \text{and} \quad \text{VarR}[S_N(k_0)] \approx 10^{-3}.$$

We thus observe that the relative variance of the optimal parameters is roughly of the same order of magnitude as the relative variance introduced in the function to minimize. Given the amount of noise present in the system, our procedure robustly identifies the optimal parameters of the microscopic distribution.

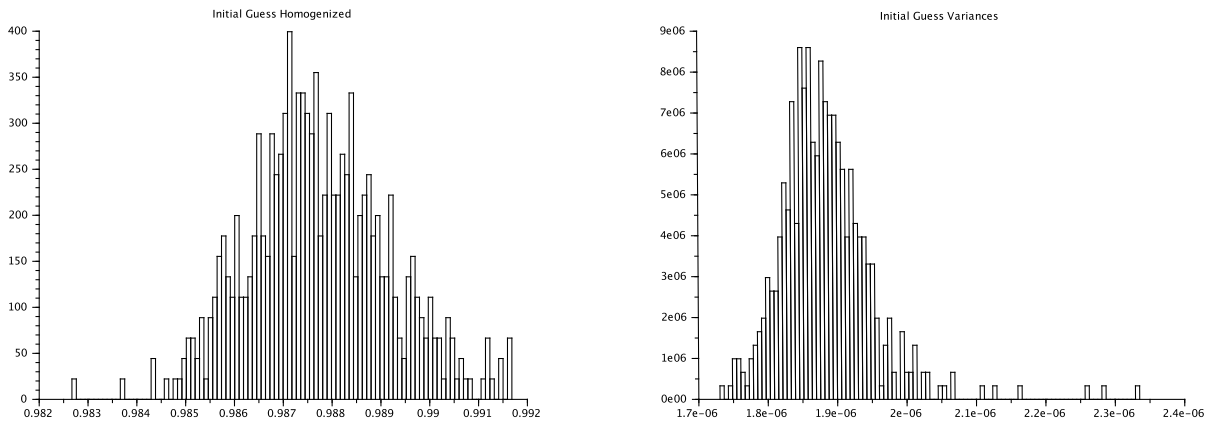


Figure 5.6 – Top: distribution of  $K_N^*(\theta_0, \omega)$ . Bottom: distribution of  $S_N(k_0, \omega)$ .

## 5.5 Appendix: Computation of the derivatives of (5.50)

We introduce

$$f(\lambda, k) := \lambda^4 \left( \sum_{i=1}^N w_i^4(k) \right)^{-1}$$

and

$$g(k) := \left( \sum_{i=1}^N w_i^8(k) \right) \left( \sum_{i=1}^N w_i^4(k) \right)^{-2}$$

where  $w_i(k)$  is defined by (5.49), and recast the function (5.50) as

$$\tilde{F}_N^{1D}(\theta) = \left( \frac{N}{K_N^{*,\text{obs}}} f(\lambda, k) - 1 \right)^2 + \left( \frac{1}{S_N^{\text{obs}}} \left( g(k) - \frac{1}{N} \right) - 1 \right)^2,$$

where we have kept implicit the dependence with respect to  $\omega$ . Computing the derivatives of  $\tilde{F}_N^{1D}$  therefore amounts to computing those of  $f$  and  $g$ .

A tedious but straightforward computation leads to the following expressions:

$$\begin{aligned} g'(k) &= \frac{8}{k} \left( \sum_i w_i^4 \right)^{-2} \left[ \frac{\sum_i w_i^8}{\sum_i w_i^4} \sum_i \ln(w_i) w_i^4 - \sum_i \ln(w_i) w_i^8 \right], \\ g''(k) &= \frac{16}{k^2} \left( \sum_i w_i^4 \right)^{-2} \left[ \sum_i \ln(w_i) w_i^8 (1 + 4 \ln w_i) - \frac{\sum_i w_i^8}{\sum_i w_i^4} \sum_i \ln(w_i) w_i^4 (1 + 2 \ln w_i) \right] \\ &\quad - \frac{32}{k^2} \left( \sum_i w_i^4 \right)^{-3} \left[ 4 \left( \sum_i w_i^8 \ln(w_i) \right) \left( \sum_i w_i^4 \ln(w_i) \right) - 3 \left( \sum_i w_i^4 \ln(w_i) \right)^2 \frac{\sum_i w_i^8}{\sum_i w_i^4} \right], \end{aligned}$$

whereas

$$\begin{aligned} \partial_\lambda f &= \frac{4}{\lambda} f(\lambda, k), \\ \partial_{\lambda\lambda}^2 f &= \frac{12}{\lambda^2} f(\lambda, k), \\ \partial_k f &= \frac{4\lambda^4}{k} \left( \sum_i \ln(w_i) w_i^4 \right) \left( \sum_i w_i^4 \right)^{-2}, \\ \partial_{kk}^2 f &= -\frac{8\lambda^4}{k^2} \left( \sum_i w_i^4 \right)^{-2} \left[ \sum_i \ln(w_i) w_i^4 (1 + 2 \ln w_i) \right] + \frac{32\lambda^4}{k^2} \left( \sum_i w_i^4 \right)^{-3} \left( \sum_i \ln(w_i) w_i^4 \right)^2, \\ \partial_{\lambda k}^2 f &= \frac{16\lambda^3}{k} \left( \sum_i \ln(w_i) w_i^4 \right) \left( \sum_i w_i^4 \right)^{-2} = \frac{4}{\lambda} \partial_k f. \end{aligned}$$

# Appendix A

## Neumann, Dirichlet and periodic approximations of the corrector problem

### A.1 Introduction

The purpose of this appendix, useful for **Chapter 4**, is to present three classical approximations of the random corrector problem (1.8), using a problem set on a truncated domain complemented by periodic, Dirichlet or Neumann boundary conditions. These three approximations are  $A_{\text{Per}}^{*,N}$  (see (A.4)–(A.7)),  $A_{\text{Dir}}^{*,N}$  (see (A.2)–(A.3)) and  $A_{\text{Neu}}^{*,N}$  (see (A.9)–(A.10)) of  $A^*$ .

For the sake of completeness, we also provide a complete proof of the fact that

$$A_{\text{Neu}}^{*,N}(\omega) \leq A_{\text{Per}}^{*,N}(\omega) \leq A_{\text{Dir}}^{*,N}(\omega), \quad \text{almost surely.} \quad (\text{A.1})$$

In section A.2, we introduce these three approximations  $A_{\text{Dir}}^{*,N}$ ,  $A_{\text{Neu}}^{*,N}$  and  $A_{\text{Per}}^{*,N}$ , and we recall the corresponding convergence results when  $N$  tends to infinity. Throughout this section,  $A$  need not be symmetric. In section A.3, we introduce a dual formulation in terms of flux. In the case when  $A$  is symmetric, it allows us to reformulate the corrector problem with Neumann and periodic boundary conditions as variational problems on the flux. In section A.4, we compare these variational problems to infer the ordering claimed in (A.1) in the case when the matrix  $A$  (and hence each of the approximations  $A_{\text{Dir}}^{*,N}$ ,  $A_{\text{Neu}}^{*,N}$  and  $A_{\text{Per}}^{*,N}$ ) is symmetric.

### A.2 Approximations of the corrector

In the celebrated article [BP04], A. Bourgeat and A. Piatnitski introduce three approximations of (1.8). They consider the corrector equation on the bounded domain  $Q_N = (0, N)^d$ , complemented with either homogeneous Dirichlet boundary conditions, periodic boundary conditions, or Neumann boundary conditions. We discuss the Dirichlet approximation in Section A.2.1, the periodic approximation in Section A.2.2, and the Neumann approximation in Section A.2.3. The two first approximations are carefully considered and the third one is rapidly processed, but it turns out the Neumann approximation is more subtle.

Throughout this appendix, we assume that  $A$  is uniformly elliptic, bounded and stationary, that is to say:

- (*boundedness*): The field  $A$  is bounded almost everywhere and almost surely by some constant  $\|A\|_{L^\infty(\Omega \times \mathbb{R}^d)}$ .
- (*ellipticity*): There exists some constant  $\alpha > 0$  such that, for almost all  $x \in \mathbb{R}^d$ , all  $p \in \mathbb{R}^d$  and almost surely in  $\omega \in \Omega$ ,

$$\alpha|p|^2 \leq p \cdot A(x, \omega)p.$$

- (*stationnarity*):  $A$  is stationary in the sense of (1.3).

### A.2.1 The Dirichlet approximation

Consider the problem

$$\begin{cases} -\operatorname{div} \left[ A(\cdot, \omega) \left( p + \nabla w_{\operatorname{Dir}}^{N,p}(\cdot, \omega) \right) \right] = 0 \text{ in } Q_N, \\ w_{\operatorname{Dir}}^{N,p}(\cdot, \omega) = 0 \text{ on } \partial Q_N. \end{cases} \quad (\text{A.2})$$

This problem is well-posed almost surely in  $\omega \in \Omega$ . The corresponding apparent homogenized matrix  $A_{\operatorname{Dir}}^{*,N}(\omega)$  is defined by

$$\forall p \in \mathbb{R}^d, \quad A_{\operatorname{Dir}}^{*,N}(\omega)p = \frac{1}{|Q_N|} \int_{Q_N} A(x, \omega) (p + \nabla w_{\operatorname{Dir}}^{N,p}(x, \omega)) dx. \quad (\text{A.3})$$

It is shown in [BP04, Theorem 2] that

$$\lim_{N \rightarrow \infty} A_{\operatorname{Dir}}^{*,N}(\omega) = A^* \text{ a.s.}$$

### A.2.2 The periodic approximation

Consider the problem

$$\begin{cases} -\operatorname{div} \left[ \tilde{A}_{\operatorname{Per}}^N(\cdot, \omega) \left( p + \nabla w_{\operatorname{Per}}^{N,p}(\cdot, \omega) \right) \right] = 0 \text{ in } \mathbb{R}^d, \\ w_{\operatorname{Per}}^{N,p}(\cdot, \omega) \text{ is } Q_N\text{-periodic,} \end{cases} \quad (\text{A.4})$$

where  $\tilde{A}_{\operatorname{Per}}^N$  is the  $Q_N$  periodic extension of  $A$ , and the first equality is to be understood in the periodic distributional sense (alternatively, one may consider Equation (A.4) on the torus  $\mathbb{R}^d/Q_N$ , the point being that, for instance, periodically repeated linear functions are not admissible as test functions nor solutions). This problem is well-posed up to the addition of a (possibly random) constant. The variational formulation of (A.4) is: find  $w_{\operatorname{Per}}^{N,p} \in V_{\operatorname{per}}(Q_N)$  such that

$$\forall \psi \in V_{\operatorname{per}}(Q_N), \quad \int_{Q_N} \nabla \psi \cdot A(x, \omega) \left( p + \nabla w_{\operatorname{Per}}^{N,p}(x, \omega) \right) dx = 0, \quad (\text{A.5})$$

where  $\cdot$  denotes the usual  $\mathbb{R}^d$  scalar product and where

$$V_{\operatorname{per}}(Q_N) = \left\{ v \in H_{\operatorname{loc}}^1(\mathbb{R}^d), v \text{ is } Q_N\text{-periodic, } \int_{Q_N} v = 0 \right\}. \quad (\text{A.6})$$

The corresponding apparent homogenized matrix  $A_{\operatorname{Per}}^{*,N}(\omega)$  is defined by

$$\forall p \in \mathbb{R}^d, \quad A_{\operatorname{Per}}^{*,N}(\omega)p = \frac{1}{|Q_N|} \int_{Q_N} A(x, \omega) (p + \nabla w_{\operatorname{Per}}^{N,p}(x, \omega)) dx. \quad (\text{A.7})$$

It is shown in [BP04, Theorem 1] that

$$\lim_{N \rightarrow \infty} A_{\operatorname{Per}}^{*,N}(\omega) = A^* \text{ almost surely.} \quad (\text{A.8})$$

### A.2.3 Neumann approximation

**Theorem A.1.** *Consider the Neumann corrector, that is*

$$\begin{cases} -\operatorname{div} A(\cdot, \omega)(p + \nabla w_{\text{Neu}}^{N,p}(\cdot, \omega)) = 0 \text{ in } Q_N, \\ A(\cdot, \omega)(p + \nabla w_{\text{Neu}}^{N,p}(\cdot, \omega)) \cdot n = p \cdot n \text{ on } \partial Q_N, \end{cases} \quad (\text{A.9})$$

the solution of which is unique in  $H^1(Q_N)$  up to the addition of a random constant. We define the random variable  $S_{\text{Neu}}^{*,N}(\omega)$  and the approximation  $A_{\text{Neu}}^{*,N}(\omega)$  of  $A^*$  by

$$S_{\text{Neu}}^{*,N}(\omega)p := \frac{1}{|Q_N|} \int_{Q_N} \left( p + \nabla w_{\text{Neu}}^{N,p}(x, \omega) \right) dx, \quad A_{\text{Neu}}^{*,N}(\omega) := S_{\text{Neu}}^{*,N}(\omega)^{-1}. \quad (\text{A.10})$$

Then,

$$\lim_{N \rightarrow \infty} A_{\text{Neu}}^{*,N}(\omega) = A^* \text{ almost surely.} \quad (\text{A.11})$$

*Proof.* There exists a solution to the problem (A.9) since the compatibility condition  $\int_{\partial Q_N} p \cdot n = 0$  is satisfied. In the sequel, we impose the condition

$$\int_{Q_N} \left( p \cdot x + w_{\text{Neu}}^{N,p}(x, \omega) \right) dx = 0 \quad (\text{A.12})$$

to enforce uniqueness. The equation (A.9) is equivalent to the variational formulation:

$$\forall \varphi \in H^1(Q_N), \quad \int_{Q_N} \nabla \varphi \cdot A(x, \omega)(p + \nabla w_{\text{Neu}}^{N,p}(x, \omega)) dx = \int_{\partial Q_N} \varphi p \cdot n.$$

To prove the convergence of  $S_{\text{Neu}}^{*,N}$ , we use the same type of arguments as in [BP04] (where there is a typo concerning which quantity converges towards  $A^*$ ).

We first rescale the equation to the domain  $Q$ . Let  $\tilde{w}_{\text{Neu}}^{N,p}(\cdot, \omega) := \frac{1}{N} w_{\text{Neu}}^{N,p}(N\cdot, \omega)$  be the rescaled corrector. It satisfies the equation

$$\begin{cases} -\operatorname{div} A_N(\cdot, \omega)(p + \nabla \tilde{w}_{\text{Neu}}^{N,p}(\cdot, \omega)) = 0 \text{ in } Q, \\ A_N(\cdot, \omega)(p + \nabla \tilde{w}_{\text{Neu}}^{N,p}(\cdot, \omega)) \cdot n = p \cdot n \text{ on } \partial Q, \end{cases} \quad (\text{A.13})$$

where  $A_N(\cdot, \omega) := A(N\cdot, \omega)$ . The variational formulation of (A.13) is given by

$$\forall \varphi \in H^1(Q), \quad \int_Q \nabla \varphi(x) \cdot A_N(x, \omega)(p + \nabla \tilde{w}_{\text{Neu}}^{N,p}(x, \omega)) dx = \int_{\partial Q} \varphi p \cdot n. \quad (\text{A.14})$$

Taking  $\varphi(x, \omega) = p \cdot x + \tilde{w}_{\text{Neu}}^{N,p}(x, \omega)$ , we obtain a  $H^1(Q)$  estimate on  $\tilde{w}_{\text{Neu}}^{N,p}$ :

$$\begin{aligned} \|p + \nabla \tilde{w}_{\text{Neu}}^{N,p}\|_{L^2(Q)}^2 &\leq \frac{1}{a_-} \int_Q (p + \nabla \tilde{w}_{\text{Neu}}^{N,p}) \cdot A_N(p + \nabla \tilde{w}_{\text{Neu}}^{N,p}) \\ &= \frac{1}{a_-} \int_{\partial Q} (p \cdot x + \tilde{w}_{\text{Neu}}^{N,p}) p \cdot n \leq \frac{C_1 C_2}{a_-} \|p \cdot x + \tilde{w}_{\text{Neu}}^{N,p}\|_{H^1(Q)}, \end{aligned}$$

where  $a_-$  is the coercivity constant of  $A$ ,  $C_1$  is continuity constant of the trace operator from  $H^1(Q)$  into  $H^{1/2}(\partial Q)$  and  $C_2 = \|p \cdot n\|_{L^2(\partial Q)}$ . We next combine this estimate with the Poincaré-Wirtinger inequality (recall (A.12)), to obtain that the function  $x \mapsto p \cdot x + \tilde{w}_{\text{Neu}}^{N,p}(x, \omega)$  is uniformly bounded in  $H^1(Q)$ .

Therefore, up to extraction, the function  $x \mapsto p \cdot x + \tilde{w}_{\text{Neu}}^{N,p}(x, \omega)$  converges weakly in  $H^1(Q)$  and almost surely to some  $v^*(\cdot, \omega) \in H^1(Q)$ . In view of [JKO94, p.14, Convergence of arbitrary solutions],  $A_N(\cdot, \omega)(p + \nabla \tilde{w}_{\text{Neu}}^{N,p}(\cdot, \omega))$  converges weakly in  $L^2(Q)$  and almost surely towards  $A^* \nabla v^*(\cdot, \omega)$ .

We are now in position to pass to the limit in (A.14). Let  $\varphi \in H^1(Q)$ ,

$$0 = \int_Q \nabla \varphi \cdot A_N(\cdot, \omega)(p + \nabla \tilde{w}_{\text{Neu}}^{N,p}(\cdot, \omega)) - \int_{\partial Q} \varphi p \cdot n \xrightarrow{N \rightarrow \infty} \int_Q \nabla \varphi \cdot A^* \nabla v^*(\cdot, \omega) - \int_{\partial Q} \varphi p \cdot n.$$

Therefore,  $v^*$  satisfies the variational formulation

$$\forall \varphi \in H^1(Q), \quad \int_Q \nabla \varphi \cdot A^* \nabla v^*(\cdot, \omega) = \int_{\partial Q} \varphi p \cdot n, \quad (\text{A.15})$$

complemented with the constraint  $\int_Q v^*(\cdot, \omega) = 0$  (this is a trivial consequence of (A.12) and weak  $L^2(Q)$  limit). This equation, whose strong formulation is

$$\begin{cases} -\operatorname{div} A^* \nabla v^* = 0 & \text{in } Q, \\ A^* \nabla v^* \cdot n = p \cdot n & \text{on } \partial Q, \end{cases} \quad (\text{A.16})$$

has a unique solution  $v^*(x) = x \cdot (A^*)^{-1} p - \int_Q x \cdot (A^*)^{-1} p$  (this solution is deterministic).

We now use the weak limit in  $H^1(Q)$ . The convergence

$$S_{\text{Neu}}^{*,N}(\omega) = \frac{1}{|Q_N|} \int_{Q_N} p + \nabla \tilde{w}_{\text{Neu}}^{N,p}(\cdot, \omega) = \int_Q p + \nabla \tilde{w}_{\text{Neu}}^{N,p}(\cdot, \omega) \xrightarrow{N \rightarrow \infty} \int_Q \nabla v^* = (A^*)^{-1} p,$$

holds almost surely. This concludes the proof. □

### A.3 Flux formulation

Assume now that  $A$  is symmetric. Then it turns out that the Neumann corrector problem (A.9) is related to a problem where the unknown is  $\sigma = A \nabla u$ . This reformulation is useful in order to establish (A.1). We first introduce the appropriate functional spaces.

#### A.3.1 The $H^{\text{div}}$ space

For any open set  $\mathcal{D}$ , we consider the space

$$H^{\text{div}}(\mathcal{D}) = \left\{ \sigma \in L^2(\mathcal{D})^d, \quad \operatorname{div} \sigma \in L^2(\mathcal{D}) \right\},$$

which is an Hilbert space for the scalar product

$$(\sigma, t)_{H^{\text{div}}(\mathcal{D})} = (\sigma, t)_{L^2(\mathcal{D})^d} + (\operatorname{div} \sigma, \operatorname{div} t)_{L^2(\mathcal{D})}.$$

Assume that the boundary of  $\mathcal{D}$  is sufficiently smooth (say e.g. Lipschitz smooth) so that the Green formula holds for  $C^1(\overline{\mathcal{D}})$  functions:

$$\forall \sigma \in C^1(\overline{\mathcal{D}})^d, \quad \forall \varphi \in C^1(\overline{\mathcal{D}}), \quad \int_{\partial \mathcal{D}} \varphi (\sigma \cdot n) = \int_{\mathcal{D}} \nabla \varphi \cdot \sigma + \int_{\mathcal{D}} \varphi \operatorname{div} \sigma.$$

it is well known that, for any  $\sigma \in H^{\text{div}}(\mathcal{D})$ , we can define the trace of  $\sigma \cdot n$  on  $\partial\mathcal{D}$  by extension of the above formula (although  $\sigma$  itself does not have a well-defined trace on  $\partial\mathcal{D}$ ). For any  $\varphi \in H^1(\mathcal{D})$ , we set

$$\langle \sigma \cdot n, \varphi \rangle = \int_{\mathcal{D}} \nabla \varphi \cdot \sigma + \int_{\mathcal{D}} \varphi \operatorname{div} \sigma. \quad (\text{A.17})$$

### A.3.2 The periodic $H_{\text{per}}^{\text{div}}$ space

We recall that  $Q = (0, 1)^d$  and  $Q_N = (0, N)^d$ . We now define the well-known space of periodic  $H^{\text{div}}$  functions as

$$H_{\text{per}}^{\text{div}}(Q_N) = \left\{ \sigma \in L_{\text{per}}^2(Q_N)^d, \quad \operatorname{div} \sigma \in L_{\text{loc}}^2(\mathbb{R}^d) \right\},$$

where we recall that

$$L_{\text{per}}^2(Q_N) = \left\{ v \in L_{\text{loc}}^2(\mathbb{R}^d), \quad v \text{ is } Q_N\text{-periodic} \right\}.$$

In  $H_{\text{per}}^{\text{div}}(Q_N)$ , we request that  $\operatorname{div} \sigma$  to be in  $L_{\text{loc}}^2(\mathbb{R}^d)$ , and not only in  $L^2(Q_N)$  (see Remark A.2 below). The space  $H_{\text{per}}^{\text{div}}(Q_N)$ , endowed with its natural scalar product  $(\cdot, \cdot)_{H^{\text{div}}(Q_N)}$  is a Hilbert space.

**Remark A.2.** *The space  $H_{\text{per}}^{\text{div}}(Q_N)$  is a distinct space from*

$$\left\{ \sigma \in L_{\text{per}}^2(Q_N)^d, \quad \operatorname{div} \sigma \in L^2(Q_N) \right\}.$$

*Indeed, in the above space, the normal trace  $\sigma \cdot n$  may jump across  $\partial Q_N$ . In contrast, in  $H_{\text{per}}^{\text{div}}(Q_N)$ , we have that  $\operatorname{div} \sigma \in L_{\text{loc}}^2(\mathbb{R}^d)$ , hence the normal trace is well-defined on  $\partial Q_N$ .*

### A.3.3 Flux formulation of the periodic problem (A.4)

For any  $\sigma \in L_{\text{loc}}^2(\mathbb{R}^d)$  and  $p \in \mathbb{R}^d$ , we consider the energy functional

$$\mathcal{E}_p(\sigma, \omega) = \frac{1}{|Q_N|} \int_{Q_N} (\sigma + p) \cdot A^{-1}(\cdot, \omega)(\sigma + p). \quad (\text{A.18})$$

We also consider the variational problem

$$\mathcal{S}_{\text{Per}}^{\star, N}(p, \omega) = \inf \left\{ \mathcal{E}_p(\sigma, \omega), \quad \sigma \in H_{\text{per}}^{\text{div}}(Q_N), \quad \int_{Q_N} \sigma = 0, \quad \operatorname{div} \sigma = 0 \text{ in } \mathbb{R}^d \right\}. \quad (\text{A.19})$$

Let

$$\mathcal{W} = \left\{ \sigma \in H_{\text{per}}^{\text{div}}(Q_N), \quad \int_{Q_N} \sigma = 0, \quad \operatorname{div} \sigma = 0 \text{ in } \mathbb{R}^d \right\}.$$

The following result is proved in [JKO94, Equation (1.64)]. We provide its proof for consistency.

**Lemma A.3.** *The variational problem (A.19) is well posed and there exists a matrix  $S_{\text{Per}}^{\star, N}(\omega)$  such that  $\mathcal{S}_{\text{Per}}^{\star}(p, \omega) = p \cdot S_{\text{Per}}^{\star, N}(\omega)p$ . Additionally, assume that  $A$  is symmetric. Then the Euler-Lagrange equations of problem (A.19) read*

$$\left\{ \begin{array}{l} \text{Find } \sigma \in \mathcal{W} \text{ such that} \\ \text{for any } h \in \mathcal{W}, \text{ we have } \int_{Q_N} h \cdot A^{-1}(\cdot, \omega)(p + \sigma) = 0. \end{array} \right. \quad (\text{A.20})$$

*Finally, it holds that  $S_{\text{Per}}^{\star, N}(\omega) = A_{\text{Per}}^{\star, N}(\omega)^{-1}$ .*



*Proof.* The space  $\mathcal{W}$  is a Hilbert space for the same norm as  $H_{\text{per}}^{\text{div}}(Q_N)$ . On  $\mathcal{W}$ , the norm is simply the  $L^2(Q_N)$  norm. Problem (A.19) consists in minimizing a strongly convex quadratic functional over  $\mathcal{W}$ . The functional is coercive according to the  $L^2(Q_N)$  norm. We can thus apply the Lax-Milgram theorem, which implies that the problem is well-posed. From the linearity of the Euler-Lagrange equation, we deduce that  $\mathcal{S}_{\text{Per}}^{*,N}(p, \omega) = p \cdot \mathcal{S}_{\text{Per}}^{*,N}(\omega)p$  for some matrix  $S_{\text{Per}}^{*,N}(\omega)$ .

It remains to show that  $S_{\text{Per}}^{*,N}(\omega) = A_{\text{Per}}^{*,N}(\omega)^{-1}$ . Let  $p \in \mathbb{R}^d$  and set  $\xi(\omega) = A_{\text{Per}}^{*,N}(\omega)^{-1}p$ . We denote by  $w_{\text{Per}}^{N,\xi}$  the periodic corrector function (defined up to the addition of a random constant) in the direction  $\xi$ , namely the solution to (A.4), that is

$$-\text{div} \left[ \tilde{A}_{\text{Per}}^N \left( \nabla w_{\text{Per}}^{N,\xi} + \xi \right) \right] = 0 \text{ in } \mathbb{R}^d, \quad w_{\text{Per}}^{N,\xi}(\cdot, \omega) \text{ is } Q_N\text{-periodic,}$$

where we recall that  $\tilde{A}_{\text{Per}}^N$  is the  $Q_N$  periodic extension of  $A$ .

We set  $\sigma_\xi = \tilde{A}_{\text{Per}}^N \left( \xi + \nabla w_{\text{Per}}^{N,\xi} \right) - p$  and claim that  $\sigma_\xi(\cdot, \omega) \in \mathcal{W}$ . We indeed have that  $\text{div} \sigma_\xi = 0$  in  $\mathbb{R}^d$ , while

$$\frac{1}{|Q_N|} \int_{Q_N} \sigma_\xi(x, \omega) dx = A_{\text{Per}}^{*,N}(\omega) \xi(\omega) - p = 0.$$

Furthermore,  $\sigma_\xi(\cdot, \omega) \in L_{\text{loc}}^2(\mathbb{R}^d)$  and is  $Q_N$  periodic, while  $\text{div} \sigma_\xi$  is in  $L_{\text{loc}}^2(\mathbb{R}^d)$ . We hence indeed have that  $\sigma_\xi(\cdot, \omega) \in \mathcal{W}$ .

We next show that  $\sigma_\xi$  satisfies the Euler-Lagrange equation (A.20). Let  $h \in \mathcal{W}$ . Then,

$$\int_{Q_N} h \cdot A^{-1}(p + \sigma_\xi) = \int_{Q_N} h \cdot \left( \xi + \nabla w_{\text{Per}}^{N,\xi} \right) = \int_{Q_N} h \cdot \nabla w_{\text{Per}}^{N,\xi},$$

using that the mean of  $h$  vanishes. We next integrate by part and use that  $\text{div} h = 0$  on  $\mathbb{R}^d$ , yielding

$$\int_{Q_N} h \cdot A^{-1}(p + \sigma_\xi) = \int_{\partial Q_N} (h \cdot n) w_{\text{Per}}^{N,\xi}.$$

We now observe that  $w_{\text{Per}}^{N,\xi}$  has a well-defined trace on  $\partial Q_N$  (since  $w_{\text{Per}}^{N,\xi} \in H_{\text{loc}}^1(\mathbb{R}^d)$ ). Likewise,  $h \cdot n$  has also a well-defined trace since  $h \in H_{\text{per}}^{\text{div}}(Q_N)$ . Since  $h$  and  $w_{\text{Per}}^{N,\xi}$  are periodic, the above right-hand side vanishes. Thus,  $\sigma_\xi$  satisfies the Euler-Lagrange equation (A.20), and is hence the unique minimizer of (A.19). As a consequence,

$$\begin{aligned} p \cdot S_{\text{Per}}^{*,N} p &= \frac{1}{|Q_N|} \int_{Q_N} (\sigma_\xi + p) \cdot A^{-1}(\sigma_\xi + p) \\ &= \frac{1}{|Q_N|} \int_{Q_N} (\xi + \nabla w_{\text{Per}}^{N,\xi}) \cdot A(\xi + \nabla w_{\text{Per}}^{N,\xi}) = \xi \cdot A_{\text{Per}}^{*,N} \xi = p \cdot (A_{\text{Per}}^{*,N})^{-1} p. \end{aligned}$$

Since  $p$  is arbitrary, this concludes the proof.  $\square$

### A.3.4 Flux formulation of the Neumann problem (A.9)

For any vector  $p \in \mathbb{R}^d$ , we define

$$\mathcal{S}_{\text{Neu}}^{*,N}(p, \omega) = \inf \left\{ \mathcal{E}_p(\sigma, \omega), \sigma \in H^{\text{div}}(Q_N), \sigma \cdot n = 0 \text{ on } \partial Q_N, \text{div } \sigma = 0 \text{ in } Q_N \right\}, \tag{A.21}$$

where  $H^{\text{div}}(Q_N)$  is defined in Section A.3.1 and  $\mathcal{E}_p(\sigma, \omega)$  is defined by (A.18).

**Lemma A.4.** *Problem (A.21) is well-posed.*

Additionally, assume that  $A$  is symmetric. Then, the solution  $\sigma_N$  of this problem is given by  $\sigma_N = A(p + \nabla w_{\text{Neu}}^{N,p}) - p$ , where  $w_{\text{Neu}}^{N,p}$  is the solution to (A.9). In addition, it holds that  $\mathcal{S}_{\text{Neu}}^{\star,N}(p, \omega) = p \cdot S_{\text{Neu}}^{\star,N}(\omega)p$ , where  $S_{\text{Neu}}^{\star,N}(\omega)$  is defined by (A.10).

*Proof.* The space

$$\mathcal{V} = \left\{ \sigma \in H^{\text{div}}(Q_N), \sigma \cdot n = 0 \text{ on } \partial Q_N, \text{div } \sigma = 0 \text{ in } Q_N \right\}$$

is closed subspace of  $H^{\text{div}}(Q_N)$  and therefore a Hilbert space. In  $\mathcal{V}$ , the  $H^{\text{div}}$  norm reduces to the  $L^2$  norm. The functional  $\mathcal{E}_p(\cdot, \omega)$  is hence coercive on  $\mathcal{V}$ . It is also continuous and strongly convex, therefore the problem (A.21) is well posed. Let  $\sigma_N^p(\cdot, \omega) \in \mathcal{V}$  be its unique solution.

In the symmetric case, the Euler-Lagrange equation of (A.21) reads

$$\forall \psi \in \mathcal{V}, \int_{Q_N} \psi \cdot A^{-1}(\sigma_N^p + p) = 0. \quad (\text{A.22})$$

We next show that  $\sigma_N^p = A(p + \nabla w_{\text{Neu}}^{N,p}) - p$ . Let  $\psi \in \mathcal{V}$ . We see that

$$\int_{Q_N} \psi \cdot A^{-1}(A(p + \nabla w_{\text{Neu}}^{N,p}) - p + p) = \int_{Q_N} \psi \cdot \nabla(p \cdot x + w_{\text{Neu}}^{N,p}) = 0,$$

where the latter equality is obtained by integration by parts. Thus,  $A(p + \nabla w_{\text{Neu}}^{N,p}) - p$  satisfies the Euler-Lagrange equations associated to problem (A.21). Therefore it is the solution. In addition,

$$\begin{aligned} \mathcal{S}_{\text{Neu}}^{\star,N}(p, \omega) &= \frac{1}{|Q_N|} \int_{Q_N} \left( (\sigma_N^p + p) \cdot A^{-1} A(p + \nabla w_{\text{Neu}}^{N,p}) \right) (x, \omega) dx \\ &= \frac{1}{|Q_N|} \int_{Q_N} p \cdot (p + \nabla w_{\text{Neu}}^{N,p}(\cdot, \omega)), \end{aligned}$$

because  $\int_{Q_N} \sigma_N^p \cdot (p + \nabla w_{\text{Neu}}^{N,p}) = 0$  by integration by parts. By definition of  $\mathcal{S}_{\text{Neu}}^{\star,N}(\omega)$  (see (A.22)) we have  $\mathcal{S}_{\text{Neu}}^{\star,N}(p, \omega) = p \cdot S_{\text{Neu}}^{\star,N}(\omega)p$ . This concludes the proof.  $\square$

## A.4 Comparison of the approximations

We now establish (A.1).

### A.4.1 Comparison of the Dirichlet and periodic approximations

We have the following result:

**Theorem A.5.** *We assume that  $A$  is a symmetric matrix. Then  $A_{\text{Dir}}^{\star,N}(\omega)$  and  $A_{\text{Per}}^{\star,N}(\omega)$  are symmetric matrices and, in the sense of symmetric matrices, we have, for any  $N$ ,*

$$A_{\text{Per}}^{\star,N}(\omega) \leq A_{\text{Dir}}^{\star,N}(\omega) \text{ almost surely.} \quad (\text{A.23})$$

*Proof.* The symmetry of  $A_{\text{Dir}}^{*,N}(\omega)$  and  $A_{\text{Per}}^{*,N}(\omega)$  is a classical result. We now prove (A.23). Let  $p \in \mathbb{R}^d$ . Using the symmetry of  $A$ , we see that

$$p \cdot A_{\text{Dir}}^{*,N}(\omega)p = \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} (p + \nabla v) \cdot A(\cdot, \omega)(p + \nabla v), \quad v \in H_0^1(Q_N) \right\}$$

while

$$p \cdot A_{\text{Per}}^{*,N}(\omega)p = \inf \left\{ \frac{1}{|Q_N|} \int_{Q_N} (p + \nabla v) \cdot A(\cdot, \omega)(p + \nabla v), \quad v \in V_{\text{per}}(Q_N) \right\}, \quad (\text{A.24})$$

where  $V_{\text{per}}(Q_N)$  is defined by (A.6). For any  $v \in H^1(Q_N)$ , let

$$J(v, \omega) = \frac{1}{|Q_N|} \int_{Q_N} (p + \nabla v) \cdot A(\cdot, \omega)(p + \nabla v).$$

Let  $v \in H_0^1(Q_N)$ . Consider the  $Q_N$ -periodic extension  $\tilde{v}$  of  $v$ , and set

$$\bar{v} = \tilde{v} - \frac{1}{|Q_N|} \int_{Q_N} \tilde{v}.$$

By construction, we see that  $\bar{v} \in H_{\text{loc}}^1(\mathbb{R}^d)$  and that  $\bar{v}$  is  $Q_N$  periodic. In addition, the mean of  $\bar{v}$  on  $Q_N$  vanishes. Hence  $\bar{v} \in V_{\text{per}}(Q_N)$  and  $J(\bar{v}) = J(v)$ . Thus

$$p \cdot A_{\text{Per}}^{*,N}(\omega)p \leq J(\bar{v}, \omega) = J(v, \omega).$$

Taking the infimum over  $v \in H_0^1(Q_N)$ , we deduce that

$$p \cdot A_{\text{Per}}^{*,N}(\omega)p \leq p \cdot A_{\text{Dir}}^{*,N}(\omega)p$$

which is valid for any  $p$ . This concludes the proof.  $\square$

#### A.4.2 Comparison of the Neumann and periodic approximations

We have the following result:

**Theorem A.6.** *We assume that  $A$  is a symmetric matrix. Then  $A_{\text{Neu}}^{*,N}(\omega)$  is a symmetric matrix and, in the sense of symmetric matrices, we have, for any  $N$ ,*

$$A_{\text{Neu}}^{*,N}(\omega) \leq A_{\text{Per}}^{*,N}(\omega) \text{ almost surely.} \quad (\text{A.25})$$

*Proof.* We deduce the properties on  $A^{*,N}$  from the properties of  $S^{*,N}$  that we now establish.

The symmetry of  $S_{\text{Neu}}^{*,N}(\omega)$  relies on the equality  $\sigma_N^p = A(p + \nabla w_{\text{Neu}}^{N,p}) - p$  from Lemma A.4. Indeed, for all  $\xi, p \in \mathbb{R}^d$ ,

$$\xi \cdot S_{\text{Neu}}^{*,N}p = \frac{1}{|Q_N|} \int_{Q_N} \xi \cdot (p + \nabla w_{\text{Neu}}^{N,p}) = \frac{1}{|Q_N|} \int_{Q_N} \xi \cdot A^{-1}(p + \sigma_N^p).$$

In view of the Euler-Lagrange equations of (A.21), namely (A.22), we have  $\int_{Q_N} \sigma_N^\xi \cdot A^{-1}(p + \sigma_N^p) = 0$ . Therefore,

$$\xi \cdot S_{\text{Neu}}^{*,N}p = \frac{1}{N^d} \int_{Q_N} (\xi + \sigma_N^\xi) \cdot A^{-1}(p + \sigma_N^p),$$

and thus  $S_{\text{Neu}}^{\star,N}(\omega)$  is symmetric.

We now prove the bound  $S_{\text{Per}}^{\star,N}(\omega) \leq S_{\text{Neu}}^{\star,N}(\omega)$  almost surely. Let  $p \in \mathbb{R}^d$ . We have that

$$p \cdot S_{\text{Neu}}^{\star,N}(\omega)p = \inf \{ \mathcal{E}_p(\sigma, \omega), \sigma \in \mathcal{V} \}$$

while

$$p \cdot S_{\text{Per}}^{\star,N}(\omega)p = \inf \{ \mathcal{E}_p(\sigma, \omega), \sigma \in \mathcal{W} \}.$$

Let  $\sigma \in \mathcal{V}$ . Consider the  $Q_N$ -periodic extension  $\bar{\sigma}$  of  $\sigma$ . We already have that  $\bar{\sigma} \in L_{\text{per}}^2(Q_N)$ . Let us show that  $\text{div } \bar{\sigma} = 0$  in  $\mathcal{D}'(\mathbb{R}^d)$ . Let  $\varphi \in \mathcal{D}(\mathbb{R}^d)$ . Then

$$\langle \text{div } \bar{\sigma}, \varphi \rangle = -\langle \bar{\sigma}, \nabla \varphi \rangle = \sum_{k \in \mathbb{Z}^d} \int_{Q_N + Nk} \bar{\sigma} \cdot \nabla \varphi = \sum_{k \in \mathbb{Z}^d} \int_{Q_N} \bar{\sigma} \cdot \nabla \varphi(\cdot + Nk)$$

where the sum in  $k$  is actually finite. Let  $\psi = \sum_{k \in \mathbb{Z}^d} \varphi(\cdot + Nk)$ . We thus have

$$\langle \text{div } \bar{\sigma}, \varphi \rangle = \int_{Q_N} \sigma \cdot \nabla \psi = - \int_{Q_N} \psi \text{div } \sigma + \langle \sigma \cdot n, \psi \rangle$$

and both terms vanish because  $\text{div } \sigma = 0$  on  $Q_N$  and  $\sigma \cdot n = 0$  on  $\partial Q_N$  for any  $\sigma \in \mathcal{V}$ .

We hence have that  $\bar{\sigma} \in H_{\text{per}}^{\text{div}}(Q_N)$  and  $\text{div } \bar{\sigma} = 0$  in  $\mathcal{D}'(\mathbb{R}^d)$ . We eventually compute that

$$e_i \cdot \int_{Q_N} \bar{\sigma} = e_i \cdot \int_{Q_N} \sigma = \int_{\partial Q_N} x_i \sigma \cdot n - \int_{Q_N} x_i \text{div } \sigma = 0,$$

using again that  $\text{div } \sigma = 0$  on  $Q_N$  and  $\sigma \cdot n = 0$  on  $\partial Q_N$ . Hence  $\bar{\sigma} \in \mathcal{W}$ , and  $\mathcal{E}_p(\bar{\sigma}, \omega) = \mathcal{E}_p(\sigma, \omega)$ . Thus

$$p \cdot S_{\text{Per}}^{\star,N}(\omega)p \leq \mathcal{E}_p(\bar{\sigma}, \omega) = \mathcal{E}_p(\sigma, \omega).$$

Taking the infimum over  $\sigma \in \mathcal{V}$ , we deduce that

$$p \cdot S_{\text{Per}}^{\star,N}(\omega)p \leq p \cdot S_{\text{Neu}}^{\star,N}(\omega)p$$

which is valid for any  $p$ . This concludes the proof of (A.24).  $\square$



# Bibliography

- [ACLB<sup>+</sup>12] A. Anantharaman, R. Costaouec, C. Le Bris, F. Legoll, and F. Thomines, *Introduction to numerical stochastic homogenization and the related computational challenges: some recent developments*, Multiscale modeling and analysis for materials simulation, Lect. Notes Ser. Inst. Math. Sci. Natl. Univ. Singap., vol. 22, World Sci. Publ., Hackensack, NJ, 2012, pp. 197–272. MR 2895600
- [ACS14] S. N. Armstrong, P. Cardaliaguet, and P. E. Souganidis, *Error estimates and convergence rates for the stochastic homogenization of Hamilton-Jacobi equations*, J. Amer. Math. Soc. **27** (2014), no. 2, 479–540. MR 3164987
- [ALB10] A. Anantharaman and C. Le Bris, *Homogénéisation d'un matériau périodique faiblement perturbé aléatoirement*, C. R. Math. Acad. Sci. Paris **348** (2010), no. 9-10, 529–534. MR 2645167 (2011g:35023)
- [ALB11] ———, *A numerical approach related to defect-type theories for some weakly random problems in homogenization*, Multiscale Modeling & Simulation **9** (2011), no. 2, 513–544.
- [ALB12] ———, *Elements of mathematical foundations for numerical approaches for weakly random homogenization problems*, Commun. Comput. Phys. **11** (2012), no. 4, 1103–1143. MR 2864078 (2012k:35028)
- [AV12] A. Abdulle and G. Vilmart, *A priori error estimates for finite element methods with numerical quadrature for nonmonotone nonlinear elliptic problems*, Numer. Math. **121** (2012), no. 3, 397–431. MR 2929073
- [Bal76] John M Ball, *Convexity conditions and existence theorems in nonlinear elasticity*, Archive for rational mechanics and Analysis **63** (1976), no. 4, 337–403.
- [BCLBL12a] X. Blanc, R. Costaouec, C. Le Bris, and F. Legoll, *Variance reduction in stochastic homogenization: the technique of antithetic variables*, Numerical Analysis of Multiscale Computations, Springer, 2012, pp. 47–70.
- [BCLBL12b] ———, *Variance reduction in stochastic homogenization using antithetic variables*, Markov Processes and Related Fields **18** (2012), no. 1, 31–66.
- [Bis11] M. Biskup, *Recent progress on the random conductance model*, Probab. Surv. **8** (2011), 294–373. MR 2861133
- [BL] M. Bornert and F. Legoll, *in prep.*

- [BL93] J. W. Barrett and W. B. Liu, *Finite element approximation of the  $p$ -Laplacian*, Math. Comp. **61** (1993), no. 204, 523–537. MR 1192966 (94c:65129)
- [BLBL06] X. Blanc, C. Le Bris, and P.-L. Lions, *Une variante de la théorie de l'homogénéisation stochastique des opérateurs elliptiques*, C. R. Math. Acad. Sci. Paris **343** (2006), no. 11-12, 717–724. MR 2284699 (2009c:60195)
- [BLBL07] ———, *Stochastic homogenization and random lattices*, Journal de mathématiques pures et appliquées **88** (2007), no. 1, 34–63.
- [BLP78] A. Bensoussan, J.-L. Lions, and G. Papanicolaou, *Asymptotic methods in periodic structures*, Studies in Math. Appl **5** (1978).
- [BO14] C. Bernardin and S. Olla, *Thermodynamics and non-equilibrium macroscopic dynamics of chains of anharmonic oscillators*, Lecture Notes available at <https://www.ceremade.dauphine.fr/~olla> (2014).
- [BP04] A. Bourgeat and A. Piatnitski, *Approximations of effective coefficients in stochastic homogenization*, Annales de l'Institut Henri Poincaré (B) Probability and Statistics, vol. 40, Elsevier, 2004, pp. 153–165.
- [CD99] D. Cioranescu and P. Donato, *An introduction to homogenization*, Oxford Lecture Series in Mathematics and its Applications, vol. 17, The Clarendon Press, Oxford University Press, New York, 1999. MR 1765047 (2001j:35019)
- [Cho89] S.-S. Chow, *Finite element error estimates for nonlinear elliptic equations of monotone type*, Numer. Math. **54** (1989), no. 4, 373–393. MR 972416 (90a:65235)
- [CLBL10] R. Costeaouec, C. Le Bris, and F. Legoll, *Variance reduction in stochastic homogenization: proof of concept, using antithetic variables*, SeMA Journal **50** (2010), no. 1, 9–26.
- [Cos12] R. Costeaouec, *Asymptotic expansion of the homogenized matrix in two weakly stochastic homogenization settings*, Applied Mathematics Research eXpress **2012** (2012), no. 1, 76–104.
- [DMM86a] G. Dal Maso and L. Modica, *Nonlinear stochastic homogenization*, Ann. Mat. Pura Appl. (4) **144** (1986), 347–389. MR 870884 (88h:49025)
- [DMM86b] ———, *Nonlinear stochastic homogenization and ergodic theory*, J. Reine Angew. Math. **368** (1986), 28–42. MR 850613 (88k:28021)
- [ES08] B. Engquist and P. E. Souganidis, *Asymptotic and numerical homogenization*, Acta Numerica **17** (2008), 147–190.
- [Fat56] I. Fatt, *The network model of porous media.*, Petrol. Trans. AIME **207** (1956), 144–159.
- [Fis96] G. S. Fishman, *Monte Carlo*, Springer Series in Operations Research, Springer-Verlag, New York, 1996, Concepts, algorithms, and applications. MR 1392474 (97g:65019)

- [GM75] R. Glowinski and A. Marroco, *Sur l'approximation, par éléments finis d'ordre un, et la résolution, par pénalisation-dualité d'une classe de problèmes de dirichlet non linéaires*, ESAIM: Mathematical Modelling and Numerical Analysis-Modélisation Mathématique et Analyse Numérique **9** (1975), no. R2, 41–76.
- [GN11] A. Gloria and S. Neukamm, *Commutability of homogenization and linearization at identity in finite elasticity and applications*, Ann. Inst. H. Poincaré Anal. Non Linéaire **28** (2011), no. 6, 941–964. MR 2859933 (2012j:35017)
- [GNO14] A. Gloria, S. Neukamm, and F. Otto, *An optimal quantitative two-scale expansion in stochastic homogenization of discrete elliptic equations*, ESAIM Math. Model. Numer. Anal. **48** (2014), no. 2, 325–346. MR 3177848
- [GNO15] ———, *Quantification of ergodicity in stochastic homogenization: optimal bounds via spectral gap on Glauber dynamics*, Invent. Math. **199** (2015), no. 2, 455–515. MR 3302119
- [GO12] A. Gloria and F. Otto, *An optimal error estimate in stochastic homogenization of discrete elliptic equations*, Ann. Appl. Probab. **22** (2012), no. 1, 1–28. MR 2932541
- [GO15] ———, *Quantitative estimates on the periodic approximation of the corrector in stochastic homogenization*, ESAIM: Proc. **48** (2015), 80–97.
- [JKO94] V. V. Jikov, S. M. Kozlov, and O. A. Oleĭnik, *Homogenization of differential operators and integral functionals*, Springer-Verlag, Berlin, 1994. MR 1329546 (96h:35003b)
- [KFG<sup>+</sup>03] T. Kanit, S. Forest, I. Galliet, V. Mounoury, and D. Jeulin, *Determination of the size of the representative volume element for random composites: statistical and numerical approach*, International Journal of Solids and Structures **40** (2003), no. 13, 3647–3679.
- [Koz79] S. M. Kozlov, *Averaging of random operators*, Matematicheskii Sbornik **151** (1979), no. 2, 188–202.
- [Koz87] ———, *Averaging of difference schemes*, Mathematics of the USSR-Sbornik **57** (1987), no. 2, 351.
- [Kre85] U. Krengel, *Ergodic theorems*, de Gruyter Studies in Mathematics, vol. 6, Walter de Gruyter & Co., Berlin, 1985, With a supplement by Antoine Brunel. MR 797411 (87i:28001)
- [Kün83] R. Künnemann, *The diffusion limit for reversible jump processes on  $d$  with ergodic random bond conductivities*, Communications in Mathematical Physics **90** (1983), no. 1, 27–68.
- [LB10] C. Le Bris, *Some numerical approaches for weakly random homogenization*, Numerical Mathematics and Advanced Applications 2009, Springer, 2010, pp. 29–45.
- [LBLM] C. Le Bris, F. Legoll, and W. Minvielle, *Special quasirandom structures: a selection approach for stochastic homogenization*.



- [LBT12] C. Le Bris and F. Thomines, *A reduced basis approach for some weakly stochastic multiscale problems*, Chinese Annals of Mathematics, Series B **33** (2012), no. 5, 657–672.
- [Liu08] J. S. Liu, *Monte Carlo strategies in scientific computing*, Springer Series in Statistics, Springer, New York, 2008. MR 2401592 (2010b:65013)
- [LM13] F. Legoll and W. Minvielle, *Variance reduction using antithetic variables for a nonlinear convex stochastic homogenization problem*, arXiv preprint arXiv:1302.0038 (2013).
- [LM15a] ———, *A control variate approach based on a defect-type theory for variance reduction in stochastic homogenization*, Multiscale Modeling & Simulation **13** (2015), no. 2, 519–550.
- [LM15b] ———, *Variance reduction using antithetic variables for a nonlinear convex stochastic homogenization problem*, Discrete and Continuous Dynamical Systems - Series S **8** (2015), no. 1, 1–27.
- [LMOS] F. Legoll, W. Minvielle, A. Obliger, and M. Simon, *In prep.*
- [LMOS15] ———, *A parameter identification problem in stochastic homogenization*, ESAIM: Proc. **48** (2015), 190–214.
- [LT94] P. Le Tallec, *Numerical methods for nonlinear three-dimensional elasticity*, Handbook of numerical analysis, Vol. III, Handb. Numer. Anal., III, North-Holland, Amsterdam, 1994, pp. 465–622. MR 1307410 (96b:73093)
- [Min15] W. Minvielle, *Thèse de l'Université Paris Est.*
- [Mou15] J.-C. Mourrat, *First-order expansion of homogenized coefficients under bernoulli perturbations*, Journal de Mathématiques Pures et Appliquées **103** (2015), no. 1, 68–101.
- [Nol14] J. Nolen, *Normal approximation for a random elliptic equation*, Probability Theory and Related Fields **159** (2014), no. 3-4, 661–700.
- [NPS12] J. Nolen, G. A. Pavliotis, and A. M. Stuart, *Multiscale modelling and inverse problems*, Numerical Analysis of Multiscale Problems, Springer, 2012, pp. 1–34.
- [PV81] G. C. Papanicolaou and S. R. S. Varadhan, *Boundary value problems with rapidly oscillating random coefficients*, Random fields, Vol. I, II (Esztergom, 1979), Colloq. Math. Soc. János Bolyai, vol. 27, North-Holland, Amsterdam-New York, 1981, pp. 835–873. MR 712714 (84k:58233)
- [Shi84] A. N. Shiriyayev, *Probability*, Graduate Texts in Mathematics, vol. 95, Springer-Verlag, New York, 1984. MR 737192 (85a:60007)
- [Tar97] L. Tartar, *Estimations of homogenized coefficients [MR0540123 (80i:35010)]*, Topics in the mathematical modelling of composite materials, Progr. Non-linear Differential Equations Appl., vol. 31, Birkhäuser Boston, Boston, MA, 1997, pp. 9–20. MR 1493038

- [Tem72] A. A. Tempel'man, *Ergodic theorems for general dynamical systems*, Trudy Moskov. Mat. Obšč. **26** (1972), 95–132. MR 0374388 (51 #10588)
- [Tho97] V. Thomée, *Galerkin finite element methods for parabolic problems*, Springer Series in Computational Mathematics, vol. 25, Springer-Verlag, Berlin, 1997. MR 1479170 (98m:65007)
- [vPDFN10] J. von Pezold, A. Dick, M. Friák, and J. Neugebauer, *Generation and performance of special quasirandom structures for studying the elastic properties of random alloys: Application to Al-Ti*, Physical Review B **81** (2010), no. 9, 094203.
- [WFBZ90] S.-H. Wei, L. G. Ferreira, J. E. Bernard, and A. Zunger, *Electronic properties of random alloys: Special quasirandom structures*, Physical Review B **42** (1990), no. 15, 9622.
- [Yur86] V. V. Yurinskii, *Averaging of symmetric diffusion in random medium*, Siberian Mathematical Journal **27** (1986), no. 4, 603–613.
- [ZWFB90] A. Zunger, S.-H. Wei, L. G. Ferreira, and J. E. Bernard, *Special quasirandom structures*, Physical Review Letters **65** (1990), no. 3, 353.