

# Traitement de la confusion spatiale par R-INLA dans le cas de données géostatistiques

Jérémy Lamouroux, INRAE  
Début de thèse : 01/02/2022  
École doctorale : EDMH

Encadrants : I. Albert, UMR MIA Paris-Saclay, département MathNum INRAE ;  
S. Leblond, C. Meyer, UMS PatriNat (OFB-CNRS-MNHN)

12 Décembre 2024



AgroParisTech

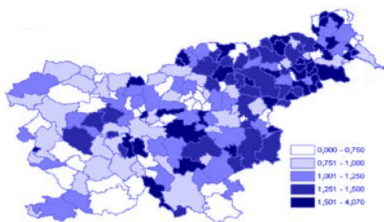


# Plan de la présentation

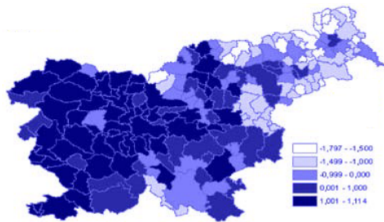
- 1 Contexte et objectifs
- 2 Méthodes
- 3 Résultats
- 4 Discussion et perspectives

# Contexte : Un cas de confusion spatiale

*Hodges and Reich, The American Statistician (2010)*



(a) Incidence de la maladie



(b) Indice socio-économique

## Modèle Null

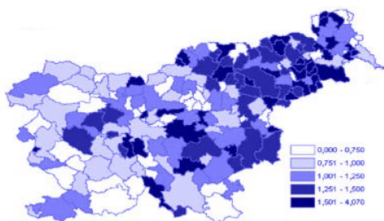
$$\underbrace{Y(s)}_{\text{Incidence de la maladie}} = \alpha + \beta_{Null} \underbrace{X(s)}_{\text{Indice socio-économique}} + \underbrace{\epsilon(s)}_{\text{Bruit aléatoire}} \rightarrow \text{Effet négatif } (\hat{\beta}_{Null} < 0) \text{ très significatif}$$

## Modèle spatial

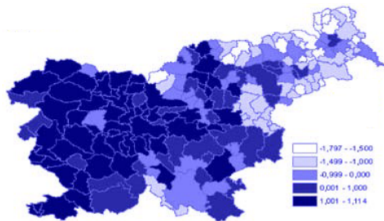
$$Y(s) = \alpha + \beta_{Spatial} X(s) + \underbrace{U(s)}_{\text{Effet spatial aléatoire}} + \epsilon(s) \rightarrow \hat{\beta}_{Spatial} \neq \hat{\beta}_{Null} \text{ et non significatif}$$

# Contexte : Un cas de confusion spatiale

*Hodges and Reich, The American Statistician (2010)*



(a) Incidence de la maladie



(b) Indice socio-économique

## Modèle Null

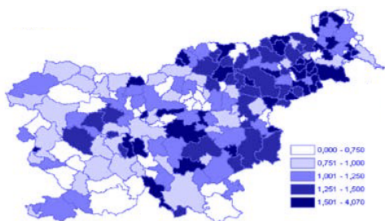
$$\underbrace{Y(s)}_{\text{Incidence de la maladie}} = \alpha + \beta_{Null} \underbrace{X(s)}_{\text{Indice socio-économique}} + \underbrace{\epsilon(s)}_{\text{Bruit aléatoire}} \rightarrow \text{Effet négatif } (\hat{\beta}_{Null} < 0) \text{ très significatif}$$

## Modèle spatial

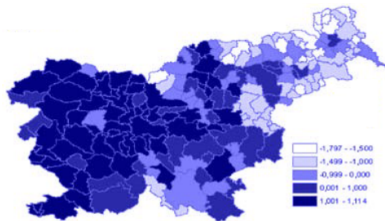
$$Y(s) = \alpha + \beta_{Spatial} X(s) + \underbrace{U(s)}_{\text{Effet spatial aléatoire}} + \epsilon(s) \rightarrow \hat{\beta}_{Spatial} \neq \hat{\beta}_{Null} \text{ et non significatif}$$

# Contexte : Un cas de confusion spatiale

*Hodges and Reich, The American Statistician (2010)*



(a) Incidence de la maladie



(b) Indice socio-économique

## Modèle Null

$$\underbrace{Y(s)}_{\text{Incidence de la maladie}} = \alpha + \beta_{Null} \underbrace{X(s)}_{\text{Indice socio-économique}} + \underbrace{\epsilon(s)}_{\text{Bruit aléatoire}} \rightarrow \text{Effet négatif } (\hat{\beta}_{Null} < 0) \text{ très significatif}$$

## Modèle spatial

$$Y(s) = \alpha + \beta_{Spatial} X(s) + \underbrace{U(s)}_{\text{Effet spatial aléatoire}} + \epsilon(s) \rightarrow \hat{\beta}_{Spatial} \neq \hat{\beta}_{Null} \text{ et non significatif}$$

## Les quatre formes de confusion spatiale

- 1 **Biais dû à un facteur confondant omis** : Présence d'une variable non mesurée avec une structure spatiale influençant à la fois la variable et la covariable (Schnell et Papadogeorgou (2020)).
  - exemple : Si une étude analyse la relation entre la pollution de l'air (variable) et des maladies respiratoires (covariable), mais néglige la proximité des zones industrielles (confondeur spatial), les résultats peuvent être faussés.
- 2 **Colinéarité des effets aléatoires** : Modification des estimations des effets fixes due à l'ajout d'effets aléatoires spatiaux colinéaires avec les covariables (Hodges et Reich (2010)).
  - exemple: Dans une régression qui inclut à la fois des variables explicatives comme le revenu et des effets aléatoires spatiaux, si le revenu varie fortement selon des régions spécifiques, les deux composants risquent d'être colinéaires.
- 3 **Biais de régularisation** : Biais dans des échantillons finis associé à l'utilisation de fonctions de régression flexibles comme les splines ou les processus gaussiens (GP) pour modéliser une fonction inconnue de l'espace (Dupont et al. (2020)).
  - exemple : Une spline utilisée pour modéliser l'effet de la distance géographique peut sous-estimer les variations locales importantes en raison d'un lissage excessif.
- 4 **Concurvité** : Difficulté à évaluer l'effet d'une variable lorsque celle-ci est une fonction lisse (ou proche d'une fonction lisse) de l'espace, lorsque le modèle inclut également une fonction lisse arbitraire de l'espace (Ramsay et al. (2003), Paciorek (2010)).
  - exemple : Étudier l'effet de la température (qui varie de manière lisse dans l'espace) tout en incluant une fonction spatiale arbitraire dans le modèle peut rendre impossible l'identification claire de l'effet de la température.

## Les quatre formes de confusion spatiale

- 1 **Biais dû à un facteur confondant omis** : Présence d'une variable non mesurée avec une structure spatiale influençant à la fois la variable et la covariable (Schnell et Papadogeorgou (2020)).
  - exemple : Si une étude analyse la relation entre la pollution de l'air (variable) et des maladies respiratoires (covariable), mais néglige la proximité des zones industrielles (confondeur spatial), les résultats peuvent être faussés.
- 2 **Colinéarité des effets aléatoires** : Modification des estimations des effets fixes due à l'ajout d'effets aléatoires spatiaux colinéaires avec les covariables (Hodges et Reich (2010)).
  - exemple: Dans une régression qui inclut à la fois des variables explicatives comme le revenu et des effets aléatoires spatiaux, si le revenu varie fortement selon des régions spécifiques, les deux composants risquent d'être colinéaires.
- 3 **Biais de régularisation** : Biais dans des échantillons finis associé à l'utilisation de fonctions de régression flexibles comme les splines ou les processus gaussiens (GP) pour modéliser une fonction inconnue de l'espace (Dupont et al. (2020)).
  - exemple : Une spline utilisée pour modéliser l'effet de la distance géographique peut sous-estimer les variations locales importantes en raison d'un lissage excessif.
- 4 **Concurvité** : Difficulté à évaluer l'effet d'une variable lorsque celle-ci est une fonction lisse (ou proche d'une fonction lisse) de l'espace, lorsque le modèle inclut également une fonction lisse arbitraire de l'espace (Ramsay et al. (2003), Paciorek (2010)).
  - exemple : Étudier l'effet de la température (qui varie de manière lisse dans l'espace) tout en incluant une fonction spatiale arbitraire dans le modèle peut rendre impossible l'identification claire de l'effet de la température.

## Les quatre formes de confusion spatiale

- 1 **Biais dû à un facteur confondant omis** : Présence d'une variable non mesurée avec une structure spatiale influençant à la fois la variable et la covariable (Schnell et Papadogeorgou (2020)).
  - exemple : Si une étude analyse la relation entre la pollution de l'air (variable) et des maladies respiratoires (covariable), mais néglige la proximité des zones industrielles (confondeur spatial), les résultats peuvent être faussés.
- 2 **Colinéarité des effets aléatoires** : Modification des estimations des effets fixes due à l'ajout d'effets aléatoires spatiaux colinéaires avec les covariables (Hodges et Reich (2010)).
  - exemple: Dans une régression qui inclut à la fois des variables explicatives comme le revenu et des effets aléatoires spatiaux, si le revenu varie fortement selon des régions spécifiques, les deux composants risquent d'être colinéaires.
- 3 **Biais de régularisation** : Biais dans des échantillons finis associé à l'utilisation de fonctions de régression flexibles comme les splines ou les processus gaussiens (GP) pour modéliser une fonction inconnue de l'espace (Dupont et al. (2020)).
  - exemple : Une spline utilisée pour modéliser l'effet de la distance géographique peut sous-estimer les variations locales importantes en raison d'un lissage excessif.
- 4 **Concurvité** : Difficulté à évaluer l'effet d'une variable lorsque celle-ci est une fonction lisse (ou proche d'une fonction lisse) de l'espace, lorsque le modèle inclut également une fonction lisse arbitraire de l'espace (Ramsay et al. (2003), Paciorek (2010)).
  - exemple : Étudier l'effet de la température (qui varie de manière lisse dans l'espace) tout en incluant une fonction spatiale arbitraire dans le modèle peut rendre impossible l'identification claire de l'effet de la température.

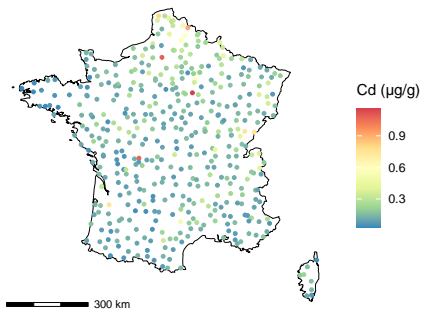


## Les quatre formes de confusion spatiale

- 1 **Biais dû à un facteur confondant omis** : Présence d'une variable non mesurée avec une structure spatiale influençant à la fois la variable et la covariable (Schnell et Papadogeorgou (2020)).
  - exemple : Si une étude analyse la relation entre la pollution de l'air (variable) et des maladies respiratoires (covariable), mais néglige la proximité des zones industrielles (confondeur spatial), les résultats peuvent être faussés.
- 2 **Colinéarité des effets aléatoires** : Modification des estimations des effets fixes due à l'ajout d'effets aléatoires spatiaux colinéaires avec les covariables (Hodges et Reich (2010)).
  - exemple: Dans une régression qui inclut à la fois des variables explicatives comme le revenu et des effets aléatoires spatiaux, si le revenu varie fortement selon des régions spécifiques, les deux composants risquent d'être colinéaires.
- 3 **Biais de régularisation** : Biais dans des échantillons finis associé à l'utilisation de fonctions de régression flexibles comme les splines ou les processus gaussiens (GP) pour modéliser une fonction inconnue de l'espace (Dupont et al. (2020)).
  - exemple : Une spline utilisée pour modéliser l'effet de la distance géographique peut sous-estimer les variations locales importantes en raison d'un lissage excessif.
- 4 **Concurvité** : Difficulté à évaluer l'effet d'une variable lorsque celle-ci est une fonction lisse (ou proche d'une fonction lisse) de l'espace, lorsque le modèle inclut également une fonction lisse arbitraire de l'espace (Ramsay et al. (2003), Paciorek (2010)).
  - exemple : Étudier l'effet de la température (qui varie de manière lisse dans l'espace) tout en incluant une fonction spatiale arbitraire dans le modèle peut rendre impossible l'identification claire de l'effet de la température.

## Contexte biologique

La concentration de Cadmium (Cd) dans les mousses



**Biomarqueurs :** Les mousses accumulent les métaux lourds présents dans l'air

**Données :** 445 mesures de concentration de Cd à travers la France

**Figure:** Carte des mesures de Cd dans les mousses en France (campagne BRAMM 2016)

# Contexte biologique

La covariable  $EMEP_{air}$

Qu'est ce que cette covariable ?

- Valeurs de concentration de Cadmium dans l'air prédites par les modèles EMEP

Les modèles du Programme européen de surveillance et d'évaluation (EMEP) ?

- Modèle de transport atmosphérique
- Réseau limité de stations de mesure (13 en France, principalement urbaines)

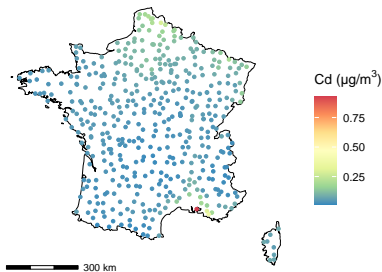


Figure: Carte de valeurs des valeurs extraite d' $EMEP_{air}$  sur les 445 coordonnées des mousses terrestre

# 1<sup>ère</sup> étape : Modèles ajustés sur nos données

## Modèle linéaire univarié (Null)

$$Y(s) = \alpha + \beta_{Null}X(s) + \epsilon(s)$$

$$\epsilon(s) \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau^2)$$

- $Y(s) = \log(Cd(s)) \in \mathbb{R}$ , log de la concentration de Cadmium dans les mousses à la position  $s$
- $X(s) \in \mathbb{R}^+$ , covariable  $EMEP_{air}$ , valeur prédite par un modèle physique
- $s \in \mathbb{R}^2$  coordonnées dans l'espace,  $(s_1, \dots, s_{445}) \in \mathbb{R}^2$  positions des observations

Prendre en compte la corrélation spatiale :

## Modèle spatial

$$Y(s) = \alpha + \beta_{Spatial}X(s) + U(s) + \epsilon(s), \quad U(s) \sim \mathcal{N}(0, \sigma^2\rho(s, s'; \phi))$$

processus gaussien avec fonction de covariance de Matérn

... On soupçonne une confusion spatiale

# 1<sup>ère</sup> étape : Modèles ajustés sur nos données

## Modèle linéaire univarié (Null)

$$Y(s) = \alpha + \beta_{Null}X(s) + \epsilon(s)$$

$$\epsilon(s) \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau^2)$$

- $Y(s) = \log(Cd(s)) \in \mathbb{R}$ , log de la concentration de Cadmium dans les mousses à la position  $s$
- $X(s) \in \mathbb{R}^+$ , covariable  $EMEP_{air}$ , valeur prédite par un modèle physique
- $s \in \mathbb{R}^2$  coordonnées dans l'espace,  $(s_1, \dots, s_{445}) \in \mathbb{R}^2$  positions des observations

Prendre en compte la corrélation spatiale :

## Modèle spatial

$$Y(s) = \alpha + \beta_{Spatial}X(s) + U(s) + \epsilon(s), \quad U(s) \sim \mathcal{N}(0, \sigma^2\rho(s, s'; \phi))$$

processus gaussien avec fonction de covariance de Matérn

... On soupçonne une confusion spatiale

## Objectifs de l'étude

- Réduire la confusion spatiale entre la covariable  $X(s)$  et l'effet spatial  $U(s)$  en utilisant les Spatially Varying Coefficient (SVC) méthodes modélisé par INLA :
  - ▷ Modèle SVC *Gelfand et al., Journal of the American Statistical Association (2003)*
- Interpréter les résultats sur les estimateurs de  $\beta$  et  $\beta(s)$
- Obtenir des cartes de prédiction

## Modèle Spatial Varying Coefficient (SVC)

On veut spatialiser le coefficients de notre covariable  $X(s)$  pour prendre en compte les variations globales et locales des effets des covariables

### Modèle Spatial Varying Coefficient (SVC)

$$Y(s) = \alpha + \beta_{SVC}(s)X(s) + U(s) + \varepsilon(s),$$

$$\varepsilon(s) \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau^2),$$

$$U(s) \sim \mathcal{N}(0, \sigma^2 \rho(s, s'; \phi)),$$

$$\beta_{SVC}(s) \sim \mathcal{N}(0, \sigma_1^2 \rho_1(s, s'; \phi_1))$$

où  $\beta_{SVC}(s)$  est l'effet du coefficient variant dans l'espace à estimer, modélisé comme un Gaussian Random Field (GRF) de moyenne zéro et une matrice de covariance de Matérn à paramétrer,  $\varepsilon(s)$  est un bruit d'erreur iid avec l'écart type,  $U(s)$  est un effet aléatoire spatial modélisé comme un GRF de moyenne zéro et une matrice de covariance de Matérn.

## Modèle Spatial Varying Coefficient + (SVC+)

On veut spatialiser le coefficients de notre covariable  $X(s)$  pour prendre ne compte les variations globales et locales des effets des covariables tout en conservant un effet fixe estimé de notre covariable

### Modèle Spatial Varying Coefficient + (SVC+)

$$Y(s) = \alpha + \beta_{SVC+} X(s) + \beta_{SVC+}(s) X(s) + U(s) + \varepsilon(s),$$

$$\varepsilon(s) \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau^2),$$

$$U(s) \sim \mathcal{N}(0, \sigma^2 \rho(s, s'; \phi)),$$

$$\beta_{SVC+}(s) \sim \mathcal{N}(0, \sigma_1^2 \rho_1(s, s'; \phi_1)),$$

où  $\beta_{SVC+}$  est l'effet fixe de la covariable à estimer,  $\beta_{SVC+}(s)$  est l'effet du coefficient variant dans l'espace



# Introduction à INLA et EDPS

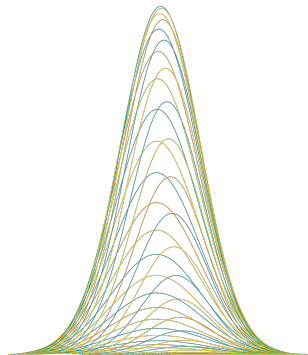
## INLA (Integrated Nested Laplace Approximation) :

- Pour l'inférence bayésienne
- Plus rapide que le MCMC
- Utilisation d'Équation aux dérivées partielles stochastiques (EDPS)

**Package inlabru:** faciliter la modélisation spatiale en utilisant INLA

**Choix des priors** pour les hyperparamètres de la Matérn ; trouver  $\phi_0$  et  $\sigma_0$  tels que

$$\mathbb{P}(\phi < \phi_0) = 0.05 \quad \text{et} \quad \mathbb{P}(\sigma > \sigma_0) = 0.05$$



# INLA

## Priors pour les paramètres et hyperparamètres

### Les paramètres concernés :

- **Paramètres de régression** :  $\beta$ , avec des priors par défaut.
- **Écart-type résiduel** :  $\tau$ , également avec des priors par défaut.

### Hyperparamètres du Matérn (Fuglstad et al. (2019)) :

- **Paramètre d'échelle spatiale (range)** :  $\phi$ 
  - $\mathbb{P}(\phi < \phi_0) = 0.05$ .
  - $\rho_0$  choisi entre  $\frac{1}{10}$  et  $\frac{2}{5}$  de la portée spatiale supposée.
- **Ecart-type** :  $\sigma$ 
  - $\mathbb{P}(\sigma > \sigma_0) = 0.05$ .
  - $\sigma_0$  fixé entre 2.5 et 40 fois l'écart-type attendu.

## EDPS et approximation des processus gaussiens

### EDPS :

Permet d'approximer un processus gaussien continu (GP) par un champ gaussien markovien aléatoire (GMRF) discret à l'aide d'un maillage.

### Équation aux dérivées partielles stochastiques (EDPS) :

$$(\kappa^2 - \Delta)^{\alpha/2}(\tau u(s)) = \mathcal{W}(s)$$

où :

- $\kappa > 0$  : paramètre d'échelle spatiale.
- $\Delta$  : opérateur laplacien.
- $\alpha$  : contrôle la régularité du processus ( $\alpha = 2$  ici).
- $\tau$  : paramètre de l'écart-type.
- $s \in \mathbb{R}^2$  est la position spatiale.
- $\mathcal{W}(s)$  : bruit blanc spatial gaussien.

**Solution** : Le processus spatial  $u(s)$  possède une fonction de covariance de type Matérn qui va être approximée.

# Approximation du Processus Spatial

Le champ spatial est approximé par une somme pondérée de fonctions de base càd un mesh définit comme suit :

$$u(s) = \sum_{k=1}^n \psi_k(s) u_k$$

## Définitions :

- $\psi_k(s)$  : Fonctions de base déterministes définies sur le maillage.
- $u_k$  : Poids gaussiens associés à chaque fonction de base.

La distribution jointe des poids  $u = \{u_1, \dots, u_n\}$  est choisie pour que  $u(s)$  approxime un processus gaussien continu.

**Objectif** : Construire un maillage triangulé qui sert de base à la représentation de l'EDPS.

# Construction du Maillage (Mesh)

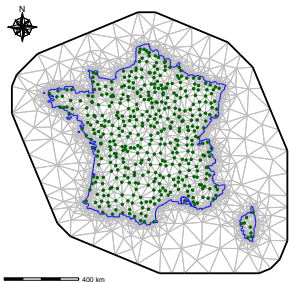


Figure: Carte du mesh de la France

## Caractéristiques :

- 3572 nœuds au total.
- Résolution plus fine dans les zones denses en observations.
- Extension pour réduire la variance près des bordures.

## Principe :

- Petits triangles : Zones d'intérêt avec des observations denses (résolution plus fine).
- Grands triangles : Zones éloignées ou sans données (réduction du coût computationnel).

**Gestion des bordures :** Le maillage est étendu au-delà de la région d'intérêt pour minimiser les effets de bordure.

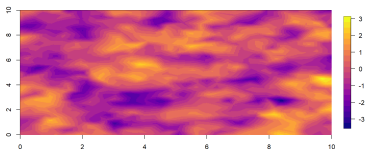
## Approche pratique :

- La frontière interne (ligne bleue) correspond au contour de la France.
- La frontière externe (ligne noire) englobe la zone étendue.
- Les points verts représentent les emplacements des observations.

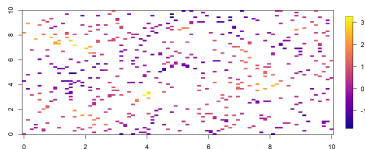
# Plan de simulation I

## Génération des données :

- Simulation de 50 jeux de données, avec 500 observations chacun.
- Création de deux processus gaussiens distincts avec une structure de covariance de Matérn spécifique
- Sélection aléatoire de 500 points dans un pavé spatial  $[0, 10] \times [0, 10] \subset \mathbb{R}^2$ .



(a) Champ spatial simulé



(b) Échantillonnage de 500 points sur le champ spatial simulé

## Plan de simulation II

### Création de la confusion spatiale :

- Mélange des deux champs gaussiens selon un certain ratio pour introduire de la colinéarité
- Construction<sup>1</sup> de la covariable  $\mathbf{X}(s)$  et de la variable réponse  $\mathbf{Y}(s)$  comme suit :

$$X(s) = 0.35z_x(s) + \epsilon_x(s), \quad \epsilon_x(s) \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_x^2)$$

$$U(s) = z_u(s) - z_x(s)$$

$$Y(s) = \beta X(s) + U(s) + \epsilon_y(s), \quad \epsilon_y(s) \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_y^2)$$

- Paramètres utilisés :  $\sigma_y = 1$ ,  $\sigma_x = 0.1$ ,  $\beta = 3$ .

---

<sup>1</sup>en s'appuyant sur Dupont et al. (2022)

## Résultats des simulations

$$\beta = 3$$

Model	$\hat{\beta}$	ESD	$\overline{SE}$	DIC	WAIC	CI
Null	1.031	0.136	0.149	1582.51	1582.529	0
Spatial	2.649	0.800	0.303	1498.956	1500.482	76
SVC				1484.652	1486.608	
SVC+	2.960	0.398	0.402	1481.641	1483.453	96

**Table:** Tableau des valeurs à étudier pour les estimateurs de  $\beta$ . A noter que lorsque qu'on parle de médiane de  $\hat{\beta}$ , il s'agit de la moyenne sur les 50 simulations de la médiane à posteriori du  $\beta$  estimé, de même pour la variance et la moyenne. Toutes les valeurs ont été arrondi à  $10^{-3}$ . Avec  $\rho$  et  $\sigma$  prior fixé à 0.05 et 3, respectivement, avec une des probabilité de 0.05.



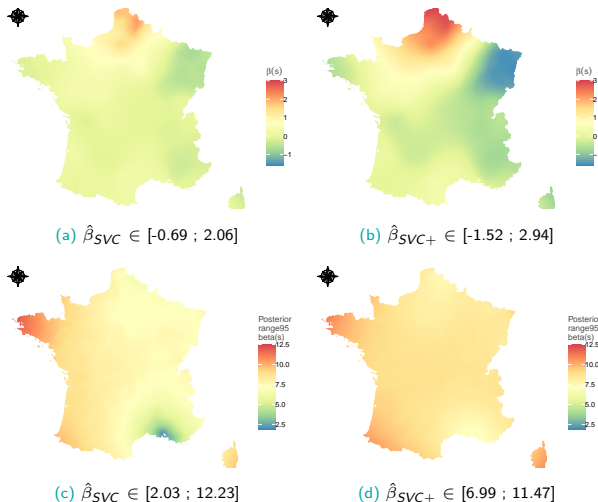
# Résultats pour les données réelles



	Null	Spatial	SVC	SVC+
$\hat{\beta}$	2.306	0.230		0.115
0.025quant	1.422	-0.820		-3.408
0.975quant	3.191	1.280		3.637
$\hat{\phi}_1$			404 [88 ; 1592]	588 [121 ; 3338]
$\hat{\phi}$		617 [288 ; 1562]	650 [264 ; 1865]	659 [259 ; 2022]
$\hat{\sigma}_1$			1.487 [0.382 ; 4.306]	2.119 [0.708 ; 5.702]
$\hat{\sigma}$		0.474 [0.300 ; 0.790]	0.448 [0.258 ; 0.821]	0.443 [0.252 ; 0.822]
$\hat{\tau}$	0.572 [0.536 ; 0.611]	0.467 [0.431 ; 0.506]	0.467 [0.431 ; 0.506]	0.467 [0.431 ; 0.506]
DIC	769.07	628.45	631.46	627.62
WAIC	775.49	629.28	632.88	627.98

**Table:** Résumé des estimations des paramètres du modèle, des critères DIC et WAIC pour l'ensemble de données sur les mousses Cd. Moyennes postérieures et intervalles crédibles à 95% de  $\hat{\beta}$ , 50 (2,5, 97,5) centiles pour les autres paramètres.

# Cartes des médianes postérieures avec leur intervalle postérieur à 95% de $\beta(s)$



**Figure:** Estimations médianes postérieures (en haut) et estimations postérieures à 95% (en bas) de  $\beta(s)$  pour le modèle SVC (à gauche) et le modèle SVC+ (à droite).

# Cartes des médianes à posteriori de $U(s)$ avec leur intervalle postérieur à 95%.

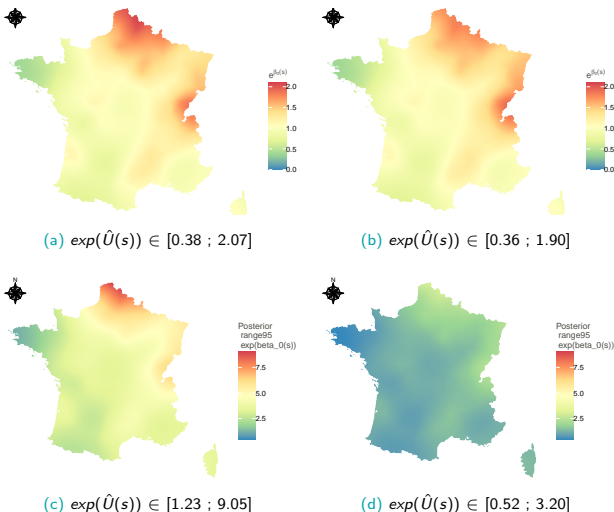
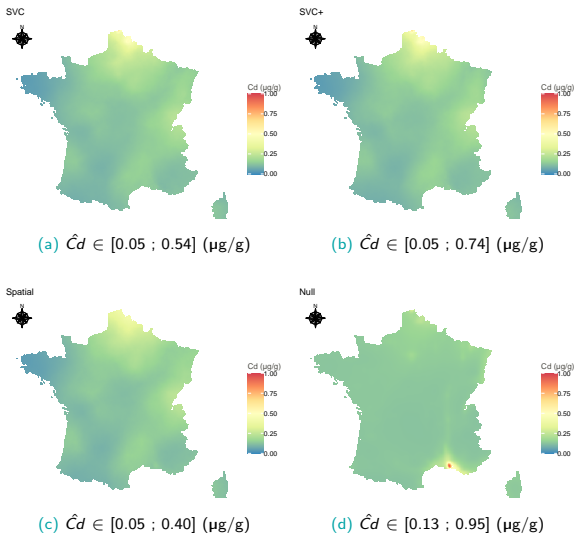


Figure: Les estimations médianes postérieures de l'effet spatial aléatoire,  $U(s)$ , vont de 0 à 3 pour le modèle SVC (à gauche) et le modèle SVC+ (à droite), avec leur intervalle de 95% (ligne 21 / 23

# Cartes de prédiction de la médiane postérieure des concentrations de Cd dans les mousses avec les modèles SVC, SVC+, spatial et Null



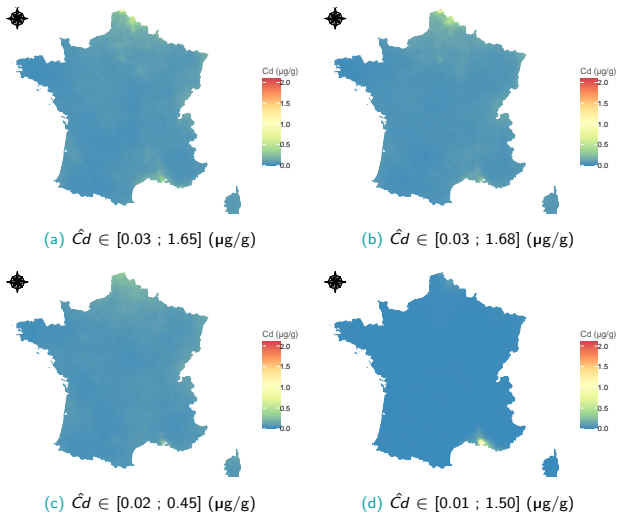
## Discussion et Perspectives

Pour aller plus loin avec le modèle SVC :

- Il serait intéressant d'explorer d'autres modèles à coefficients variables dans l'espace, tels que la régression pondérée géographiquement (GWR) (Fotheringham et al., 2009, qui peut capturer les variations locales dans les relations entre les covariables et les résultats, ou également le filtrage spatial par vecteur propre (ESF) (Mu-rakami and Griffith, 2015)
- Développer ce travail en mettant en œuvre des modèles SVC multivariés, ce qui permettrait de saisir des effets d'interaction complexes

Merci !

# Cartes de prédiction



Figure