Estimation sous contraintes pour des données tronquées Application au cas des altitudes d'arrêts en Haute-Savoie

Aurore Lavigne, Liliane Bel, Nicolas Eckert et Eric Parent

journée RESSTE, 11 mai 2015





• • = • • = •

Introduction

Cartes d'aléas/risques, des outils indispensables à la gestion du risque d'avalanches

Exemple de carte de risque : probabilité annuelle de décès (+/- proportionnelle à la probabilité d'atteinte)



Proposer une nouvelle méthode de construction de cartes d'aléa à l'échelle régionale, s'appuyant sur la dépendance spatiale entre couloirs. Objectifs :

- Démontrer la dépendance spatiale entre les altitudes d'arrêts moyenne des couloirs.
- Utiliser cette dépendance spatiale pour prédire la distribution des altitudes d'arrêts sur un nouveau couloir.

<日

<</p>

Altitudes d'arrêts enregistrées par l'EPA

Altitudes d'arrêts en Haute-Savoie

- 389 couloirs avec au moins une donnée
- 5331 avalanches
- Période d'étude : 1925-2012
- En 2002, un seuil d'observation s_c est défini pour chaque couloir => troncature à droite



d'avalanches

伺下 イヨト イヨト

Couloirs

Haute-Savoie

Données topologiques

- Altitude de la vallée au pied du couloir h_c.
- Exposition de la zone de départ.

en

Introduction

Exemples de tracés EPA



5 / 25

Des données hétérogènes

Activité des couloirs

- ▶ 162 couloirs avec moins de 5 avalanches.
- ► 51 couloirs avec plus de 30 événements.
- 177 avalanches enregistrées pour le couloir le plus actif.

Variabilités de la variance inter-couloir

- Dépends de la topographie et des caractéristiques de la zone d'arrêt.
- Les écart-types s'étendent de 0! à 250 m.
- L'erreur d'observation est d'environ 10 m 50 m.

Différence de dénivelé entre le seuil d'observation et la vallée

- autour de 300 m,
- s'étend de 0 à 1000 m.
- Les modèles de dépassement de seuil ne peuvent pas être utilisés.

< □ > < □ > < □ > < □ > < □ > < □ >

Introduction

Plan

Introduction

Modélisation spatio-temporelle des altitudes d'arrêts

Problème d'estimation liées aux données tronquées

Résultats

э

イロト イボト イヨト イヨト

Plan

Introduction

Modélisation spatio-temporelle des altitudes d'arrêts

Problème d'estimation liées aux données tronquées

Résultats

э

Modélisation des observations par une distribution tronquée

- ► Y_{cti} altitude d'arrêt de l'avalanche i, sur le couloir c, l'année t
- s_c seuil d'observation du couloir c
- μ_{ct} espérance de la distribution complète (non tronquée)
- σ_c^2 variance du couloir *c* pour la distribution complète

$$f(y_{cti}|\mu_{ct},\sigma_c^2) = \frac{1}{\Phi\left(\frac{s_c-\mu_{ct}}{\sigma_c}\right)} \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left(\frac{-1}{2\sigma_c^2}(y_{cti}-\mu_{ct})^2\right) \mathbf{1}_{y_{cti} < s_c}$$

・ロト ・ 同ト ・ ヨト ・ ヨト

Modélisation des observations par une distribution tronquée

- ► Y_{cti} altitude d'arrêt de l'avalanche i, sur le couloir c, l'année t
- s_c seuil d'observation du couloir c
- µ_{ct} espérance de la distribution complète (non tronquée)
- σ_c^2 variance du couloir *c* pour la distribution complète

$$f(y_{cti}|\mu_{ct},\sigma_c^2) = \frac{1}{\Phi\left(\frac{s_c-\mu_{ct}}{\sigma_c}\right)} \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left(\frac{-1}{2\sigma_c^2}(y_{cti}-\mu_{ct})^2\right) \mathbf{1}_{y_{cti} < s_c}$$

On tient compte de la troncature et de l'hétéroscédasticité.

Modélisation spatio-temporelle des altitudes d'arrêts

Un modèle spatio-temporel additif

$$\mu_{ct} = \alpha + B_t + C_c. \tag{1}$$

イロト イヨト イヨト ・

- α altitude d'arrêt moyenne
- ► *B_t* terme temporel, modélisé par une marche aléatoire à l'ordre deux
- C_c terme spatial

On ajoute les contraintes d'identifiabilité : $\sum_t B_t = \sum_c C_c = 0$

Modèle spatial

Le couloir est réduit à un point c dont les coordonnées appartiennent au domaine spatial \mathcal{D} . Les termes spatiaux C_c sont modélisés par un processus gaussien, avec une covariance exponentielle et une pépite. Dépendance spatiale

$$\operatorname{cov}(C_c, C_{c'}) = \begin{cases} \tau^2 \exp(-\frac{h_{cc'}}{\phi}) \text{ if } c \neq c' \\ \rho^2 + \tau^2 \text{ otherwise,} \end{cases}$$

avec

- $h_{cc'}$ distance entre c et c',
- ϕ paramètre de portée, τ^2 pallier, et ρ^2 pépite.

Régression

$$E(C_c) = a + bh_c + c_{o_c} \begin{cases} h_c \text{ altitude de la vallée} \\ o_c \in \{N, S\} \text{ exposition du couloir} \end{cases}$$

イロト 不得下 イヨト イヨト

Plan

Introduction

Modélisation spatio-temporelle des altitudes d'arrêts

Problème d'estimation liées aux données tronquées

Résultats

э

イロト イボト イヨト イヨト

Problèmes d'estimation engendrés par la troncature

Cas d'une variable normale tronquée à droite

- Espérance et variance d'une distribution gaussienne tronquée sont corrélées.
- La vraisemblance est distordue et présente une large zone plate.



Log-vraisemblance d'un échantillon Gaussien tronqué à droite à 0.1.

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

Conséquences sur l'inférence fréquentiste

- La variance de l'estimateur du maximum de vraisemblance est grande
- Ajout d'un biais pour diminuer la variance : vraisemblance restreinte (Cope, 2011), vraisemblance pénalisée (A'Hearn, 2004)

Inférence bayésienne

Priors vagues

- $[lpha] \propto {
 m cte}$, a, b, $c \sim N(0, 10^7)$
- $\phi \sim U[0, 30000]$
- $\tau^2 \sim \text{IG}(0.001, 0.001)$ et $\rho^2 \sim \text{IG}(0.001, 0.001)$

Prior informatif, élicité

▶ $\sigma_c^2 \sim \text{IG}$ avec pour espérance et écart-type 25² tronquée avant 30² et après 316².

・ 同 ト ・ ヨ ト ・ ヨ ト

Inférence bayésienne

Priors vagues

- $[lpha] \propto {
 m cte}$, a, b, $c \sim N(0, 10^7)$
- $\phi \sim U[0, 30000]$
- $\tau^2 \sim \text{IG}(0.001, 0.001)$ et $\rho^2 \sim \text{IG}(0.001, 0.001)$

Prior informatif, élicité

▶ $\sigma_c^2 \sim \text{IG}$ avec pour espérance et écart-type 25² tronquée avant 30² et après 316².

On tient compte de l'erreur de mesure.

A (1) < A (2) < A (2) </p>

Échantillonneur de Gibbs

Objectif : tirer dans les distributions *a posteriori* conditionnelles complètes des paramètres et des variables latentes.

Deux principales difficultés :

- Covariance exponentielle.
- Troncature

Utilisation d'une variable auxiliaire Z_{cti} , correspondant à l'observation sous la loi complète, telle que $\mathbf{P}(Y_{cti} < y) = \mathbf{P}(Z_{cti} < z)$ (Griffiths, 2004).

$$Z_{cti} = \mu_{ct} + \sigma \Phi^{-1} \left(\frac{\Phi(\frac{Y_{cti} - \mu_{ct}}{\sigma_c})}{\Phi(\frac{s_c - \mu_{ct}}{\sigma_c})} \right)$$

Problème rencontré

```
Dès la première itération de l'algorithme :
Chaine 1 iteration 1
Z[4969] is not a number
mu = 2014.806636
seuil = 1000
```

```
y[i] = 995
sig = 118.9297427
```

Sur un couloir bien documenté :

Altitude d'arrêt	1240	1250	1260	1270	1280	1290	1300	1305
Nb d'obs	2	13	9	7	12	4	24	3

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Solution : contraindre $\alpha + C_c$ à un domaine raisonnable

Nous savons que :

- > au moins 50% des avalanches sont observées sur chaque couloir,
- l'altitude de la vallée h_c est une borne inférieure pour toutes les avalanches.

・ロト ・四ト・モート・モート

Solution : contraindre $\alpha + C_c$ à un domaine raisonnable

Nous savons que :

- ▶ au moins 50% des avalanches sont observées sur chaque couloir,
- l'altitude de la vallée h_c est une borne inférieure pour toutes les avalanches.

Ainsi, $h_c \leq \alpha + C_c \leq s_c$.

Tirage de C_c , partie 1 : contraintes d'identifiabilités

On cherche à tirer C dans $[C|\theta, Y]$, connaissant la contraintes $\sum_{c=1}^{n_c} C_c = 0.$ (θ : autres paramètres et variables)

Idée (Gelman, 2005) : Tirer C_1, \dots, C_{n_c-1} dans $C|\theta, Y, \mathbf{1}'_c C = 0$ et prendre $C_{n_c} = -\sum_{c=1}^{n_c-1} C_c$.

Si $C| heta, Y \sim \mathcal{N}(m, \Sigma)$ alors la loi jointe de C| heta, Y et de $\mathbf{1}_c C = 0$ s'écrit

$$\begin{bmatrix} C|\theta, Y \\ \mathbf{1}'_{c}C \end{bmatrix} \sim N\left(\begin{bmatrix} m \\ \mathbf{1}'_{c}m \end{bmatrix}, \begin{bmatrix} \Sigma & \Sigma\mathbf{1}'_{c} \\ \mathbf{1}'_{c}\Sigma & \mathbf{1}'_{c}\Sigma\mathbf{1}'_{c} \end{bmatrix}\right)$$

La distribution conditionnelle est alors

$$C|\theta, Y, \mathbf{1}'_{c}C = 0 \sim \mathcal{N}(m_{0}, \Sigma_{0}) \text{ with } \begin{cases} m_{0} = m - \frac{\Sigma \mathbf{1}_{c}\mathbf{1}'_{c}m}{\mathbf{1}'_{c}\Sigma \mathbf{1}_{c}} \\ \Sigma_{0} = \Sigma - \frac{\Sigma \mathbf{1}_{c}\mathbf{1}'_{c}\Sigma}{\mathbf{1}'_{c}\Sigma \mathbf{1}_{c}} \end{cases}$$

▲□▶ ▲圖▶ ▲国▶ ▲国▶ - ヨー のへの

Tirage de C_c , partie 2 : $h_c - \alpha < C_c < s_c - \alpha$

Idées (Rodriguez-Yam et al, 2004) :

- ► Ecrire les contraintes sous forme vectoriel $\mathbf{V}C < v$, \mathbf{V} de dimension $2n_c \times n_c 1$.
- Enlever la dépendance entre les composantes de C en considérant w = LC avec L telle que LΣ₀L' = I_{nc-1}
- Tirer successivement et conditionnellement aux autres, chaque composant de w en respectant ses contraintes VL⁻¹w < v.</p>

Pour ne pas recalculer $VL^{-1}w$ à chaque composante on remarque que pour la contrainte *i* sur la composante *j*₀ :

$$(\mathbf{V}L^{-1}w)_{i} = \underbrace{\sum_{j\neq j_{0}}^{n_{c}-1} \left(\sum_{l=1}^{n_{c}-1} \mathbf{V}_{il}L_{lj}^{-1}\right)w_{j}}_{u_{i}^{-j_{0}}} + \underbrace{\left(\sum_{l=1}^{n_{c}-1} \mathbf{V}_{il}L_{lj_{0}}^{-1}\right)}_{u_{i}^{j_{0}}}w_{j_{0}} < v_{i}.$$

La contrainte *i* pour la composante j_0 de **w** s'écrit $w_{j_0} < \frac{v_i - u_i^{-j_0}}{u_i^{j_0}}$.

Algorithme pour tirer C_1, \dots, C_{n_c-1} dans sa distribution tronquée

- Calculer le vecteur $\mathbf{u} = \mathbf{VC}$;
- Initialiser la matrice L, son inverse L^{-1} , et le vecteur w = LC;
- Pour chaque j dans $\{1, \cdots n_c 1\}$
 - Calculer le vecteur de dimension $2n_c$, $\mathbf{u}^j = (u_1^j, \cdots u_k^j)$, et $\mathbf{u}^{-j} = \mathbf{u} \mathbf{u}^j w_j$;
 - ► Trouver l'intervalle [a_j, b_j] dans lequel tirer w_j pour satisfaire les 2n_c contraintes;
 - ► Tirer w_j dans la distribution normale tronquée N((Lm₀)_j, 1) de support [a_j, b_j].
 - Mettre à jour **u** comme $\mathbf{u} = \mathbf{u}^{-j} + \mathbf{u}^j w_j$;
- Obtenir C avec C = L⁻¹w.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

Résultats

Plan

Introduction

Modélisation spatio-temporelle des altitudes d'arrêts

Problème d'estimation liées aux données tronquées

Résultats

э

<ロト <問ト < 目と < 目と

Résultats

Résultats-Spatial

- Terme spatial représente 98.5% de la variabilité totale
- Dans le terme spatial, le terme de régression représente 86% de la variabilité.
- Portée effective estimée à 10 km.





	portée	pallier	pépite	cte	alt. val.	exposition
	ϕ	$ au^2$	ρ^2	а	Ь	CS
esp.	3037.61 m	4644.89 m ²	2426.96 m ²	-935.06 m	0.83	7.02 m
éc. type	1790.80 m	$1192.07 \ { m m}^2$	616.54 m ²	31.02 m	0.03	9.40 m

Résultats

Résultats-Temporel

 Terme spatial représente seulement 1.5% de la variabilité spatio-temporelle



	Espérance	Ecart-type		
δ_0	137.7 m ²	33.1 m ²		
δ_1	0.88 <i>m</i> ²	1.09 <i>m</i> ²		

Conclusion

- Identification d'un problème de convergence de l'algorithme de Gibbs
 - Troncature des données
 - > Distribution des données presque discrètes à cause de l'erreur de mesure
- Solution : ajout de contraintes
 - Connaissance a priori
 - Peu restrictives
 - Impossible de mesure son effet sur les estimateurs dans ce cas

<日

<</p>

A'Hearn B (2004) A restricted maximum likelihood estimator for truncated height samples. Economics & Human Biology 2(1) :5–19

Cope EW (2011) Penalized likelihood estimators for truncated data. Journal of Statistical Planning and Inference 141(1):345–358

Griffiths W (2004) A Gibbs' sampler for the parameters of a truncated multivariate normal distribution. In : Becker R, Hurn S (eds) Contemporary issues in economics and econometrics : theory and application, Edward Elgar Pub, pp 75–91

Rodriguez-Yam G, Davis R, Scharf L (2004) Efficient Gibbs sampling of truncated multivariate normal with application to constrained linear regression. Unpublished Manuscript

イロト 不得 トイヨト イヨト