

Estimation non-paramétrique de mesures de risque pour des lois conditionnelles à queues lourdes avec application à des extrêmes pluviométriques

par

Jonathan EL METHNI

en collaboration avec

Laurent GARDES & Stéphane GIRARD

Congrès SMAI 2015 Les Karellis
Lundi 8 Juin 2015



- 1 Introduction
 - Théorie des valeurs extrêmes
 - Problématique
 - Mesures de risque
- 2 Cadre de travail et résultats asymptotiques
- 3 Applications à un jeu de données hydrologique
- 4 Conclusions et perspectives

- **But de la théorie des valeurs extrêmes** \implies étudier et caractériser le comportement des valeurs extrêmes d'un échantillon de variables aléatoires.
- L'exemple d'application historique est l'hydrologie \implies Gumbel années 50.

La quantité de pluie est modélisée par une variable aléatoire Y ayant pour fonction de survie $\bar{F}(y) = 1 - F(y) = \mathbb{P}(Y \geq y)$.

On dispose de $Y_{1,n} \leq \dots \leq Y_{n,n}$ un échantillon ordonné de quantités de pluies annuelles.

Les hydrologues souhaitent estimer le niveau de pluie H qui est atteint ou dépassé en moyenne une fois sur T années *i.e.* on veut estimer H tel que

$$1/T = \mathbb{P}(Y \geq H) = \bar{F}(H)$$

autrement dit on veut estimer

$$H = \bar{F}^{-1}(1/T)$$

- H sera appelé niveau de retour correspondant à une période de retour T .
- C'est la quantité standard d'intérêt dans les études environnementales.

Le quantile d'ordre $\alpha \in]0, 1[$ de la fonction de survie est défini par

$$q(\alpha) = \bar{F}^{-1}(\alpha).$$

\implies Le niveau de retour $H = \bar{F}^{-1}(1/T)$ est donc un quantile d'ordre $\alpha = 1/T$.

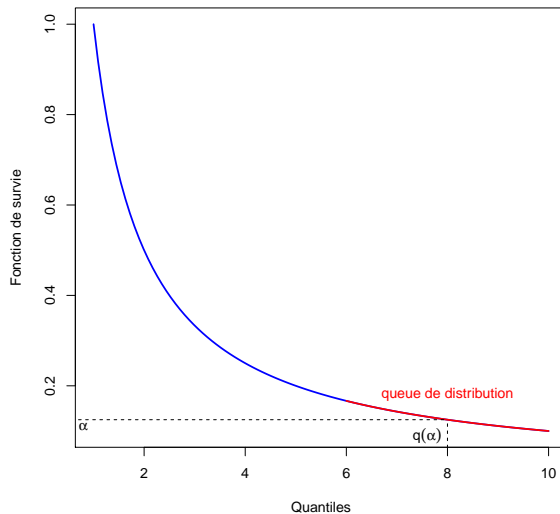
- Que se passe-t-il si la période de retour T est supérieure à la durée d'observation ?
Autrement dit que se passe-t-il si

$$T > n \iff \alpha = 1/T < 1/n \xrightarrow{n \rightarrow \infty} 0 \quad ?$$

On souhaite estimer des quantiles dit extrêmes $q(\alpha_n)$ d'ordre α_n définis par

$$q(\alpha_n) = \bar{F}^{-1}(\alpha_n) \quad \text{avec} \quad \alpha_n \rightarrow 0 \quad \text{quand} \quad n \rightarrow \infty$$

- Un niveau de retour dont la période de retour est supérieure à la durée d'observation est un **quantile extrême**.



Que se passe-t-il si $q(\alpha_n) > Y_{n,n}$?

On peut montrer que l'on a

$$\mathbb{P}(q(\alpha_n) > Y_{n,n}) = \exp(-n\alpha_n(1 + o(1)))$$

- 1er cas :

$$\text{Si } n\alpha_n \rightarrow \infty \text{ alors } \mathbb{P}(q(\alpha_n) > Y_{n,n}) \rightarrow 0.$$

Un estimateur naturel est la statistique d'ordre $Y_{n-\lfloor n\alpha_n \rfloor+1,n}$ (où $\lfloor \cdot \rfloor$ est la fonction partie entière).

- 2ème cas :

$$\text{Si } n\alpha_n \rightarrow 0 \text{ alors } \mathbb{P}(q(\alpha_n) > Y_{n,n}) \rightarrow 1.$$

Dans ce cas, on ne peut pas estimer $q(\alpha_n)$ en inversant simplement la fonction de répartition empirique :

$$\hat{F}_n(y) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{Y_i \leq y\} \quad \text{car} \quad \hat{F}_n(y) = 1 \quad \text{pour} \quad y \geq Y_{n,n}.$$

Le comportement de $Y_{n,n}$ est caractérisé par sa fonction de répartition $F_{Y_{n,n}}(y) = F^n(y)$ qui est une loi dégénérée.

Théorème

Soit $(Y_n)_{n \geq 1}$ une suite de variables aléatoires indépendantes et identiquement distribuées de fonction de répartition F . S'il existe deux suites normalisantes réelles $(a_n)_{n \geq 1} > 0$ et $(b_n)_{n \geq 1} \in \mathbb{R}$ et une loi non-dégénérée \mathcal{H}_γ telles que

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\frac{Y_{n,n} - b_n}{a_n} \leq y \right) = \lim_{n \rightarrow \infty} F^n(a_n y + b_n) = \mathcal{H}_\gamma(y),$$

alors à une translation et un changement d'échelle près on a

$$\mathcal{H}_\gamma(y) = \exp \left(-(1 + \gamma y)_+^{-1/\gamma} \right),$$

où $\gamma \in \mathbb{R}$ et $z_+ = \max(0, z)$.

- \mathcal{H}_γ est appelée fonction de répartition de la loi des valeurs extrêmes.
- Si F vérifie le théorème de Fisher-Tippett-Gnedenko on dit que F appartient au domaine d'attraction de \mathcal{H}_γ .
- Cette loi dépend du seul paramètre de forme γ appelé **indice des valeurs extrêmes ou indice de queue**.

Trois domaines d'attraction

Selon le signe de γ , on distingue trois domaines d'attraction :

- si $\gamma < 0$, on dit que F appartient au domaine d'attraction de **Weibull**. Il contient des lois dont la fonction de survie n'a **pas de queue de distribution**.
- si $\gamma = 0$, on dit que F appartient au domaine d'attraction de Gumbel. Il contient les lois dont la fonction de survie est à décroissance exponentielle, *i.e.* les lois à queues légères.
- si $\gamma > 0$, on dit que F appartient au domaine d'attraction de **Fréchet**. Il contient les lois dont la fonction de survie est à décroissance polynomiale, *i.e.* les lois à **queues lourdes**.

Fréchet ($\gamma > 0$)	Gumbel ($\gamma = 0$)	Weibull ($\gamma < 0$)
Pareto	Normale	Uniforme
Student	Exponentielle	Beta
Burr	Log-normale	ReverseBurr
Chi-deux	Gamma	
Fréchet	Weibull	
Log-gamma	Logistique	
Log-logistique	Gumbell	
Cauchy		

Toutes les lois appartenant au domaine d'attraction de Fréchet peuvent se réécrire

$$\bar{F}(y) = y^{-1/\gamma} \ell(y) \quad \text{avec } \gamma > 0,$$

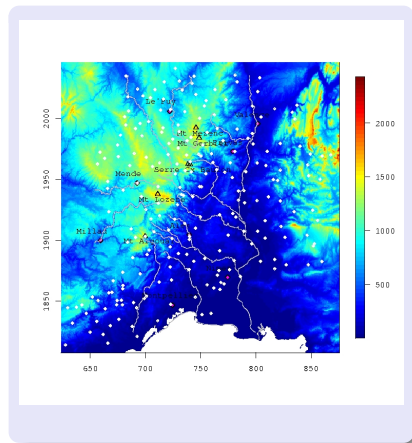
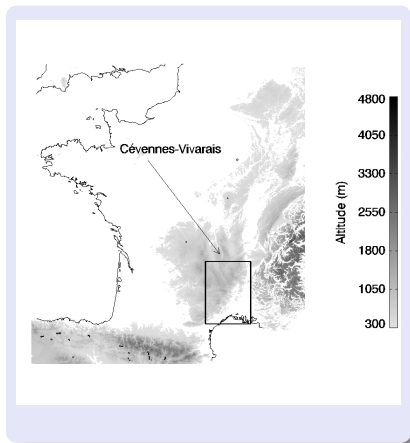
et ℓ est une **fonction à variations lentes** (f.v.l) à l'infini i.e. $\forall \lambda \geq 1$,

$$\lim_{y \rightarrow \infty} \frac{\ell(\lambda y)}{\ell(y)} \rightarrow 1.$$

- Le paramètre γ contrôle le comportement de la queue de la fonction de survie et donc celui des valeurs extrêmes.

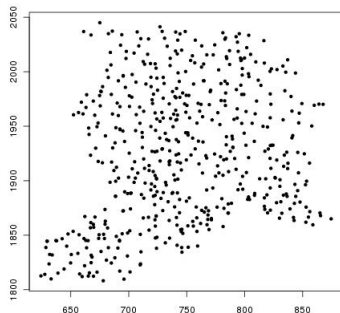
Dans de nombreuses applications, la variable d'intérêt est mesurée simultanément avec une covariable.

La région Cévennes-Vivarais



Horizontalement on a la longitude (en km), verticalement la latitude (en km) et en échelle de couleurs l'altitude (en m). Quelques stations d'observations (losanges blancs).

Les 523 stations d'intérêt



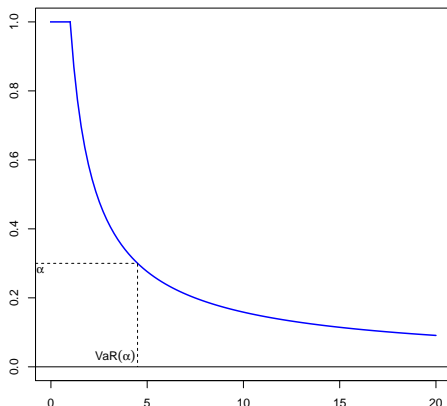
- Données fournies par Météo-France.
- Y : hauteur de pluie journalière en mm.
- $X = \{\text{longitude, latitude, altitude}\}$.

- 1958 \implies 2000 soit 43 ans de données.
- 523 stations notées $\{x_t; t = 1, \dots, 523\}$.
- Nombre total d'observation = 5 513 734.

But \implies obtenir des cartes d'estimation de mesures de risque des pluies journalières correspondant à une période de retour de 100 ans en tout point de la région.

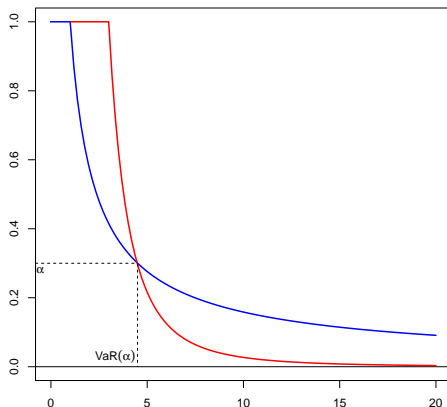
- Soit Y une variable aléatoire désignant un montant de pertes. La Value-at-Risk au niveau $\alpha \in]0, 1[$ notée $\text{VaR}(\alpha)$ introduite en 1993 est définie par

$$\text{VaR}(\alpha) := \bar{F}^{-1}(\alpha).$$



- La $\text{VaR}(\alpha)$ est le quantile d'ordre α de la fonction de survie de la v.a. Y .

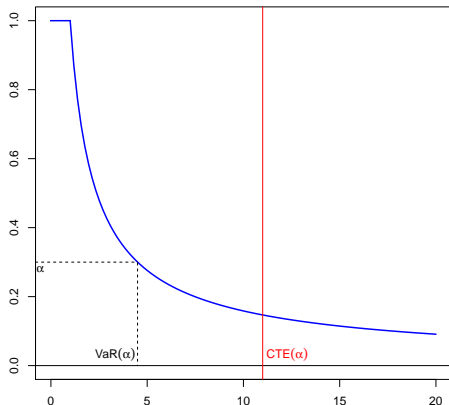
- Soit Y_1 et Y_2 deux v.a. de pertes ayant pour fonctions de survies associées \bar{F}_1 et \bar{F}_2 .



- Un des principaux reproches fait à la VaR est que des v.a à **queues légères** et à **queues lourdes** (Embrechts *et al.* [1997]) peuvent avoir la même $\text{VaR}(\alpha)$.

- La Conditional Tail Expectation au niveau $\alpha \in]0, 1[$ notée $CTE(\alpha)$ est définie par

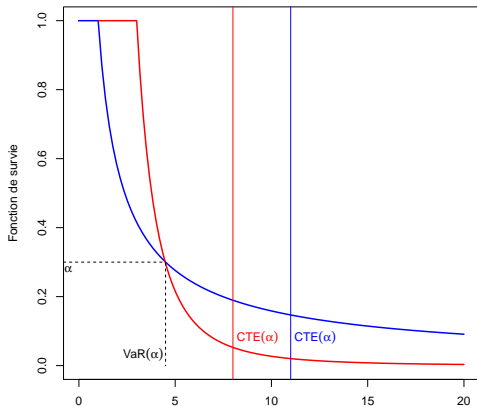
$$CTE(\alpha) := \mathbb{E}(Y|Y > VaR(\alpha)).$$



- La $CTE(\alpha)$, donne des informations sur la distribution de Y au delà de la $VaR(\alpha)$ et donc contrairement à la $VaR(\alpha)$, sur l'épaisseur de la queue de distribution.

- La Conditional Tail Expectation au niveau $\alpha \in]0, 1[$ notée $CTE(\alpha)$ est définie par

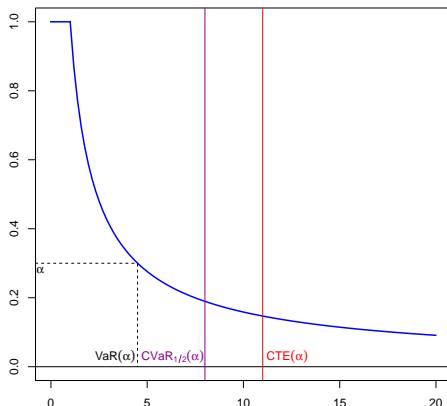
$$CTE(\alpha) := \mathbb{E}(Y|Y > VaR(\alpha)).$$



- La $CTE(\alpha)$, donne des informations sur la distribution de Y au delà de la $VaR(\alpha)$ et donc contrairement à la $VaR(\alpha)$, sur l'épaisseur de la queue de distribution.

- La Conditional-Value-at-Risk au niveau $\alpha \in]0, 1[$ notée $CVaR_\lambda(\alpha)$ introduite par Rockafellar et Uryasev [2000] est définie par

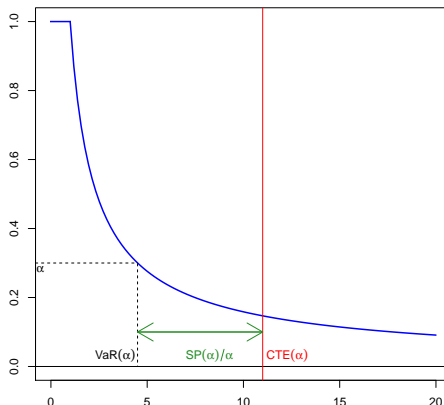
$$CVaR_\lambda(\alpha) := \lambda VaR(\alpha) + (1 - \lambda)CTE(\alpha) \quad \text{avec} \quad 0 \leq \lambda \leq 1.$$



- On peut remarquer que $CVaR_1(\alpha) = VaR(\alpha)$ et que $CVaR_0(\alpha) = CTE(\alpha)$.

- La mesure de risque Stop-loss Premium reinsurance avec un niveau de rétention égal à $\text{VaR}(\alpha)$ (voir Cai et Tan [2007]) est définie par :

$$\text{SP}(\alpha) := \alpha (\text{CTE}(\alpha) - \text{VaR}(\alpha)).$$



- Cette mesure de risque permet ainsi de mettre en évidence les cas dangereux.

L'apport nouveau de nos travaux consiste en l'ajout de deux difficultés supplémentaires dans le cadre de l'estimation de mesures de risque.

① On ajoute la présence d'une **covariable aléatoire** $X \in \mathbb{R}^p$.

② On s'intéresse à l'estimation de mesures de risque dans le cas de **pluies extrêmes**.

⇒ Pour cela on remplace l'ordre fixé $\alpha \in]0, 1[$ par une suite $\alpha_n \rightarrow 0$ quand la taille de l'échantillon $n \rightarrow \infty$.

En notant par $\bar{F}(\cdot|x)$ la fonction de survie conditionnelle de Y sachant $X = x$, on définit la Regression Value-at Risk par :

$$\text{RVaR}(\alpha_n|x) := \bar{F}^{-1}(\alpha_n|x),$$

et la Regression Conditional Tail Expectation par :

$$\text{RCTE}(\alpha_n|x) := \mathbb{E}(Y|Y > \text{RVaR}(\alpha_n|x), X = x).$$

⇒ Les mesures de risque dépendent alors seulement de RVaR et de RCTE.

$$\begin{aligned} \text{RCVaR}_\lambda(\alpha_n|x) &= \lambda \text{RVaR}(\alpha_n|x) + (1 - \lambda) \text{RCTE}(\alpha_n|x), \\ \text{RSP}(\alpha_n|x) &= \alpha_n (\text{RCTE}(\alpha_n|x) - \text{RVaR}(\alpha_n|x)). \end{aligned}$$

⇒ On veut estimer toutes les mesures de risque mentionnées.

Il nous faut obtenir la loi jointe asymptotique de $\left\{ \widehat{\text{RCTE}}_n(\alpha_n|x), \widehat{\text{RVaR}}_n(\alpha_n|x) \right\}$.

On définit le moment conditionnel d'ordre $a \geq 0$ de Y sachant $X = x$ par

$$\varphi_a(y|x) = \mathbb{E}(Y^a \mathbb{I}\{Y > y\} | X = x),$$

où $\mathbb{I}\{\cdot\}$ est la fonction indicatrice.

En remarquant que $\varphi_0(y|x) = \bar{F}(y|x)$ on obtient

$$\begin{aligned} \text{RVaR}(\alpha_n|x) &= \varphi_0^{-1}(\alpha_n|x), \\ \text{RCTE}(\alpha_n|x) &= \frac{1}{\alpha_n} \varphi_1(\varphi_0^{-1}(\alpha_n|x)|x). \end{aligned}$$

Estimateur de $\varphi_a(\cdot|x)$:

Pour estimer le moment d'ordre $a \geq 0$ de Y sachant $X = x$, on utilise un estimateur à noyau défini pour $(x, y) \in \mathbb{R}^p \times \mathbb{R}$ par

$$\hat{\varphi}_{a,n}(y|x) = \frac{\sum_{i=1}^n K\left(\frac{x - X_i}{h_n}\right) Y_i^a \mathbb{I}\{Y_i > y\}}{\sum_{i=1}^n K\left(\frac{x - X_i}{h_n}\right)}.$$

- La fonction K est appelée noyau. C'est une densité bornée sur \mathbb{R}^p , de support S inclu dans la boule unité de \mathbb{R}^p .
- (h_n) est une suite non aléatoire telle que $h_n \rightarrow 0$ quand $n \rightarrow \infty$ appelée paramètre de lissage.

Estimateur de $\varphi_a^{-1}(\cdot|x)$:

Comme $\hat{\varphi}_{a,n}(\cdot|x)$ est une fonction décroissante on donc peut définir un estimateur de $\varphi_a^{-1}(\cdot|x)$ par

$$\hat{\varphi}_{a,n}^{-1}(\alpha|x).$$

(F) On suppose que la fonction de survie conditionnelle de Y sachant $X = x$ est à queue lourde et admet une fonction de densité.

Ce qui revient à supposer que

$$\forall y > 0, \quad \text{on a } \bar{F}(y|x) = y^{-1/\gamma(x)} \ell(y|x)$$

où dans ce contexte,

- $\gamma(\cdot)$ est une fonction inconnue et positive de la covariable x et sera appelée **indice de queue conditionnel** puisqu'elle contrôle la lourdeur de la queue de la loi conditionnelle de Y sachant $X = x$.
- $\ell(\cdot|x)$ est une fonction à variations lentes à l'infini. On a (à x fixé), pour tout $\lambda > 0$,

$$\lim_{y \rightarrow \infty} \frac{\ell(\lambda y|x)}{\ell(y|x)} = 1.$$

Cette hypothèse revient à supposer que la loi conditionnelle de Y sachant $X = x$ est à queue lourde.

Théorème 1

Supposons que **(F)**, soit vérifiée et que pour tout $x \in \mathbb{R}^p$ tel que $g(x) > 0$ et $0 < \gamma(x) < 1/2$ on ait $(\alpha_n)_{n \geq 1}$ telle que $\alpha_n \rightarrow 0$ et $nh_n^p \alpha_n \rightarrow \infty$ quand $n \rightarrow \infty$. alors le vecteur aléatoire

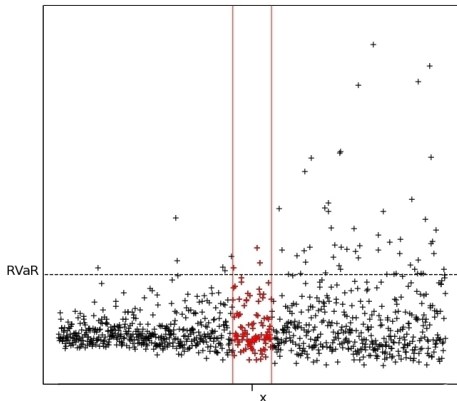
$$\sqrt{nh_n^p \alpha_n} \left\{ \left(\frac{\widehat{\text{RCTE}}_n(\alpha_n|x)}{\text{RCTE}(\alpha_n|x)} - 1 \right), \left(\frac{\widehat{\text{RVaR}}_n(\alpha_n|x)}{\text{RVaR}(\alpha_n|x)} - 1 \right) \right\}$$

est asymptotiquement Gaussien, centré, avec une matrice de variance-covariance

$$\Sigma(x) = \gamma^2(x) \frac{\|K\|_2^2}{g(x)} \begin{pmatrix} \frac{2(1-\gamma(x))}{1-2\gamma(x)} & 1 \\ 1 & 1 \end{pmatrix}.$$

- $\Sigma(x)$ est proportionnelle à $\gamma^2(x) \implies$ si $\gamma(x)$ augmente (*i.e.* plus la queue est lourde) plus la variance de nos estimateurs augmente.
- La densité $g(x)$ de la covariable est au dénominateur de la matrice de variance-covariance asymptotique \implies moins il y a de point (*i.e.* plus la densité est faible) plus la variance des estimateurs sera grande.

- $nh_n^p \alpha_n \rightarrow \infty$: condition nécessaire et suffisante pour qu'il y ait presque sûrement au moins une observation dans la région $B(x, h_n) \times [\text{RVaR}(\alpha_n|x), +\infty[$ de $\mathbb{R}^p \times \mathbb{R}$.



Si $\alpha_n = \alpha$ fixé on retrouve la condition de normalité asymptotique classique : $nh_n^p \rightarrow \infty$.

Si $h_n = h$ fixé on retrouve la condition de normalité asymptotique classique : $n\alpha_n \rightarrow \infty$.

- Dans le Théorème 1, la condition $nh_n^p \alpha_n \rightarrow \infty$ nous donne une restriction sur l'ordre des mesures de risque que l'on peut estimer.
- Cette restriction est due à l'utilisation de nos estimateurs à noyau qui ne permettent pas d'extrapoler au-delà de la plus grande observation dans la boule $B(x, h_n)$.
- Par conséquent, α_n doit correspondre à un ordre de quantile extrême se trouvant dans l'échantillon.

Definition

Soient $(\alpha_n)_{n \geq 1}$ et $(\beta_n)_{n \geq 1}$ deux suites positives telles que $\alpha_n \rightarrow 0$, $\beta_n \rightarrow 0$ et $0 < \beta_n < \alpha_n$. On peut ainsi définir un estimateur permettant d'extrapoler par

$$\widehat{\text{RCTE}}_n^W(\beta_n | x) = \widehat{\text{RCTE}}_n(\alpha_n | x) \left(\frac{\alpha_n}{\beta_n} \right)^{\hat{\gamma}_n(x)}$$

- Dans le Théorème 1, la condition $nh_n^p \alpha_n \rightarrow \infty$ nous donne une restriction sur l'ordre des mesures de risque que l'on peut estimer.
- Cette restriction est due à l'utilisation de nos estimateurs à noyau qui ne permettent pas d'extrapoler au-delà de la plus grande observation dans la boule $B(x, h_n)$.
- Par conséquent, α_n doit correspondre à un ordre de quantile extrême se trouvant dans l'échantillon.

Definition

Soient $(\alpha_n)_{n \geq 1}$ et $(\beta_n)_{n \geq 1}$ deux suites positives telles que $\alpha_n \rightarrow 0$, $\beta_n \rightarrow 0$ et $0 < \beta_n < \alpha_n$. On peut ainsi définir un estimateur permettant d'extrapoler par

$$\widehat{\text{RCTE}}_n^W(\beta_n | x) = \widehat{\text{RCTE}}_n(\alpha_n | x) \left(\frac{\alpha_n}{\beta_n} \right)^{\hat{\gamma}_n(x)}$$

où

- $\widehat{\text{RCTE}}_n(\alpha_n | x)$ est l'estimateur à noyau précédent.

- Dans le Théorème 1, la condition $nh_n^p \alpha_n \rightarrow \infty$ nous donne une restriction sur l'ordre des mesures de risque que l'on peut estimer.
- Cette restriction est due à l'utilisation de nos estimateurs à noyau qui ne permettent pas d'extrapoler au-delà de la plus grande observation dans la boule $B(x, h_n)$.
- Par conséquent, α_n doit correspondre à un ordre de quantile extrême se trouvant dans l'échantillon.

Definition

Soient $(\alpha_n)_{n \geq 1}$ et $(\beta_n)_{n \geq 1}$ deux suites positives telles que $\alpha_n \rightarrow 0$, $\beta_n \rightarrow 0$ et $0 < \beta_n < \alpha_n$. On peut ainsi définir un estimateur permettant d'extrapoler par

$$\widehat{\text{RCTE}}_n^W(\beta_n|x) = \widehat{\text{RCTE}}_n(\alpha_n|x) \left(\frac{\alpha_n}{\beta_n} \right)^{\hat{\gamma}_n(x)}$$

où

- $\widehat{\text{RCTE}}_n(\alpha_n|x)$ est l'estimateur à noyau précédent.
- $(\alpha_n/\beta_n)^{\hat{\gamma}_n(x)}$ est le terme permettant l'extrapolation.

- Dans le Théorème 1, la condition $nh_n^p \alpha_n \rightarrow \infty$ nous donne une restriction sur l'ordre des mesures de risque que l'on peut estimer.
- Cette restriction est due à l'utilisation de nos estimateurs à noyau qui ne permettent pas d'extrapoler au-delà de la plus grande observation dans la boule $B(x, h_n)$.
- Par conséquent, α_n doit correspondre à un ordre de quantile extrême se trouvant dans l'échantillon.

Definition

Soient $(\alpha_n)_{n \geq 1}$ et $(\beta_n)_{n \geq 1}$ deux suites positives telles que $\alpha_n \rightarrow 0$, $\beta_n \rightarrow 0$ et $0 < \beta_n < \alpha_n$. On peut ainsi définir un estimateur permettant d'extrapoler par

$$\widehat{\text{RCTE}}_n^W(\beta_n | x) = \widehat{\text{RCTE}}_n(\alpha_n | x) \left(\frac{\alpha_n}{\beta_n} \right)^{\hat{\gamma}_n(x)}$$

où

- $\widehat{\text{RCTE}}_n(\alpha_n | x)$ est l'estimateur à noyau précédent.
- $(\alpha_n / \beta_n)^{\hat{\gamma}_n(x)}$ est le terme permettant l'extrapolation.
- $\hat{\gamma}_n(x)$ est un estimateur de l'indice de queue conditionnel.

Théorème 2

Supposons que les hypothèses du **Théorème 1** soient vérifiées. Considérons $\hat{\gamma}_n(x)$ un estimateur de l'indice de queue conditionnel tel que

$$\sqrt{nh_n^p \alpha_n} (\hat{\gamma}_n(x) - \gamma(x)) \xrightarrow{d} \mathcal{N}(0, v^2(x)),$$

avec $v(x) > 0$. Si $(\beta_n)_{n \geq 1}$ est une suite positive telle que $\beta_n \rightarrow 0$, $\beta_n/\alpha_n \rightarrow 0$ quand $n \rightarrow \infty$, alors pour tout $x \in \mathbb{R}^p$, on a

$$\frac{\sqrt{nh_n^p \alpha_n}}{\log(\alpha_n/\beta_n)} \left(\frac{\widehat{\text{RCTE}}_n^W(\beta_n|x)}{\text{RCTE}(\beta_n|x)} - 1 \right) \xrightarrow{d} \mathcal{N}(0, v^2(x)).$$

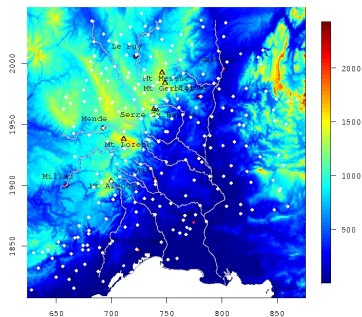
- La condition $\beta_n/\alpha_n \rightarrow 0$ permet d'extrapoler et donc de choisir un ordre β_n arbitrairement petit.
- Daouia et al. [2011] ont établi la normalité asymptotique de

$$\widehat{\text{RVaR}}_n^W(\beta_n|x) = \widehat{\text{RVaR}}_n(\alpha_n|x) \left(\frac{\alpha_n}{\beta_n} \right)^{\hat{\gamma}_n(x)}$$

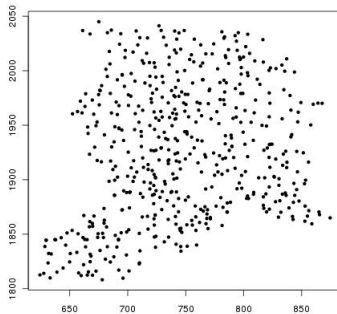
Applications à un jeu de données hydrologique

Y : hauteur de pluie journalière en mm. $X = \{\text{longitude, latitude, altitude}\}$. 1958 \implies 2000.

La région Cévennes-Vivarais



Les 523 stations d'intérêt

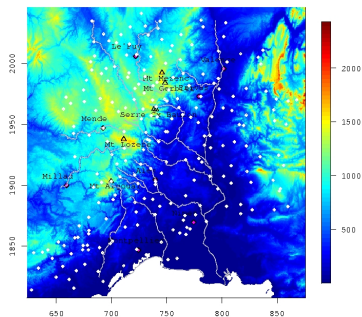


But \implies estimer $\text{RVaR}(\beta_n|x)$ et $\text{RCTE}(\beta_n|x)$ pour $\beta_n = 1/(100 * 365.25)$.

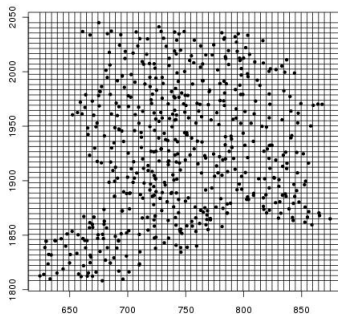
Applications à un jeu de données hydrologique

Y : hauteur de pluie journalière en mm. $X = \{\text{longitude, latitude, altitude}\}$. 1958 \implies 2000.

La région Cévennes-Vivarais



Grille 200 \times 200 points

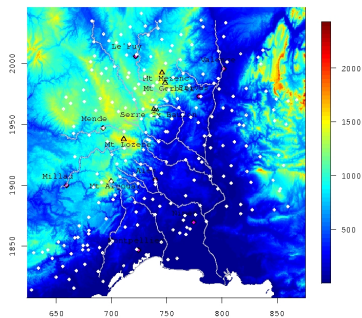


But \implies estimer $\text{RVaR}(\beta_n|x)$ et $\text{RCTE}(\beta_n|x)$ pour $\beta_n = 1/(100 * 365.25)$.

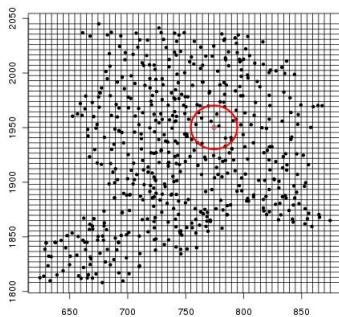
Applications à un jeu de données hydrologique

Y : hauteur de pluie journalière en mm. $X = \{\text{longitude, latitude, altitude}\}$. 1958 \implies 2000.

La région Cévennes-Vivarais



Travail dans $B(x, h_n)$



But \implies estimer $\text{RVaR}(\beta_n|x)$ et $\text{RCTE}(\beta_n|x)$ pour $\beta_n = 1/(100 * 365.25)$.

Procédure afin de choisir (h_n, α_n)

Nos estimateurs dépendent des deux paramètres de contrôle :

- h_n paramètre de lissage : problème récurrent en statistique non-paramétrique.
- α_n nombre de statistiques d'ordre utilisées : problème classique du compromis biais/variance en statistique des valeurs extrêmes.

Si α_n grand \implies biais important car on sort de la queue de distribution.

Si α_n petit \implies variance importante car on utilise peu d'observations.

\implies Mise en place d'une procédure afin de choisir (h_n, α_n) .

Notre procédure est basée sur l'estimation de la fonction $\gamma(x)$ puisqu'elle contrôle :

- la lourdeur de la queue de la loi,
- ainsi que l'extrapolation.

Le principe de notre procédure est de choisir le couple empirique $(h_{emp}, \alpha_{emp}) \in \mathcal{H} \times \mathcal{A}$ pour lequel deux estimations différentes de l'indice $\gamma(x_t)$ en chaque station t coïncident.

- Hydrologues \implies important de combiner des informations locales et régionales.

- Dans un cadre sans covariable

Soit $(\alpha_n)_{n \geq 1}$ une suite positive telle que $\alpha_n \rightarrow 0$, l'estimateur de Hill [1975] est défini par :

$$\hat{\gamma}_{n, \alpha_n} = \frac{1}{\lfloor n\alpha_n \rfloor - 1} \sum_{i=1}^{\lfloor n\alpha_n \rfloor - 1} \log(Y_{n-i+1, n}) - \log(Y_{n-\lfloor n\alpha_n \rfloor + 1, n})$$

où $\lfloor \cdot \rfloor$ est la fonction partie entière.

- Adaptation de l'estimateur de Hill à la présence d'une covariable

Soit $(\alpha_n)_{n \geq 1}$ une suite positive telle que $\alpha_n \rightarrow 0$. Un estimateur à noyau de type Hill (Gardes et Girard [2008]) est donné par

$$\hat{\gamma}_{n, \alpha_n}(x) = \frac{\sum_{j=1}^J (\log(\widehat{\text{RVaR}}_n(\tau_j \alpha_n | x)) - \log(\widehat{\text{RVaR}}_n(\tau_1 \alpha_n | x)))}{\sum_{j=1}^J \log(\tau_1 / \tau_j)},$$

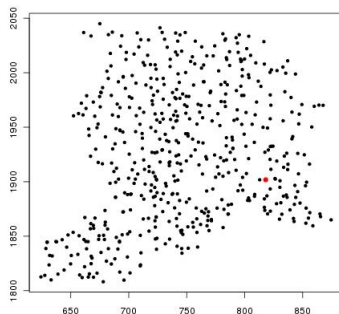
où $J \geq 1$ et $(\tau_j)_{j \geq 1}$ est une suite de poids strictement positive décroissante.

- Double boucle sur $\mathcal{H} = \{h_i; i = 1, \dots, M\}$ et sur $\mathcal{A} = \{\alpha_j; j = 1, \dots, R\}$.
- Boucle sur toutes les stations $\{x_t; t = 1, \dots, 523\}$.

Procédure de type validation croisée pour choisir h_n et α_n : Etape 1

- Double boucle sur $\mathcal{H} = \{h_i; i = 1, \dots, M\}$ et sur $\mathcal{A} = \{\alpha_j; j = 1, \dots, R\}$.
- Boucle sur toutes les stations $\{x_t; t = 1, \dots, 523\}$.

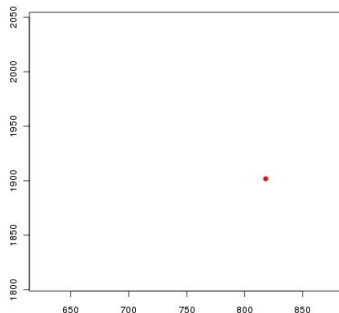
On se place en une station x_t



Procédure de type validation croisée pour choisir h_n et α_n : Etape 1

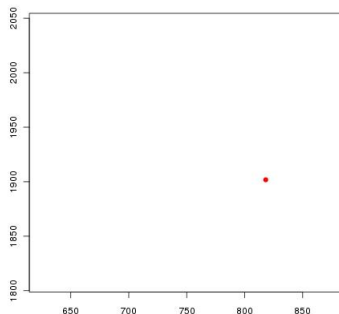
- Double boucle sur $\mathcal{H} = \{h_i; i = 1, \dots, M\}$ et sur $\mathcal{A} = \{\alpha_j; j = 1, \dots, R\}$.
- Boucle sur toutes les stations $\{x_t; t = 1, \dots, 523\}$.

On retire toutes les autres stations



- Double boucle sur $\mathcal{H} = \{h_i; i = 1, \dots, M\}$ et sur $\mathcal{A} = \{\alpha_j; j = 1, \dots, R\}$.
- Boucle sur toutes les stations $\{x_t; t = 1, \dots, 523\}$.

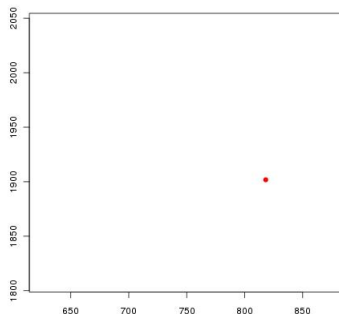
On retire toutes les autres stations



- On estime $\gamma > 0$ en utilisant l'estimateur de Hill classique.
- Il dépend seulement de α_j .
- Les α_j sont choisis de sorte que l'on soit dans la queue de distribution $\max_{j \in \{1, \dots, R\}} (\alpha_j) < 0.1$

- Double boucle sur $\mathcal{H} = \{h_i; i = 1, \dots, M\}$ et sur $\mathcal{A} = \{\alpha_j; j = 1, \dots, R\}$.
- Boucle sur toutes les stations $\{x_t; t = 1, \dots, 523\}$.

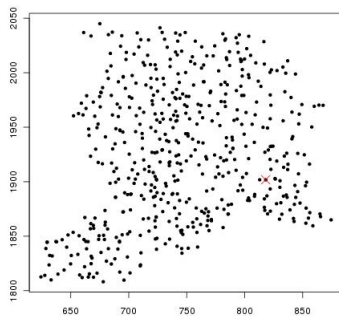
On retire toutes les autres stations



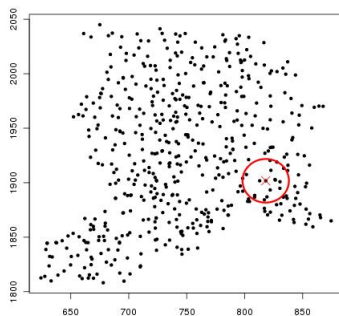
- On estime $\gamma > 0$ en utilisant l'estimateur de Hill classique.
- Il dépend seulement de α_j .
- Les α_j sont choisis de sorte que l'on soit dans la queue de distribution $\max_{j \in \{1, \dots, R\}} (\alpha_j) < 0.1$

\implies On obtient $\hat{\gamma}_{n,t,\alpha_j}$

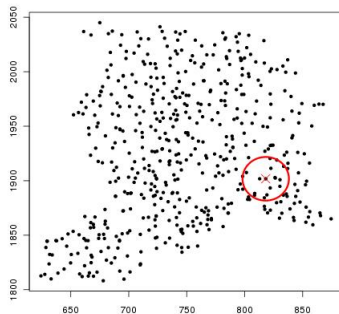
On retire cette station x_t



On travaille dans $B(x_t, h_i) \setminus \{x_t\}$

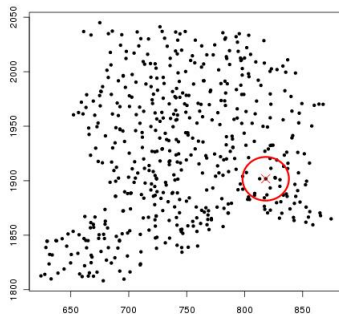


On travaille dans $B(x_t, h_i) \setminus \{x_t\}$



- On estime $\gamma(x) > 0$ en utilisant l'estimateur de Hill conditionnel.
- Il dépend de α_j et de h_i .
- Les h_i sont choisis de sorte qu'il y ait au moins une station dans $B(x_t, h_i) \setminus \{x_t\}$.

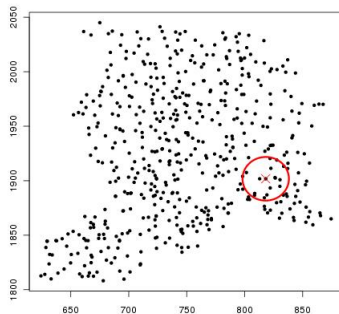
On travaille dans $B(x_t, h_i) \setminus \{x_t\}$



- On estime $\gamma(x) > 0$ en utilisant l'estimateur de Hill conditionnel.
- Il dépend de α_j et de h_i .
- Les h_i sont choisis de sorte qu'il y ait au moins une station dans $B(x_t, h_i) \setminus \{x_t\}$.

\implies On obtient $\hat{\gamma}_{n, h_i, \alpha_j}(x_t)$

On travaille dans $B(x_t, h_i) \setminus \{x_t\}$



- On estime $\gamma(x) > 0$ en utilisant l'estimateur de Hill conditionnel.
- Il dépend de α_j et de h_i .
- Les h_i sont choisis de sorte qu'il y ait au moins une station dans $B(x_t, h_i) \setminus \{x_t\}$.

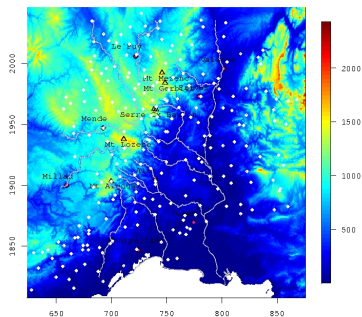
\implies On obtient $\hat{\gamma}_{n, h_i, \alpha_j}(x_t)$

$$(h_{emp}, \alpha_{emp}) = \arg \min_{(h_i, \alpha_j) \in \mathcal{H} \times \mathcal{A}} \text{mediane} \{(\hat{\gamma}_{n, t, \alpha_j} - \hat{\gamma}_{n, h_i, \alpha_j}(x_t))^2, t \in \{1, \dots, 523\}\}.$$

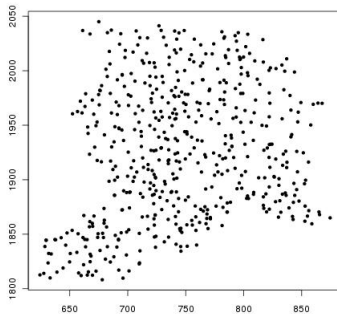
Applications à un jeu de données hydrologique

Y : hauteur de pluie journalière en mm. $X = \{\text{longitude, latitude, altitude}\}$. 1958 \implies 2000.

La région Cévennes-Vivarais



Les 523 stations d'intérêt

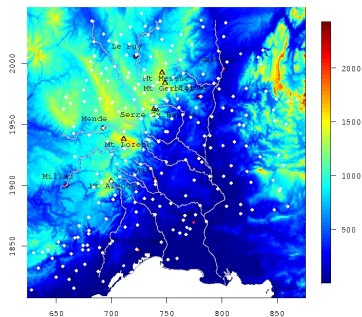


\implies Résultats de la procédure : $(h_{emp}, \alpha_{emp}) = (24, 1/(3 * 365.25))$.

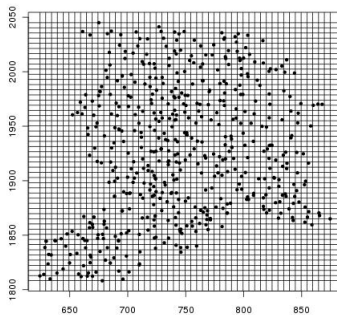
Applications à un jeu de données hydrologique

Y : hauteur de pluie journalière en mm. $X = \{\text{longitude, latitude, altitude}\}$. 1958 \implies 2000.

La région Cévennes-Vivarais



Grille 200 \times 200 points

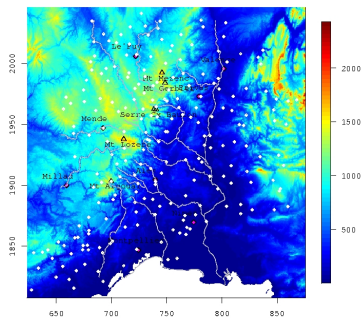


\implies Résultats de la procédure : $(h_{emp}, \alpha_{emp}) = (24, 1/(3 * 365.25))$.

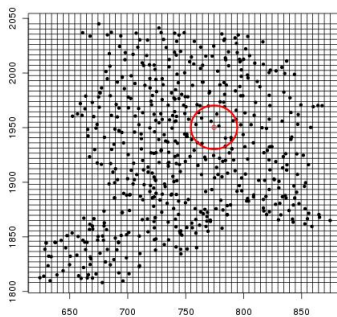
Applications à un jeu de données hydrologique

Y : hauteur de pluie journalière en mm. $X = \{\text{longitude, latitude, altitude}\}$. 1958 \implies 2000.

La région Cévennes-Vivarais



Travail dans $B(x, h_n)$

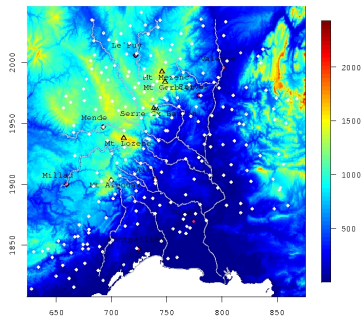


\implies Résultats de la procédure : $(h_{emp}, \alpha_{emp}) = (24, 1/(3 * 365.25))$.

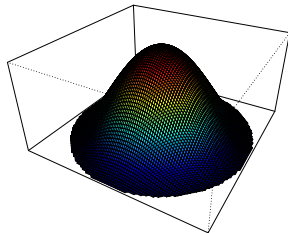
Applications à un jeu de données hydrologique

Y : hauteur de pluie journalière en mm. $X = \{\text{longitude, latitude, altitude}\}$. 1958 \implies 2000.

La région Cévennes-Vivarais

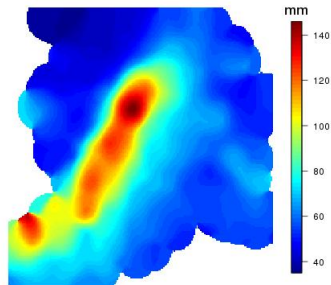


Noyau bi-quadratique

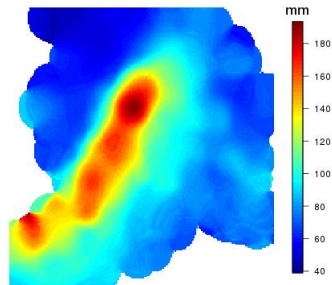


\implies Résultats de la procédure : $(h_{emp}, \alpha_{emp}) = (24, 1/(3 * 365.25))$.

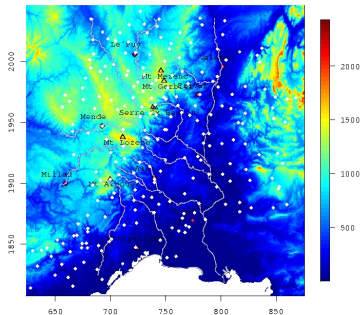
$$\widehat{\text{RVaR}}_n(1/(3 * 365.25)|x)$$



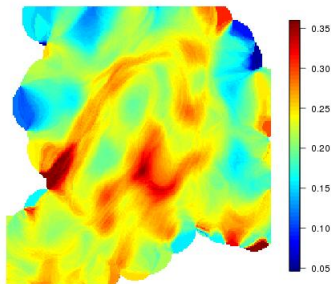
$$\widehat{\text{RCTE}}_n(1/(3 * 365.25)|x)$$



La région Cévennes-Vivarais

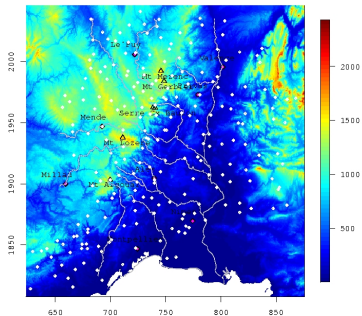


$\hat{\gamma}_{n,(1/(3*365.25))}(x)$

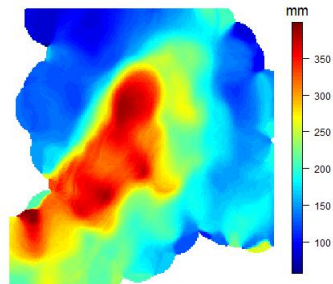


Mesures de risque extrapolées dans la région Cévennes-Vivarais

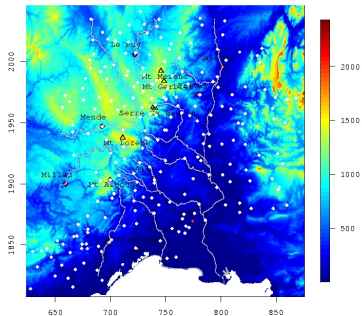
La région Cévennes-Vivarais



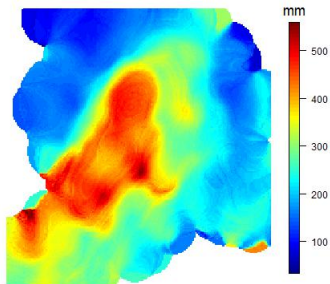
$$\widehat{\text{RVaR}}_n^W (1/(100 * 365.25)|x)$$



La région Cévennes-Vivarais



$$\widehat{RCTE}_n^W (1/(100 * 365.25)|x)$$



Points clés

- Estimer VaR, CTE, CVaR et SP dans le cas de pluies extrêmes en présence d'une covariable aléatoire dans le cas d'une loi à queue lourde.
- Extrapoler ces mesures de risque à des ordres arbitrairement petits.
- Application à un jeu de données réelles issu de l'hydrologie.

Commentaires

- + Cartes obtenues cohérentes du point de vue des hydrologues.
- + **Nouvel outil** dans la prévention des risques en hydrologie.
- + Propriétés théoriques similaires au cas univarié non extrême et/ou sans covariable.
- Fléau de la grande dimension.

Perspectives

- Estimer toutes les mesures de risque présentées dans tous les domaines d'attraction.
⇒ Cas : $\gamma(x) \in \mathbb{R}$.

Les travaux de cette présentation ont fait l'objet d'un article de recherche qui a été publié en 2014 :










El Methni, J., Gardes, L. and Girard, S. (2014). Nonparametric estimation of extreme risk measures from conditional heavy-tailed distributions, *Scandinavian Journal of Statistics*, **41**(4), 988–1012.

Sur ma page web personnelle vous trouverez les liens vers le Preprint sur Hal ou vers ma Thèse En Ligne TEL



El Methni, J. (2013). Contributions à l'estimation de quantiles extrêmes. Applications à des données environnementales.
<http://tel.archives-ouvertes.fr/tel-00924293/fr/>.

Merci de votre attention

-  Cai, J. and Tan, K. (2007). Optimal retention for a stop-loss reinsurance under the VaR and CTE risk measures, *Astin Bulletin*, **37**(1), 93–112.
-  Daouia, A., Gardes, L., Girard, S. and Lekina, A. (2011). Kernel estimators of extreme level curves, *Test*, **20**, 311–333.
-  Gardes, L. and Girard, S. (2008). A moving window approach for nonparametric estimation of the conditional tail index, *Journal of Multivariate Analysis*, **99**, 2368–2388.
-  Hill, B.M. (1975). A simple general approach to inference about the tail of a distribution, *The Annals of Statistics*, **3**, 1163–1174.
-  Embrechts, P., Kluppelberg, C. and Mikosh, T. (1997). Modelling Extremal Events, *Springer editions*.
-  Rockafellar, R.T. and Uryasev, S. (2000). Optimization of conditional value-at-risk. *Journal of Risk*, **2**, 21–42.
-  Valdez, E.A. (2005). Tail conditional variance for elliptically contoured distributions, *Belgian Actuarial Bulletin*, **5**, 26–36.