

Mélange de distributions des valeurs extrêmes généralisées

Pascal Alain Dkengne Sielenou

Travail en collaboration avec Stéphane Girard (INRIA)

Réunion de synchronisation du 26 octobre 2023

Definition (Melange des distributions de probabilités)

Une loi de probabilités est dite **loi de mélange** si sa fonction de répartition est une **moyenne pondérée algébrique** ou une **moyenne pondérée géométrique** de plusieurs fonctions de répartition.

Exemple

Considérons une suite de p fonctions de répartition $F_j, j = 1, \dots, p$ et un vecteur $\omega = (\omega_1, \dots, \omega_p) \in [0, 1]^p$ tel que $\sum_{j=1}^p \omega_j = 1$. Les lois de mélange moyennes pondérées algébriquement et géométriquement ont respectivement les formes (1) et (2) ci-dessous

$$F_S(x; \omega) = \sum_{j=1}^p \omega_j F_j(x), \quad (1)$$

$$F_P(x; \omega) = \prod_{j=1}^p F_j^{\omega_j}(x). \quad (2)$$

Definition (Mélange des distributions GEV)

- Désignons par $\omega = (\omega_1, \dots, \omega_p) \in [0, 1]^p$ un vecteur tel que $\sum_{j=1}^p \omega_j = 1$.
- Désignons par G_j une fonction de répartition de la loi GEV.
- Désignons par $\Theta = (\Theta_j, j = 1, \dots, p)$ où $\Theta_j = (\gamma_j, \sigma_j, \mu_j)$ est un vecteur des paramètres de la distribution GEV nommée G_j .

On définit le modèle de mélange G_P des lois GEV nommées G_j par

$$G_P(x; \omega, \Theta) = \prod_{j=1}^p G_j^{\omega_j}(x; \Theta_j). \quad (3)$$

Les fonctions G_j sont explicitement définies par

$$G_j(x) = G(x; \gamma_j, \sigma_j, \mu_j) = \exp \left\{ - \left[1 + \gamma_j \left(\frac{x - \mu_j}{\sigma_j} \right) \right]^{-\frac{1}{\gamma_j}} \right\}, \quad (4)$$

sur l'ensemble $\left\{ x \in \mathbb{R} : 1 + \gamma_j \left(\frac{x - \mu_j}{\sigma_j} \right) > 0 \right\}$, où $\gamma_j \neq 0$, $\mu_j \in \mathbb{R}$, $\sigma_j > 0$.

Theorem (Stabilité de la famille $\{G_{\mathbb{P}}(\cdot; \omega, \Theta)\}$)

Pour tout entier positif m et pour tout réel x , la propriété suivante est satisfaite

$$G_{\mathbb{P}}^m(x; \omega, \Theta) = \prod_{j=1}^p [G_j(x; \Theta_j(m))]^{\omega_j} = G_{\mathbb{P}}(x; \omega, \Theta(m)). \quad (5)$$

Ici, $\Theta(m) = (\Theta_j(m), j = 1, \dots, p)$ où $\Theta_j(m) = (\gamma_j(m), \sigma_j(m), \mu_j(m))$

avec $\gamma_j(m) = \gamma_j$, $\sigma_j(m) = \sigma_j m^{\gamma_j}$, $\mu_j(m) = \mu_j + \sigma_j \left(\frac{m^{\gamma_j} - 1}{\gamma_j} \right)$.

La propriété (5) montre que si la loi d'une v.a. X appartient à la famille des probabilités $\{G_{\mathbb{P}}(\cdot; \omega, \Theta)\}$, alors la loi du maximum de m copies indépendantes de X appartient également à cette même famille de probabilités.

Distribution des extrêmes et mélange des lois GEV

- Soit X une v.a. de fonction de répartition F et de borne supérieure x_F .
- Soient b_1, \dots, b_p une suite de p entiers positifs suffisamment grands.
- On suppose que pour tout $j = 1, \dots, p$ et pour toute grande valeur $x \in \mathbb{R}$, l'équivalence suivante est satisfaite

$$(\mathbb{P}\{X \leq x\})^{b_j} = (F(x))^{b_j} \sim G_j(x; \Theta_j), \quad (6)$$

où G_j est une distribution GEV de paramètre $\Theta_j = (\gamma_j, \sigma_j, \mu_j) \in \mathbb{R}^3$.

Alors, quel que soit le vecteur $\omega = (\omega_1, \dots, \omega_p) \in [0, 1]^p$ tel que $\sum_{j=1}^p \omega_j = 1$ et pour toute grande valeur $x \in \mathbb{R}$, on peut faire l'approximation suivante

$$\mathbb{P}\{X \leq x\} = F(x) \sim \prod_{j=1}^p [G_j(x; \Theta_j(b_j))]^{\omega_j} = G_{\mathbb{P}}(x; \omega, \Theta(b)). \quad (7)$$

Ici, $b = (b_1, \dots, b_p)$, $\Theta(b) = (\Theta_j(b_j), j = 1, \dots, p)$ où $\Theta_j(b_j) = (\gamma_j(b_j), \sigma_j(b_j), \mu_j(b_j))$

avec $\gamma_j(b_j) = \gamma_j$, $\sigma_j(b_j) = \sigma_j b_j^{-\gamma_j}$, $\mu_j(b_j) = \mu_j + \sigma_j \left(\frac{b_j^{-\gamma_j} - 1}{\gamma_j} \right)$.

Modélisation des valeurs extrêmes (1/2)

Soit $\mathcal{X} = (x_1, \dots, x_n)$ un échantillon d'une v.a. X de distribution de probabilités F .

Estimation du paramètre \ominus

- 1 Soit $b = \{b_j \in \mathbb{N}^*, j = 1, \dots, p\}$ un ensemble de tailles de blocs assez grandes.
- 2 Pour chaque taille de blocs $b_j \in b$, partitionner l'échantillon \mathcal{X} en $n(b_j) = \lfloor n/b_j \rfloor$ blocs disjoints contenant b_j observations consécutives.
- 3 Désignons par $\mathbf{z}_{b_j} = (z_{b_j,1}, \dots, z_{b_j,n(b_j)})$ l'échantillon des maximums où $z_{b_j,i}$ est le maximum des observations du i -th bloc de taille b_j .
- 4 Soit $\widehat{\Theta}_j = (\widehat{\gamma}_j, \widehat{\sigma}_j, \widehat{\mu}_j)$ les paramètres de la loi GEV nommée G_j estimés sur l'échantillon des maximums \mathbf{z}_{b_j} .
- 5 Pour des grandes valeurs de $x \in \mathbb{R}$, la formule (7) permet de faire l'approximation $\mathbb{P}\{X \leq x\} \approx G_{\mathbb{P}}(x; \omega, \widehat{\Theta}(b)) = \prod_{j=1}^p [G_j(x; \widehat{\Theta}_j(b_j))]^{\omega_j}$, où les composantes du vecteur $\widehat{\Theta}_j(b_j) = (\widehat{\gamma}_j(b_j), \widehat{\sigma}_j(b_j), \widehat{\mu}_j(b_j))$ constituant le paramètre $\widehat{\Theta}(b)$ s'écrivent
$$\widehat{\gamma}_j(b_j) = \widehat{\gamma}_j, \quad \widehat{\sigma}_j(b_j) = \widehat{\sigma}_j b_j^{-\widehat{\gamma}_j}, \quad \mu_j(b_j) = \widehat{\mu}_j + \widehat{\sigma}_j \left(\frac{b_j^{-\widehat{\gamma}_j} - 1}{\widehat{\gamma}_j} \right).$$

Modélisation des valeurs extrêmes (2/2)

Estimation du paramètre ω

Le vecteur ω des poids de la loi des extrêmes $G_P(x; \omega, \widehat{\Theta}(b))$ peut être estimé en résolvant le problème d'optimisation suivant

$$\widehat{\omega} = \arg \min_{\omega} \left\{ \sum_{x \in X, x > x(\alpha)} [F_{X,n}(x) - G_P(x; \omega, \widehat{\Theta}(b))]^2 \right\}, \quad (8)$$

où $F_{X,n}$ est la fonction de répartition empirique de la v.a. X estimée sur l'échantillon X de même que le quantile empirique $x(\alpha)$ d'ordre $\alpha > 0.5$.

Estimation des quantiles extrêmes

Soit $\alpha \in [0, 1]$ tel que α tend vers 0.

Le quantile extrême $x(\alpha)$ défini par $\mathbb{P}\{X > x(\alpha)\} = \alpha$ peut être estimé par une quantité $\widehat{x}(\alpha)$ qui est solution numérique de l'équation $G_P(x; \widehat{\omega}, \widehat{\Theta}(b)) = 1 - \alpha$.

Conclusion : *Ce travail explique comment combiner plusieurs modèles GEV pour obtenir un modèle assez précis dans le calcul des quantiles extrêmes d'une v.a.*

Modèles de mélange des distributions GEV : Cas IID

Soit X une variable aléatoire ayant une distribution de probabilité **inconnue**.

- Soit $b = \{b_j \in \mathbb{N}^*, j = 1, \dots, p\}$ un ensemble de p tailles de blocs assez grandes.
- Soit $(\gamma_j, \sigma_j, \mu_j)$ le vecteur des paramètres de la loi GEV nommée $G_j(\cdot)$ caractérisant la distribution des maximums de b_j obs. consécutives de la v.a. X .

Definition (Lois GEV normalisées)

La loi GEV normalisée $G_j^{1/b_j}(\cdot)$ caractérisant la distribution des grandes obs. de la v.a. X est une loi GEV $G_j(\cdot; \gamma_j(b_j), \sigma_j(b_j), \mu_j(b_j))$ dont les trois paramètres sont définis par :

- $\gamma_j(b_j) = \gamma_j$ (paramètre de forme),
- $\sigma_j(b_j) = \sigma_j (1/b_j)^{\gamma_j}$ (paramètre d'échelle),
- $\mu_j(b_j) = \mu_j + \sigma_j \left(\frac{(1/b_j)^{\gamma_j} - 1}{\gamma_j} \right)$ (paramètre de position).

Modèles de mélange des distributions GEV : Cas IID

Definition (Modèle de mélange $G_P(\cdot; \omega)$)

Le modèle de mélange $G_P(\cdot; \omega)$ est défini pour tout $x \in \mathbb{R}$ par

$$G_P(x; \omega) = \prod_{j=1}^p G_j^{\omega_j}(x; \gamma_j(b_j), \sigma_j(b_j), \mu_j(b_j)), \quad (9)$$

où $\omega = (\omega_1, \dots, \omega_p) \in [0, 1]^p$ est un vecteur des poids.

Estimation des poids

$$\hat{\omega} = \arg \min_{\omega} \left\{ \sum_{x \in \mathcal{X}, x > x_{n,\alpha}} [F_{X,n}(x) - G_P(x; \omega)]^2 \right\}, \quad (10)$$

où $F_{X,n}$ est la fonction de répartition empirique de la v.a. X estimée sur un échantillon $\mathcal{X} = \{x_1, \dots, x_n\}$ de même que le quantile empirique $x_{n,\alpha}$ d'ordre $\alpha > 0.5$.

Modèles de mélange des distributions GEV : Cas IID

Definition (Modèle de mélange $G_M(\cdot; \omega_\gamma, \omega_\sigma, \omega_\mu)$)

Soit $G(\cdot)$ la fonction de répartition de la loi GEV. Le modèle de mélange $G_M(\cdot; \omega_\gamma, \omega_\sigma, \omega_\mu)$ est défini pour tout $x \in \mathbb{R}$ par

$$G_M(x; \omega_\gamma, \omega_\sigma, \omega_\mu) = G\left(x; \sum_{j=1}^p \omega_{\gamma,j} \cdot \gamma_j(b_j), \sum_{j=1}^p \omega_{\sigma,j} \cdot \sigma_j(b_j), \sum_{j=1}^p \omega_{\mu,j} \cdot \mu_j(b_j)\right), \quad (11)$$

où $\omega_\gamma = (\omega_{\gamma,1}, \dots, \omega_{\gamma,p}) \in [0, 1]^p$, $\omega_\sigma = (\omega_{\sigma,1}, \dots, \omega_{\sigma,p}) \in [0, 1]^p$ et $\omega_\mu = (\omega_{\mu,1}, \dots, \omega_{\mu,p}) \in [0, 1]^p$ sont trois vecteurs des poids.

Estimation des poids

$$(\widehat{\omega}_\gamma, \widehat{\omega}_\sigma, \widehat{\omega}_\mu) = \arg \min_{\omega_\gamma, \omega_\sigma, \omega_\mu} \left\{ \sum_{x \in \mathcal{X}, x > x_{n,\alpha}} [F_{X,n}(x) - G_M(x; \omega_\gamma, \omega_\sigma, \omega_\mu)]^2 \right\}, \quad (12)$$

où $F_{X,n}$ est la fonction de répartition empirique de la v.a. X estimée sur un échantillon $\mathcal{X} = \{x_1, \dots, x_n\}$ de même que le quantile empirique $x_{n,\alpha}$ d'ordre $\alpha > 0.5$.

Mélange des distributions GEV : Cas Stationnaire

Soit X_t , $t = 1, 2, \dots$ une série temporelle **stationnaire**.

- Soit $b = \{b_j \in \mathbb{N}^*, j = 1, \dots, p\}$ un ensemble de p tailles de blocs assez grandes.
- Soit $(\gamma_j, \sigma_j, \mu_j)$ le vecteur des paramètres de la loi GEV nommée $G_j(\cdot)$ caractérisant la distribution des maximums de b_j obs. consécutives de la série X_t .
- Soit $\theta_j \in [0, 1]$ l'**indice extremal** associé au seuil défini par le quantile empirique $x_{n,1/b_j}$ d'ordre $1/b_j$ estimé sur une séquence $X_n = \{x_1, \dots, x_n\}$ de la série X_t .

Indice extremal

Un indice extremal θ quantifie le degré de dépendance entre l'occurrence des valeurs extrêmes consécutives. Cette dépendance est **forte** lorsque θ tend vers 0 et **faible** lorsque θ tend vers 1 .

Définition (Lois GEV normalisées)

La loi GEV normalisée $G_j^{\theta_j/b_j}(\cdot)$ caractérisant la distribution des grandes obs. de X_t est une loi GEV $G_j(\cdot; \gamma_j(b_j, \theta_j), \sigma_j(b_j, \theta_j), \mu_j(b_j, \theta_j))$ dont les trois paramètres sont définis par : $\gamma_j(b_j, \theta_j) = \gamma_j$, $\sigma_j(b_j, \theta_j) = \sigma_j (\theta_j/b_j)^{\gamma_j}$, $\mu_j(b_j, \theta_j) = \mu_j + \sigma_j \left(\frac{(\theta_j/b_j)^{\gamma_j} - 1}{\gamma_j} \right)$.

Mélange des distributions GEV : Cas Stationnaire

Definition (Modèle de mélange $G_P(\cdot; \omega)$)

Le modèle de mélange $G_P(\cdot; \omega)$ est défini pour tout $x \in \mathbb{R}$ par

$$G_P(x; \omega) = \prod_{j=1}^p G_j^{\omega_j}(x; \gamma_j(b_j, \theta_j), \sigma_j(b_j, \theta_j), \mu_j(b_j, \theta_j)), \quad (13)$$

où $\omega = (\omega_1, \dots, \omega_p) \in [0, 1]^p$ est un vecteur des poids.

Definition (Modèle de mélange $G_M(\cdot; \omega_\gamma, \omega_\sigma, \omega_\mu)$)

Soit $G(\cdot)$ la fonction de répartition de la loi GEV. Le modèle de mélange $G_M(\cdot; \omega_\gamma, \omega_\sigma, \omega_\mu)$ est défini pour tout $x \in \mathbb{R}$ par

$$G_M(x; \omega_\gamma, \omega_\sigma, \omega_\mu) = G\left(x; \sum_{j=1}^p \omega_{\gamma,j} \cdot \gamma_j(b_j, \theta_j), \sum_{j=1}^p \omega_{\sigma,j} \cdot \sigma_j(b_j, \theta_j), \sum_{j=1}^p \omega_{\mu,j} \cdot \mu_j(b_j, \theta_j)\right), \quad (14)$$

où $\omega_\gamma = (\omega_{\gamma,1}, \dots, \omega_{\gamma,p}) \in [0, 1]^p$, $\omega_\sigma = (\omega_{\sigma,1}, \dots, \omega_{\sigma,p}) \in [0, 1]^p$ et $\omega_\mu = (\omega_{\mu,1}, \dots, \omega_{\mu,p}) \in [0, 1]^p$ sont trois vecteurs des poids.

Mélange des distributions GEV : Cas Non-Stationnaire

Soit X_t , $t = 1, 2, \dots$ une série temporelle **non-stationnaire**.

Soit $Y_t = (Y_{1,t}, \dots, Y_{q,t})$ une série temporelle de q covariables pour la série X_t .

Soit x_1, \dots, x_n une séquence de n obs. de la série X_t .

On suppose que chaque obs. x_ℓ est associée à un vecteur de q cov. $y_\ell = (y_{1,\ell}, \dots, y_{q,\ell})$.

- Soit $b = \{b_j \in \mathbb{N}^*, j = 1, \dots, p\}$ un ensemble de p tailles de blocs assez grandes.
- Soit $(\gamma_j(y_t), \sigma_j(y_t), \mu_j(y_t))$ le vecteur des paramètres de la loi GEV nommée $G_j(\cdot | Y_t = y_t)$ caractérisant la distribution conditionnelle des maximums de b_j obs. consécutives de la série X_t .
- Soit $\theta_j \in [0, 1]$ l'**indice extremal** associé au seuil défini par le quantile empirique $x_{n,1/b_j}$ d'ordre $1/b_j$ estimé sur une séquence $X_n = \{x_1, \dots, x_n\}$ de la série X_t .

Structure des paramètres

- $\mu_j(y_t) = \mu_{0,j} + \mu_{1,j} f_1(y_t) + \dots + \mu_{q,j} f_q(y_t),$
- $\sigma_j(y_t) = \exp \left\{ \phi_{0,j} + \phi_{1,j} g_1(y_t) + \dots + \phi_{q,j} g_q(y_t) \right\},$
- $\gamma_j(y_t) = \gamma_{0,j} + \gamma_{1,j} h_1(y_t) + \dots + \gamma_{q,j} h_q(y_t),$

où f_ℓ, g_ℓ, h_ℓ sont des fonctions continues de supports dans \mathbb{R}^q et à valeurs dans \mathbb{R} .

Definition (Lois GEV normalisées)

Soit $y_t = (y_{1,t}, \dots, y_{q,t}) \in \mathbb{R}^q$ un vecteur constitué des valeurs potentielles des q covariables associées à la série X_t . La loi GEV normalisée $G_j^{\theta_j/b_j}(\cdot | Y_t = y_t)$ caractérisant la distribution conditionnelle des grandes obs. de X_t est la loi GEV

$G_j(\cdot; \gamma_j(b_j, \theta_j, y_t), \sigma_j(b_j, \theta_j, y_t), \mu_j(b_j, \theta_j, y_t))$ dont les trois paramètres sont définis par :

- $\gamma_j(b_j, \theta_j, y_t) = \gamma_j(y_t)$ (paramètre de forme),
- $\sigma_j(b_j, \theta_j, y_t) = \sigma_j(y_t) \cdot (\theta_j/b_j)^{\gamma_j(y_t)}$ (paramètre d'échelle),
- $\mu_j(b_j, \theta_j, y_t) = \mu_j(y_t) + \sigma_j(y_t) \cdot \left(\frac{(\theta_j/b_j)^{\gamma_j(y_t)} - 1}{\gamma_j(y_t)} \right)$ (paramètre de position).

Mélange des distributions GEV : Cas Non-Stationnaire

Definition (Modèle de mélange $G_P(\cdot | Y_t = y_t; \omega(y_t))$)

Soit $y_t = (y_{1,t}, \dots, y_{q,t}) \in \mathbb{R}^q$ un vecteur constitué des valeurs potentielles des q covariables associées à la série X_t . Le modèle de mélange $G_P(\cdot | Y_t = y_t; \omega(y_t))$ est défini pour tout $x \in \mathbb{R}$ par

$$G_P(x | Y_t = y_t; \omega(y_t)) = \prod_{j=1}^p G_j^{\omega_j(y_t)}(x; \gamma_j(b_j, \theta_j, y_t), \sigma_j(b_j, \theta_j, y_t), \mu_j(b_j, \theta_j, y_t)), \quad (15)$$

où $\omega(y_t) = (\omega_1(y_t), \dots, \omega_p(y_t)) \in [0, 1]^p$ est un vecteur des poids.

Estimation des poids $\omega_j(y_t)$

$$\widehat{\omega}(y_t) = \arg \min_{\omega} \left\{ \sum_{x \in \mathcal{X}, x > x_{n,\alpha}} [F_{X,n}(x | Y_t \in \mathcal{V}(y_t, k)) - G_P(x | Y_t = y_t; \omega(y_t))]^2 \right\}, \quad (16)$$

où $F_{X,n}$ est la fonction de répartition empirique de la v.a. X estimée sur un échantillon $\mathcal{X} = \{x_1, \dots, x_n\}$ de même que le quantile empirique $x_{n,\alpha}$ d'ordre $\alpha > 0.5$.

Mélange des distributions GEV : Cas Non-Stationnaire

Definition (Modèle de mélange $G_M(\cdot | Y_t = y_t; \omega_\gamma(y_t), \omega_\sigma(y_t), \omega_\mu(y_t))$)

Soit $G(\cdot)$ la fonction de répartition de la loi GEV. Soit $y_t = (y_{1,t}, \dots, y_{q,t}) \in \mathbb{R}^q$ un vecteur constitué des valeurs potentielles des q covariables associées à la série X_t . Le modèle de mélange $G_M(\cdot | Y_t = y_t; \omega_\gamma(y_t), \omega_\sigma(y_t), \omega_\mu(y_t))$ est défini pour tout $x \in \mathbb{R}$ par

$$G_M(x | Y_t = y_t; \omega_\gamma(y_t), \omega_\sigma(y_t), \omega_\mu(y_t)) = G\left(x; \sum_{j=1}^p \omega_{\gamma,j}(y_t) \cdot \gamma_j(b_j, \theta_j, y_t), \sum_{j=1}^p \omega_{\sigma,j}(y_t) \cdot \sigma_j(b_j, \theta_j, y_t), \sum_{j=1}^p \omega_{\mu,j}(y_t) \cdot \mu_j(b_j, \theta_j, y_t)\right), \quad (17)$$

où $\omega_\gamma(y_t) = (\omega_{\gamma,1}(y_t), \dots, \omega_{\gamma,p}(y_t)) \in [0, 1]^p$, $\omega_\sigma(y_t) = (\omega_{\sigma,1}(y_t), \dots, \omega_{\sigma,p}(y_t)) \in [0, 1]^p$ et $\omega_\mu(y_t) = (\omega_{\mu,1}(y_t), \dots, \omega_{\mu,p}(y_t)) \in [0, 1]^p$ sont trois vecteurs des poids.

Mélange des distributions GEV : Cas Non-Stationnaire

Estimation des poids $\omega_\gamma(y_t)$, $\omega_\sigma(y_t)$, $\omega_\mu(y_t)$

$$\begin{aligned} (\widehat{\omega}_\gamma(y_t), \widehat{\omega}_\sigma(y_t), \widehat{\omega}_\mu(y_t)) &= \arg \min_{\omega_\gamma, \omega_\sigma, \omega_\mu} \left\{ \sum_{x \in \mathcal{X}, x > x_{n,\alpha}} [F_{X,n}(x) - \right. \\ &\quad \left. G_M(x | Y_t = y_t; \omega_\gamma(y_t), \omega_\sigma(y_t), \omega_\mu(y_t))]^2 \right\}, \quad (18) \end{aligned}$$

où $F_{X,n}$ est la fonction de répartition empirique de la v.a. X estimée sur un échantillon $\mathcal{X} = \{x_1, \dots, x_n\}$ de même que le quantile empirique $x_{n,\alpha}$ d'ordre $\alpha > 0.5$.