

LABORATOIRE DES SCIENCES DU CLIMAT  
ET DE L'ENVIRONNEMENT

---

# Upscaling et Downscaling dans la modélisation hydrologique

---

*Auteur*

MATHIS DERONZIER  
MINES SAINT-ÉTIENNE

*Maîtres de stage*

EMMANUEL MOUCHE  
C.E.A.  
MATHIEU VRAC  
C.E.A.

STAGE DE RECHERCHE DE MASTER 2

Avril-Septembre  
2021

# Table des matières

<b>1</b>	<b>Downscaling</b>	<b>3</b>
1.1	Introduction à la problématique du downscaling . . . . .	3
1.2	Cumulative Distribution Function transform (CDFt) . . . . .	4
1.2.1	Quantile-Quantile . . . . .	4
1.2.2	CDFt . . . . .	5
1.3	Transport optimal . . . . .	6
1.3.1	Problématique . . . . .	6
1.3.2	Résolution du problème dans le cas fini et downscaling . . . . .	7
1.4	Analyse des résultats obtenus par downscaling . . . . .	7
1.4.1	Quelques outils mathématiques . . . . .	7
1.4.2	distance de Kolmogorov-Smirnov . . . . .	8
1.4.3	distance de Cramér-von Mises . . . . .	8
<b>2</b>	<b>Upscaling des modèles hydrologiques</b>	<b>9</b>
2.1	Généralités sur les interactions atmosphère - surface continentale - sol . . . . .	9
2.1.1	Les écoulements . . . . .	10
2.1.2	Transferts d'eau entre le sol et l'atmosphère . . . . .	11
2.2	Les concepts hydrologiques . . . . .	12
2.2.1	Quelques définitions . . . . .	12
2.2.2	Les équations pour modéliser l'écoulement . . . . .	12
2.2.3	La loi de Darcy et l'équation de diffusivité en milieu poreux . . . . .	13
2.3	Les modèles . . . . .	14
2.3.1	Modèle de surface continentale Orchidée . . . . .	14
<b>3</b>	<b>Prédictions climatiques et modélisation hydrologiques</b>	<b>14</b>
3.1	Analyse de la structure spatiale des données NARRs . . . . .	15
3.2	Présentation des résultats de prédictions climatiques . . . . .	19
3.2.1	La dégradation . . . . .	19
3.3	Les résultats de la dégradation et du downscaling sur les données . . . . .	19
3.3.1	tracé des boxplot . . . . .	19
3.3.2	affichage des données prédites en fonction des données réelles . . . . .	19
3.3.3	Différence analyse Root Mean Square error et distance de Cramér-von Mises . . . . .	21
<b>4</b>	<b>Indexes</b>	<b>23</b>
4.1	Indexe 1 : La statistique de Cramér-von Mises . . . . .	23
4.2	Indexe 2 : Projection conique conforme de Lambert . . . . .	24
4.3	Indexe 3 : Classification des populations de débit . . . . .	26

# 1 Downscaling

Le downscaling (voir par exemple Vrac et al. (2012) et Ayar et al. (2016)) est une méthode statistique utilisée dans les sciences du climat permettant d'améliorer les prédictions des modèles climatiques. À partir des données obtenues par un modèle climatique (modèle de circulation général, modèle climatique régional) et des données observées, on cherche à observer et corriger les biais systématiques introduits par les modèles. Le nom "downscaling" vient du domaine d'application de cette méthode. On passe d'un modèle à grande échelle à des données observées à petite échelle. Cette méthode est très utile dans la pratique où l'on a des simulateurs donnant des informations sur des maillage de grande distance de grille ( $\sim 200km$ ). Dans notre cas, nous utilisons le downscaling pour prédire les variables climatiques de *précipitation* et d'*évapotranspiration* sur le bassin du **Little Washita** un bassin d'une centaine de kilomètres carrés. Nous avons testé les principales méthodes de downscaling à partir de deux jeux de données : les données NARRs (de taille de grille de  $30 \times 30km^2$ ) ainsi que celles de l'IPSL (de taille de grille de  $200 \times 200km^2$ ), nous décrivons ce travail dans la section 3.

Pour formuler rigoureusement l'approche du downscaling nous introduisons des hypothèses communément admises dans les sciences du climat. On suppose que les variables étudiées sont des variables aléatoires réelles dépendantes du temps et de l'espace. On appelle  $\mathcal{M}(\Omega, \mathbb{R})$  l'espace des variables aléatoires réelles et  $\mathcal{S}(\mathbb{R}^3)$  la sphère unité dans  $\mathbb{R}^3$ , on suppose de plus que l'on peut faire correspondre chaque point de la terre à un point de la sphère unité. Nous travaillerons par la suite sur la sphère unité que l'on considérera être la terre.

**Définition 1.** Pour une variable quantitative  $V$  à valeur dans  $\mathbb{R}$ , on appelle  $\mathcal{T}_V$  la fonction donnant les valeurs réelles de cette variable sur la terre à un moment donné, formellement (en considérant la terre comme une sphère  $\mathcal{S}(\mathbb{R}^3)$ ) nous avons

$$\mathcal{T}_V : \mathbb{R}_+ \times \mathcal{S}(\mathbb{R}^3) \rightarrow \mathcal{M}(\Omega, \mathbb{R}). \quad (1)$$

Alors,  $\mathcal{T}_V(t, x)$  est la valeur de la variable au temps  $t$  au point de coordonnée  $x$  sur terre.

**Définition 2.** On appelle simulateur de variable quantitative  $V$  à valeur dans  $\mathbb{R}$ , une fonction  $S_V$  satisfaisant :

$$S_V : \mathbb{R}_+ \times \mathcal{S}(\mathbb{R}^3) \rightarrow \mathcal{M}(\Omega, \mathbb{R}). \quad (2)$$

On peut alors estimer la qualité des simulations en mesurant une distance entre la réalisation  $\mathcal{T}_V([0, T])$  et celle de  $S_V([0, T])$ . Le travail du downscaling est de trouver des transformations sur les variables aléatoires  $S_V(t, x)$  pour minimiser ces distances. Les géostatistiques essaient d'étudier la structure de covariance entre les variables aléatoires pour améliorer les résultats des prédictions (voir par exemple Lindgren et al. (2011)).

## 1.1 Introduction à la problématique du downscaling

Le downscaling consiste donc à effectuer des transformations sur des variables aléatoire réelles. Pour étudier ces variables aléatoires et construire les transformation on commence par introduire ici les notions de **fonction de répartition** et de **fonction de répartition empirique** ainsi que celle de **fonction de densité**. Ces notions sont centrales dans les méthodes de downscaling que nous allons expliquer.

**Définition 3.** Soit  $X$  une variable aléatoire réelle, on appelle **fonction de répartition de  $X$** ,  $\mathcal{F}_X : \mathbb{R} \rightarrow [0, 1]$  la fonction vérifiant

$$\mathcal{F}_X(x) = P(X \leq x). \quad (3)$$

**Définition 4.** Soient  $X_1, X_2, \dots, X_n$ ,  $n$  réalisations indépendantes d'une variable aléatoire réelle  $X$ , on appelle **fonction de répartition empirique de  $X$** , la fonction  $\mathcal{F}_n : \mathbb{R} \rightarrow [0, 1]$  définie par

$$\mathcal{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[X_i, +\infty)}(x). \quad (4)$$

**Définition 5.** Soit  $X$  une variable aléatoire réelle, on note  $f_X$  la **fonction de densité de  $X$** ,  $f_X : \mathbb{R} \rightarrow [0, \infty[$  la fonction vérifiant

$$f_X(x) = \lim_{t \rightarrow 0} \frac{P(x \leq X \leq x+t)}{t}. \quad (5)$$

Remarquons qu'une variable aléatoire ne possède pas nécessairement de fonction de répartition (notamment les variables aléatoires à valeurs discrètes), on peut cependant les étudier dans la théorie des distributions.

Commençons par établir notre problématique dans le cas le plus simple où l'on cherche à prédire une variable  $V$  dans l'avenir alors que nous connaissons ses réalisations dans le passé à un endroit donné  $x$ . On peut alors considérer deux processus aléatoires à valeurs dans  $\mathbb{R}$ ,  $(X_t)_{t \in \mathbb{N}} = (\mathcal{S}_V(t, x))_{t \in \mathbb{N}}$  et  $(Y_t)_{t \in \mathbb{N}} = (\mathcal{T}_V(t, x))_{t \in \mathbb{N}}$ .

La problématique à laquelle nous cherchons de répondre est la suivante : connaissant  $X_1, X_2, \dots, X_n$  et  $Y_1, Y_2, \dots, Y_n$  les réalisations jusqu'au temps  $n$  ainsi que  $X_{i_1}, X_{i_2}, \dots, X_{i_m}$ , on cherche une fonction  $G : \mathbb{R} \rightarrow \mathbb{R}$  telle que les tirages  $G(X_{i_1}), \dots, G(X_{i_m})$  et  $Y_{i_1}, \dots, Y_{i_m}$  soient proches du point de vu de leur loi (nous éclaircirons ce point par la suite dans la section 1.4).

Autrement dit, en appelant  $X$  et  $Y$  les réalisations de  $(X_t)_{t \in \mathbb{N}}$  et  $(Y_t)_{t \in \mathbb{N}}$  sur  $\{1, \dots, n\}$  et  $X'$  et  $Y'$  les réalisations de  $(X_t)_{t \in \mathbb{N}}$  et de  $(Y_t)_{t \in \mathbb{N}}$  sur  $\{i_1, \dots, i_m\}$ , on cherche à définir  $G_{X,Y}$  à partir de  $X, Y$  tel que  $G_{X,Y}$  minimise

$$d(\mathcal{F}_{G_{X,Y}(X')}, \mathcal{F}_{Y'}),$$

où  $d$  est une distance définie sur les fonctions. Nous voulons aussi que  $G_{X,Y}$  respecte certaines propriétés. Une des propriété qui nous intéresse est celle de la consistance de notre transformation.

**Définition 6.** Soient  $X$  et  $Y$  deux variables aléatoires réelles et  $G_{X,Y} : \mathbb{R} \rightarrow \mathbb{R}$  une transformation, on dit que  $G_{X,Y}$  est **consistante vis à vis de  $X$  et de  $Y$**  si elle vérifie

$$\mathcal{F}_{G_{X,Y}(X)} = \mathcal{F}_Y. \quad (6)$$

Dans la suite les transformations que nous considérerons satisferont toujours l'équation (6).

## 1.2 Cumulative Distribution Function transform (CDFt)

Nous allons ici présenter l'algorithme principalement étudié et utilisé dans ce stage, l'algorithme CDF-t. Nous commencerons par présenter l'algorithme de quantile-quantile permettant de comprendre l'esprit des transformations  $G$  affectées aux processus aléatoires. Puis nous décrirons l'algorithme de CDFt-t.

### 1.2.1 Quantile-Quantile

Le quantile-quantile consiste simplement à définir  $G_{X,Y}$  la transformation permettant de passer de la fonction de distribution de  $X$  à celle de  $Y$ .

**Proposition 1.** Soit  $X$  et  $Y$  deux variables aléatoires réelles ayant des fonctions de répartition  $\mathcal{F}_X$  et  $\mathcal{F}_Y$  continues, alors  $\mathcal{F}_Y^{-1}(\mathcal{F}_X(X))$  et  $Y$  suivent la même loi.

*Démonstration.* Montrons que  $\mathcal{F}_Y^{-1} \circ \mathcal{F}_X(X)$  et  $Y$  possède la même fonction de densité.

$$\mathcal{F}_{\mathcal{F}_Y^{-1} \circ \mathcal{F}_X(X)}(y) = \mathbb{P}(\mathcal{F}_Y^{-1}(\mathcal{F}_X(X)) \leq y) = \mathbb{P}(\mathcal{F}_X(X) \leq \mathcal{F}_Y(y)),$$

comme  $\mathcal{F}_X(X)$  suit une loi uniforme sur  $[0, 1]$  si  $\mathcal{F}_X$  est continue cette égalité se réécrit

$$= \mathbb{P}(\mathcal{U}(0, 1) \leq \mathcal{F}_Y(y)) = \mathcal{F}_Y(y).$$

□

Le principe de l'algorithme **quantile-quantile** est de calculer la transformation  $G = \mathcal{F}_Y^{-1} \circ \mathcal{F}_X$ . Ainsi,  $\mathcal{F}_{G(X)} = \mathcal{F}_Y$  et l'on définit alors  $G_{X,Y} = G$ . L'une des limites de cette méthode est que le support de  $f_{G(X)}$  est inclus dans celui de  $f_Y$ , alors les valeurs prises par  $G_{X,Y}(X')$  seront incluses dans le support de  $f_Y$ . Nous aimerions que  $X'$  ait aussi une influence sur le support des valeurs prises. Ce qui n'est pas le cas dans la transformation  $G$  que nous avons considérée. C'est en partie cette observation qui a motivé l'étude d'autres transformations.

## 1.2.2 CDFt

L'algorithme de **Cumulative Distribution Function transfer** (CDFt) vise à remédier au problème des bornes en appliquant des transformations sur les lois  $X$  et  $Y$ . On considère les réalisations  $X_1, \dots, X_n$  et  $Y_1, \dots, Y_n$  ainsi que  $X'_1, \dots, X'_m$ .

**CDFt avec support égale :**

On peut par exemple faire en sorte que la transformation sur les prédictions conserve le support de celles-ci en d'autres mots on veut

$$\min_{i \in I} (X'_i) = \min_{i \in I} (G_{X,Y}(X'_i)) \quad \text{et} \quad \max_{i \in I} (X'_i) = \max_{i \in I} (G_{X,Y}(X'_i)).$$

Il suffit de transformer la loi de  $Y$  sur  $[0, 1]$ , de réaliser un quantile-quantile puis d'effectuer la transformation inverse. Concrètement, nous posons

$$\tilde{Y} = \frac{(Y - \min_{i \in I} Y_i)}{\max_{i \in I} Y_i - \min_{i \in I} Y_i}.$$

Alors quelque soit  $i$  dans  $\mathcal{I}$  on a  $\tilde{Y}_i \in [0, 1]$ . En appelant  $G = \mathcal{F}_{X, \tilde{Y}}$  la transformation quantile-quantile de  $X$  à  $\tilde{Y}$  on a  $G_{X,Y}(X'_i) \in [0, 1]$ , la transformation s'imposant naturellement est alors

$$G_{X,Y} = (\max_{i \in I} Y_i - \min_{i \in I} Y_i)G + \min_{i \in I} Y_i.$$

C'est l'une des premières idées qui viennent en tête lorsqu'on parle de transformations que l'on peut appliquer sur nos variables. Notons que l'on a aussi une grosse base de données et que l'on peut aussi essayer d'utiliser des méthodes de régression linéaire ou de machine learning pour prédire les max et les min de  $Y'$  en fonction de  $X$ . Nous allons illustrer l'usage d'une de ces méthodes dans la prochaine partie 1.2.2.

**CDFt avec méthode de prédiction**

Supposons que l'on ait des fonctions estimant la variance et de la moyenne de  $Y'$  à partir de  $X'$ , s'exprimant sous la forme  $\bar{f}(X') = \bar{Y}'$  ainsi que  $\bar{g}(X') = \bar{\sigma}(Y')$  (on suppose ici  $X'$  et  $Y'$  comme des suites de variables aléatoires)<sup>1</sup>. Ces fonctions peuvent être obtenues par des méthodes de régression linéaire sur des sous échantillons des données de  $X$  et  $Y$  soit après les avoir ordonnées si  $X_i$  et  $Y_i$  sont indépendants soit en fonction de leur indice.

On aimerait alors que  $G_{X,Y}$  conserve la moyenne ainsi que la variance. Formellement, on voudrait définir  $G_{X,Y}$  telle que quelque soient  $\{i_1, \dots, i_m\}$  un ensemble d'entiers consécutifs et  $X'$  et  $Y'$  les vecteurs des variables aléatoires des  $X_{i_j}$  et  $Y_{i_j}$  sur cet ensemble, on ait :

$$\overline{G_{X,Y}(X')} = \bar{Y}', \quad (7)$$

ainsi que

$$\sigma(G_{X,Y}(X')) = \sigma(Y'), \quad (8)$$

et que  $G$  respecte la condition de consistance (6). Concrètement, nous allons faire des transformations sur les variables aléatoires pour avoir ces conditions là. En appelant  $G_{X,Y} = \mathcal{F}_Y^{-1} \circ \mathcal{F}_X$  on peut définir la transformation  $G_{X,Y,X'}$

---

1. On définit pas  $\bar{X}$  la moyenne  $X$  et par  $\sigma(X)$  son écart type

$$G_{X,Y}(X') = \overline{g(X')} \frac{\mathcal{F}_Y^{-1} \circ \mathcal{F}_X(X') - \overline{\mathcal{F}_Y^{-1} \circ \mathcal{F}_X(X')}}{\sigma(\mathcal{F}_Y^{-1} \circ \mathcal{F}_X(X'))} + \overline{f(X')}.$$

On peut vérifier facilement que  $G_{X,Y}$  ainsi défini respecte bien les propriétés (6) (7) et (8). L'idée est alors de trouver les fonctions  $f$  et  $g$  par des méthodes de regression linéaires. Remarquons aussi que lorsque la loi est bornée cela ajoute une condition supplémentaire à  $G_{X,Y}$  et l'on ne peut pas nécessairement avoir les conditions sur la moyenne et la variance. Il faut alors choisir parmi l'une des conditions (7) et (8), c'est par exemple le cas pour les précipitations qui ne peuvent être négatives.

**Note 1.** *Remarquons que l'efficacité des transformations que nous avons décrites repose en partie sur la stationnarité des lois suivies par les variables aléatoires aléatoires au cours du temps. Les questions de stationnarité ont été abordées dans Maraun (2012), Christensen et al. (2008) ainsi que Nahar et al. (2017).*

### 1.3 Transport optimal

Nous voulons à la fois prédire la précipitation et l'évapotranspiration, on peut alors considérer une seule variable aléatoire dans  $\mathbb{R}^2$ . On peut généraliser l'idée utilisée précédemment pour trouver une méthode permettant de corriger les biais statistiques introduits par les modèles de prédictions. Cette méthode a d'autant plus d'intérêt que la variable utile dans les modèles hydrologique est

$$\text{pluie entrant dans le sol} = \text{précipitation} - \text{évapotranspiration}.$$

Comme l'objectif final est de prédire des résultats hydrologiques sur le bassin du Little Washita, il semble particulièrement pertinent de considérer la loi conjointe (précipitation, évapotranspiration). La théorie généralisant cette idée est la théorie du **transport optimal**.

La problématique du transport optimal a premièrement été introduite en par Gaspard Monge en 1781 puis a été développée par Kantorovitch en 1971 et ses travaux pour l'allocation des ressources lui ont valu un prix nobel d'économie en 1975.

#### 1.3.1 Problématique

Précédemment nous avons cherché à définir une transformation  $G$  de la variable aléatoire  $X$  telle  $G(X)$  suive la même loi que  $Y$ . En considérant les fonctions de densité de  $X$  et de  $Y$  on a cherché à ce que la fonction de densité de  $G(X)$  soit la même que celle de  $Y$ . On peut considérer une fonction de densité comme une mesure sur l'espace sur lequel on travail on a alors transformé une mesure  $f_X$  en une autre mesure  $f_Y$ . Comme ces deux mesures sont de mesure totale égale à 1, on peut dire que d'une certaine manière chaque "poids" de la mesure  $f_X$  a été déplacé vers un poids de la mesure  $f_Y$ . L'idée du transport optimal est de trouver les déplacements naturels des poids d'une fonction de densité à une autre.

Nous présenterons ici la formulation établie par Kantorovitch dans les années 70 qui a l'avantage d'inclure celle de Monge. Le livre Villani (2003) donne un cours de référence internationale sur les problématiques de transport optimal.

Considérons deux fonctions de répartitions pour des variables aléatoires  $U$  et  $V$  à valeurs dans  $A$  et  $B$ , on appelle ces fonctions  $\mathcal{F}$  et  $\mathcal{G}$  et on appelle  $f$  et  $g$  leurs fonctions de densités. On cherche alors une mesure  $\pi$  sur  $A \times B$  satisfaisant

$$\int_B d\pi(x, y) = f(x), \quad \int_A d\pi(x, y) = g(y),$$

de plus on veut que  $\pi$  satisfaisant l'équation précédente minimise la quantité

$$\mathcal{I}[\pi] = \int_{A \times B} d(x, y) d\pi(x, y),$$

où  $d$  est une certaine distance définissant le coût de transport de  $x$  à  $y$ . Dans notre cas  $U$  et  $V$  sont des variables aléatoires à valeurs dans  $\mathbb{R}^2$ . Nous voyons que le choix de la distance a une influence majeure sur la mesure obtenue. Alors on peu voir  $d\pi(x, y)$  comme la quantité déplacée de  $x$  à  $y$ .

### 1.3.2 Résolution du problème dans le cas fini et downscaling

La résolution de ce problème dans le cas fini a été traité de nombreuses fois. On utilisera les idées développées dans le papier Robin et al. (2019). Appelons  $\pi \in \mathbb{R}^{m \times n}$  une matrice de transfert de poids dans le cas fini. On a  $X_1, \dots, X_n$  ainsi que  $Y_1, \dots, Y_m$  des réalisations de  $X$  et de  $Y$  et on cherche alors une matrice  $\pi$  telle que

$$\sum_{j=1}^n \pi_{i,j} = P(X = X_i) \text{ et } \sum_{i=1}^m \pi_{i,j} = P(Y = Y_j),$$

avec  $\pi$  minimisant

$$\mathcal{I}[\pi] = \sum_{i,j} d(X_i, Y_j) \pi_{i,j}.$$

Le papier utilise la norme euclidienne comme distance, l'obtention de cette solution peut se faire par un algorithme de simplexe, voir par exemple Huang and Chen (2012). Pour corriger le biais d'estimation on peut alors pour chaque  $x$  tirés récupérer (par interpolation linéaire ou méthode de krigeage) le  $\pi(x, \cdot)$ . D'après la construction de  $\pi$ , on peut alors normaliser la fonction  $y \mapsto \pi(x, y)$  et tirer aléatoirement un point selon la loi ainsi trouvée (voir Robin et al. (2019) pour plus de détails).

### 1.4 Analyse des résultats obtenus par downscaling

Concrètement nous allons faire de la validation croisée, nous allons apprendre sur 50% de nos données et faire nos prédictions. La partie à laquelle nous allons nous intéresser ici est la distance que nous utilisons pour évaluer nos prédictions. Contrairement à la manière habituelle de faire, consistant à estimer une distance entre chaque point prédit (souvent RMSE), en climatologie nous cherchons à comprendre la tendance générale. En effet, le paradigme d'évaluation en prévisions climatiques sur plusieurs années n'a pas l'ambition de prédire ponctuellement chaque prévision, mais il a pour objectif de décrire une tendance générale. On s'intéresse alors à des informations plus générales, c'est à dire que l'on travaille sur les lois de répartitions. Il faut alors réfléchir à des normes ou des distance pour évaluer la qualité de nos prédictions.

Dans notre cas nous faisons des tests non-paramétriques, c'est à dire que l'on ignore tout des lois que nous comparons. Différentes méthodes pour tester l'égalité de lois sont connues, nous n'en développerons que deux. Le mémoire Éthier (2011) donne une présentations de principaux tests statistiques permettant d'évaluer si oui ou non à partir des réalisations  $X_1, \dots, X_n, Y_1, \dots, Y_n$  de deux lois inconnues sont les mêmes.

Nous posons habituellement en statistiques deux hypothèses :

$$\mathcal{H}_0 : \mathcal{F}_X = \mathcal{F}_Y \text{ et } \mathcal{H}_1 : \mathcal{F}_X \neq \mathcal{F}_Y,$$

où les égalités sur les lois sont en norme  $L^p$ . On suppose  $\mathcal{H}_0$  et on définit une statistique sur  $\|\mathcal{F}_X - \mathcal{F}_Y\|_{L^p(\mathbb{R})}$  permettant à partir de nos observations d'accepter ou de rejeter l'hypothèse  $\mathcal{H}_0$ . Les tests de Kolmogorov-Smirnov et Cramer-von Mises utilisent à-peu-près cette idée.

#### 1.4.1 Quelques outils mathématiques

**Proposition 2.** Soient  $X_1, \dots, X_n$   $n$  réalisations d'une variable aléatoire réelle  $X$  et  $\mathcal{F}_n$  sa fonction de répartition empirique nous avons

$$E[\mathcal{F}_n] = F_X,$$

alors la fonction de répartition empirique est un estimateur sans biais de la lois de  $F$ .

*Démonstration.* C'est en effet évident puisque  $\mathbb{1}_{[X_i, +\infty)}(x)$  suit une lois de Bernoulli de paramètre  $\mathcal{F}(x)$  alors

$$E\left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[X_i, +\infty)}(x)\right] = \frac{1}{n} \sum_{i=1}^n E[\mathbb{1}_{[X_i, +\infty)}] = \mathcal{F}_X(x).$$

□

**Théorème 1.1.** (*Glivenko-Cantelli*) Soient  $\mathcal{F}_X$  et  $\mathcal{F}_n$  respectivement la fonction de répartition et la fonction de répartition empirique. Alors

$$\|\mathcal{F}_X - \mathcal{F}_n\|_\infty \xrightarrow[n \rightarrow \infty]{prob} 0 \quad (9)$$

*Démonstration.* (Cas où  $\mathcal{F}_X$  est continue) On commence par remarquer que quelque soit  $x$  dans  $\mathbb{R}$ ,

$$F_n(x) \xrightarrow[n \rightarrow \infty]{p.s.} \mathcal{F}_X(x)$$

d'après la loi forte des grands nombres et la proposition (2). Pour  $q$  dans  $\mathbb{Q}$ , on définit

$$\Omega_q = \{\omega \in \Omega \mid \lim_{n \rightarrow \infty} \mathcal{F}_n(q) = \mathcal{F}_X(q)\},$$

d'après ce que nous avons dit, sa mesure pour la probabilité  $P(\Omega_q) = 1$  comme  $\mathbb{Q}$  est dénombrable nous avons

$$P\left(\bigcap_{q \in \mathbb{Q}} \Omega_q\right) = 1.$$

Alors, comme  $\mathbb{Q}$  est dense dans  $\mathbb{R}$  et que  $\mathcal{F}_X$  et les  $(\mathcal{F}_n)_{n \in \mathbb{N}}$  sont continues on peut assurer que

$$\|\mathcal{F}_X - \mathcal{F}_n\|_\infty \xrightarrow[n \rightarrow \infty]{prob} 0.$$

□

Le cas où  $\mathcal{F}_X$  n'est pas continue est géré par Durrett (2019) (ex 7.2 chap 1). On voit d'après le théorème 1.1 que la fonction de répartition empirique est le bon estimateur de la fonction de répartition.

#### 1.4.2 distance de Kolmogorov-Smirnov

Le test d'ajustement de Kolmogorov-Smirnov est l'un des plus utilisé pour tester l'égalité de deux lois de probabilités. Dans le contexte de l'égalité de lois de probabilité, la statistique de test est

$$K_{n,m} = \sqrt{\frac{nm}{n+m}} \|\mathcal{F} - \mathcal{F}_n\|_\infty.$$

La suite de variables aléatoires  $K_{n,m}$  converge vers une variable aléatoire  $K$  dont la fonction de survie est donnée par :

$$Q(x) = P(K > x) = \sum_{j=1}^{\infty} (-1)^{j-1} \exp(-2(jx)^2) \quad (10)$$

On peut alors l'approximer avec les premiers termes de la série pour construire le test statistique. La démonstration de ce théorème peut être trouvée dans le livre Fisz (1963)(chap 12.5). Nous voyons cependant que ce test est sensible aux données aberrantes, nous privilégierons alors le test de **Cramér-von Mises** (section 1).

#### 1.4.3 distance de Cramér-von Mises

On considère ici les deux fonctions de répartition  $\mathcal{F}_n$  et  $\mathcal{G}_m$  continues des variables aléatoires  $X$  et  $Y$ . Nous voulons tester les hypothèses

$$\mathcal{H}_0 : \mathcal{F}_n = \mathcal{G}_m \quad \text{et} \quad \mathcal{H}_1 : \mathcal{F}_n \neq \mathcal{G}_m.$$

Dans les tests proposés les statistiques sont construites à partir de deux échantillons indépendants des tests dans les cas où les variables  $X$  et  $Y$  sont indépendantes peuvent être trouvés dans Éthier (2011). On définit aussi la fonction de répartition empirique  $\mathcal{F}_n$ .



Nous avons donc  $n$  réalisations de  $X$  et  $m$  réalisations de  $Y$  de lois de répartition  $\mathcal{F}$  et  $\mathcal{G}$ . La statistique du test est définie par

$$C_{n,m} = \frac{nm}{n+m} \int_{\mathbb{R}} [\mathcal{F}_n(x) - \mathcal{G}_m(x)]^2 d\mathcal{H}_{m,n}(x), \quad (11)$$

avec,

$$\mathcal{H}_{n,m} = \frac{n}{n+m} \mathcal{F}_n + \frac{m}{n+m} \mathcal{G}_m. \quad (12)$$

$\mathcal{H}_{n,m}$  est alors la fonction de répartition empirique d'une variable aléatoire  $Z$  construite à partir des  $n+m$  réalisations indépendantes  $X_1, \dots, X_n, Y_1, \dots, Y_m$ . On peut simplement réécrire la valeur  $C_{n,m}$

$$C_{n,m} = \frac{mn}{(m+n)^2} \sum_{i=1}^{m+n} (\mathcal{F}_n(Z_i) - \mathcal{G}_m(Z_i))^2 \quad (13)$$

**Lemme 1.** *On peut simplifier cette formule en supposant que les  $(X_i)_{i \in \{1, \dots, n\}}$  et  $(Y_i)_{i \in \{1, \dots, m\}}$  sont triés on a :*

$$C_{n,m} = \frac{1}{nm(m+n)} \left[ n \sum_{i=1}^n (R_{X_i} - i)^2 + m \sum_{i=1}^m (R_{Y_i} - i)^2 \right] - \frac{4nm - 1}{6(m+n)}. \quad (14)$$

où  $R_{X_i}$  est le rang de  $X_i$  dans  $X_1, \dots, X_n, Y_1, \dots, Y_m$  autrement dit

$$R_{Z_i} = \text{Card}(\{z \in Z, z \leq Z_i\}).$$

**Note 2.** *La démonstration de cette formule se trouve dans l'indexe 1 et la formulation de cette égalité diffère de celle contenue dans le mémoire Éthier (2011)(sec 2.3.2) qui contient une erreur.*

Cette formulation a le bon goût de nous indiquer que la statistique ne dépend pas de la loi. On peut alors calculer simplement sa statistique sous l'hypothèse  $\mathcal{H}_0$  pour la loi uniforme sur  $[0, 1]$  et ainsi retrouver les quantiles présentés dans l'article de Büning (2002). Nous voyons que l'idée du test est aussi de pondérer la différence des fonctions de répartition empiriques par les observations (l'intégration selon  $\mathcal{H}_{m,n}$ ). Ainsi, si le test de Kolmogorov-Smirnov est sensible aux outliers, celui-ci l'est beaucoup moins lorsque les échantillons sont suffisamment grands.

## 2 Upscaling des modèles hydrologiques

Comme nous l'avons vu dans la précédente partie, il existe différentes échelles d'étude. L'enjeu du downscaling est de passer de la grande échelle à la petite échelle. L'objectif de l'upscaling inverse. La différence entre ces deux méthodes est que le domaine d'application du downscaling est le traitement des données alors que celui de l'upscaling est la modélisation.

L'upscaling est une réflexion sur le changement d'échelle des modèles. L'objectif est d'avoir des modèles efficaces et applicable à grande échelle. On cherche alors à simplifier les modèles tout en gardant des propriétés intéressantes. Par exemple l'upscaling des équations de Navier-Stokes sont un enjeu majeur de l'hydrologie (voir 2.2.2), on cherche à les simplifier pour pouvoir les résoudre à grandes échelles. Le temps de calcul informatique pour résoudre les équations de Navier-Stokes sur de grand domaine par toutes les méthodes explose bien que le temps de résolution algorithmique est classiquement  $O(n^3)$  où  $n$  est le nombre de mailles de notre modèle. Pour résoudre ces problèmes, des modèles simplificateurs ont vu le jour comme les modèles de colonne, modèles qui ne repose à première vue sur aucune loi physique mais qui fonctionne empiriquement.

### 2.1 Généralités sur les interactions atmosphère - surface continentale - sol

Cette section reprend brièvement les principaux mécanismes entrant en compte dans les modélisations hydrologiques. Ces mécanismes seront présentés brièvement mais des sources seront données pour les lecteurs voulant étudier plus en détails ces mécanismes. Cette section s'inspire très largement de la thèse Maquin (2016) sur les modèles hydrologiques de colonne appliqués sur le bassin du Little Washita.

### 2.1.1 Les écoulements

Les processus hydrologiques sont classiquement étudiés à l'échelle du bassin versant (voir). En hydrologie, le bassin versant est une unité géographique définie par les limites topographiques que sont les lignes de crête. L'ensemble des écoulements converge vers les dépressions, formant ainsi un réseau hydrographique qui se dirige vers le point bas du bassin versant, l'exutoire.

Il possède son équivalent en mathématiques, soit  $(E, d)$  un espace métrique et  $S : E \rightarrow E$  un endomorphisme sur  $E$ . On définit le bassin d'attraction d'un point  $a$  qu'on appelle  $B(a)$  l'ensemble des points  $x$  dans  $E$  tels que la suite  $(x_n)_{n \in \mathbb{N}} = (S^n(x))_{n \in \mathbb{N}}$  converge vers  $a$ . En considérant qu'il existe une fonction  $S$  définissant la trajectoire d'une goutte d'eau déposée au point  $x$  les deux définitions ont la même signification.

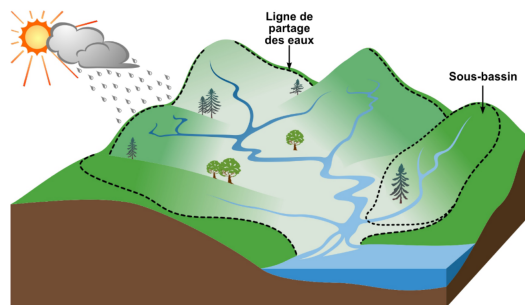


Figure 1 – Bassin versant (Source :<http://rqes-gries.ca/>).

À l'échelle du bassin versant, on distingue deux types d'écoulements : les écoulements de subsurface, les écoulements de surface.

#### Les écoulements de subsurface :

La notion d'écoulement de subsurface se rapporte à l'écoulement de l'eau dans les pores du sol. L'écoulement de subsurface dépend de plusieurs paramètres comme les caractéristiques du sol (la porosité, la perméabilité), la saturation en eau du sol la topographie et le climat (précipitation, évaporation, transpiration). Ces écoulements sont traités par les équations de la mécanique des fluides (voir section à citer) et pour plus de détails ?.

#### Les écoulements de surface :

Les écoulements de surface, aussi qualifiés de ruissellement, sont la conséquence de deux phénomènes distincts. Le ruissellement peut apparaître lorsque le sol est saturé en surface. En effet, lorsque le sol est saturé, l'eau ne peut s'y infiltrer (Cappus (1960)). Cette condition de saturation à la surface du sol peut être la conséquence d'une nappe affleurant la surface, la zone satisfaisant cette propriété est appelée zone de suintement. Cela arrive naturellement lors d'épisode pluvieux pour les nappes peu profondes. Le ruissellement peut aussi être causé par de fortes précipitations, ainsi le débit surfacique peut devenir supérieur à la quantité d'infiltration et ainsi créer un ruissellement (voir la figure). La quantité d'infiltration décroît exponentiellement lors d'évènements pluvieux (Horton (1933)).

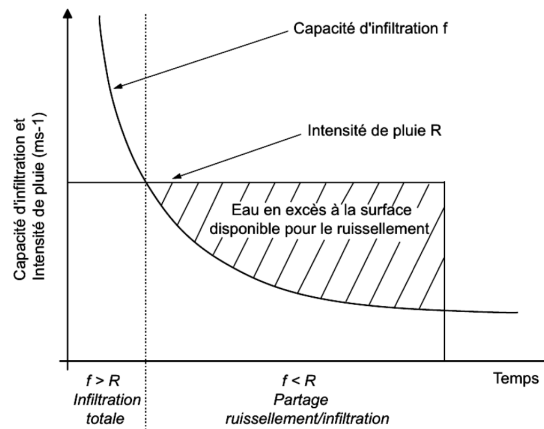


Figure 2 – Estimation du ruissellement en fonction du temps modèle de Horton

### 2.1.2 Transferts d'eau entre le sol et l'atmosphère

La végétation constitue le lien entre l'atmosphère et le sol. Les végétaux transfèrent de l'eau dans les deux sens, via les racines et la canopée. Il y a aussi des interactions directes entre le sol et l'atmosphère. Les trois principaux processus décrits sont l'évaporation du sol, la transpiration des végétaux ainsi que l'évaporation de l'eau interceptée par la canopée (mécanisme visant à conserver l'eau).

#### la transpiration :

La "transpiration" des plantes consiste en une libération de vapeur d'eau par les plantes dans l'atmosphère. Ce phénomène constitue une réponse passive à l'environnement atmosphérique dû à l'existence d'un gradient de pression positif de l'atmosphère à la canopée, on parle alors de demande atmosphérique. La description des processus d'évaporations ont été traités dans la thèse Maquin (2016).

#### Évaporation :

Sur les surfaces de sol non recouvertes de végétation (sol nu), l'eau présente dans le sol, à proximité de la surface, peut s'évaporer. Ce phénomène apparaît en présence d'un gradient de pression de vapeur d'eau entre le sol et l'atmosphère et d'un apport d'énergie. L'évaporation effective dépend de l'état hydrique de la surface du sol, l'énergie pour extraire l'eau du sol augmentant à mesure que le sol s'assèche et des propriétés conductrices du sol (voir Hillel (2003)).

#### Pertes par interception :

Lors d'un épisode pluvieux, une partie de l'eau incidente est interceptée par le feuillage. Il s'agit du phénomène dit d'interception. Cette eau présente sur la canopée peut ensuite s'évaporer directement. On désigne ce processus d'évaporation sur la canopée comme les pertes par interception. L'importance de ce flux d'eau dépend de l'ampleur du feuillage et de la capacité de stockage d'eau de la canopée, c'est-à-dire de l'épaisseur maximale de la lame d'eau par unité de surface de feuillage

#### l'évapotranspiration potentielle :

On désigne par « évapotranspiration potentielle » la quantité d'eau maximale que l'atmosphère peut extraire via les trois processus décrits précédemment. Elle correspond ainsi à la demande atmosphérique évoquée auparavant. L'évapotranspiration potentielle correspond à l'évaporation d'une surface saturée en eau. Elle dépend de paramètres atmosphériques comme l'humidité de l'air, le vent et la température. Ce

taux potentiel a la propriété de majorer la somme des flux de transpiration, d'évaporation et des pertes par interception.

## 2.2 Les concepts hydrologiques

Nous allons ici introduire les principale notions à l'étude hydrologique des sols. Plusieurs caractéristiques définissent un sol, mais avant d'étudier en détail ce qui définit un sol, il est important de comprendre que l'étude hydro d'un sol est simplement un bilan d'eau dans celui-ci. Il faut alors commencer par déterminer ce qui rentre et ce qui sort. L'estimation de ces quantités est l'objet d'étude du downscaling qui cherche à prévoir la précipitation et l'évapotranspiration (l'eau drainée par les plantes recouvrant la surface).

### 2.2.1 Quelques définitions

**Définition 7.** On appelle **porosité totale**  $\omega$  la valeur définie par

$$\omega = \frac{\text{Volume des vides}}{\text{Volume total de la roche}}. \quad (15)$$

On appelle aussi **indice des vides**  $e$  la valeur définie par

$$e = \frac{\text{Volume des vides}}{\text{Volume du solide plein}}. \quad (16)$$

On peut passer d'une formule à l'autre par la relation

$$e\omega = e - \omega,$$

Mais l'on utilise toujours la notion de porosité en hydrologie.

L'on peut trouver des méthodes de mesure de la porosité d'un sol dans l'ouvrage (De Marsily (1986)). On dit aussi que le sol n'est pas saturé lorsque l'eau n'a pas pris tout l'espace disponible, on parle alors de saturation volumique.

**Définition 8.** On parle de **saturation volumique**  $\theta$ , la saturation définie par le rapport

$$\theta = \frac{\text{Volume d'eau contenu}}{\text{Volume total}}, \quad (17)$$

on a  $0 \leq \theta \leq \omega$ . Et la **saturation volumique**  $s$

$$s = \frac{\text{Volume d'eau contenu}}{\text{Volume total des pores}}. \quad (18)$$

En fonction de la saturation volumique les échelles de temps et les forces mises en action ne sont pas les mêmes.

### 2.2.2 Les équations pour modéliser l'écoulement

On commence par rappeler les équations essentielles à la dynamique des fluides. L'équation de conservation de la matière :

$$\text{div}(\rho \vec{u}) + \frac{\partial \rho}{\partial t} = 0. \quad (19)$$

Où  $\rho$  est la masse volumique et  $\vec{u}$  le vecteur vitesse du fluide. On écrit maintenant l'équation de Navier-Stokes

$$\frac{\partial p}{\partial x^i} - \left(\zeta + \frac{\mu}{3}\right) \frac{\partial}{\partial x^i} (\text{div} \vec{u}) - \mu \nabla^2 u^i = \rho \left(F^i - \frac{du^i}{dt}\right), \quad (20)$$

$\zeta$  coefficient de viscosité du volume, (très souvent négligeable de vant  $\mu$ ) [ $ML^{-1}T^{-1}$ ],  
 $\mu$  coefficient de viscosité dynamique, [ $ML^{-1}T^{-1}$ ]  
 $\nabla^2$  le laplacien,  
 $F^i$  composante des forces à distance par unité de masse,  
 $i$  un vecteur unitaire de l'espace  $3D$ .

En milieu poreux les équations de Navier-Stokes deviennent difficilement applicables car le milieu dans lequel s'écoule le fluide dépend lui-même de l'écoulement du fluide. On donne des hypothèse simplificatrices pour résoudre l'équation de Navier-Stokes.

#### **Hypothèses simplificatrice :**

En milieu poreux on peut émettre de nombreuses hypothèses simplificatrices qui permettent de simplifier l'équation de Navier Stokes, on commence par supposer les écoulement permanents

$$\frac{\partial u^i}{\partial t} = 0,$$

on suppose aussi que le fluide est incompressible ( $\rho$  constant) alors

$$\text{div}(\rho \vec{u}) = -\frac{\partial \rho}{\partial t} = 0,$$

et finalement,

$$\text{div} \vec{u} = 0.$$

Dans ces hypothèses l'équation de Navier-Stokes devient

$$\frac{\partial p}{\partial x^i} - \mu \nabla^2 u^i - \rho F^i = 0. \quad (21)$$

Plusieurs méthodes numériques permettent de trouver des solutions à ce problèmes notamment la méthode de Galerkin ou des méthodes de différences finies voir Allaire (2005) (chap 2, chap 6).

Remarquons que ces équations ne prennent pas en compte la porosité du milieu, l'équation (19) peut être modifié en prenant en compte  $\omega$  le coefficient de porosité et un terme source  $q$  lié à la matière créant des interstices sans fluide (le coefficient est compté négativement). Alors l'équation de conservation s'écrit finalement

$$\text{div}(\vec{U}) + \frac{\partial}{\partial t}(\theta) + q = 0. \quad (22)$$

Remarquons que l'on considère la porosité  $w$  comme continue nous étudions des éléments de longueurs  $dx$  suffisamment petite pour que les équations de continuité (??) et (22) soient considérée vrais et suffisamment grandes pour que l'on puisse considérer une porosité moyenne dans un élément de volume.

### **2.2.3 La loi de Darcy et l'équation de diffusivité en milieu poreux**

Henry Darcy alors qu'il étudiait les fontaines de la ville de Dijon (1856) établit expérimentalement que le débit d'eau s'écoulant à travers un massif de sable peut se calculer

$$Q = KA \frac{\Delta h}{L}. \quad (23)$$

$A$  est la section du massif sableux

$\Delta h$  la perte de charge de l'eau entre le sommet et la base du massif sableux

$K$  est une constante dépendant du milieu poreux, baptisée coefficient de perméabilité

$L$  est l'épaisseur du massif sableux.

On appelle  $U = Q/A$  la **vitesse de filtration** d'un sol. À partir des équations de Navier-Stokes on sait que les causes du déplacement du fluide sont dûs au gradient de pression ainsi qu'aux forces extérieures. La loi de Darcy peut alors s'exprimer sous la forme générale

$$\vec{U} = \frac{k}{\mu} (\vec{\nabla} p + \rho g \vec{\nabla} z). \quad (24)$$

On peut réécrire cette équation

$$\vec{U} = -K (\vec{\nabla} h + \vec{\nabla} z),$$

avec  $K = k(\mu\rho g)^{-1} [LT^{-1}]$  et  $h = p(\rho g)^{-1} [L]$ . En posant

$$H = h + z,$$

on peut alors injecter l'équation de Darcy (24) dans l'équation de conservation de la masse (22) pour finalement obtenir

$$-\vec{\nabla} \cdot (K \vec{\nabla} H) + \frac{\partial \theta}{\partial t} + q = 0. \quad (25)$$

On définit parfois des lois de porosité liées à la grandeur  $f(H) = \theta$  la quantité d'eau est liée à la cote piézométrique et l'on pose

$$S_s(H) \frac{\partial H}{\partial t} = \frac{\partial \theta}{\partial t},$$

$S_s$  est appelé le coefficient d'emménagement. On appelle l'équation ainsi obtenu en injectant la formule sur la porosité, **l'équation de Richard généralisée**

$$S_s(H) \frac{\partial H}{\partial t} - \vec{\nabla} \cdot (K(\theta) \vec{\nabla} H) + q = 0. \quad (26)$$

Ceux sont avec ces équations que les géologues travaillent, nous pouvons remarquer que cette simplification des équations peut être considérée comme un processus d'upscaling. Notons que les équations de Darcy ont été justifiées dans les travaux de Matheron et Marle à partir de l'intégration dans un milieu réel des équations de Navier.

## 2.3 Les modèles

Il existe plusieurs modèles permettant de prédire les écoulements d'eau. Les premiers modèles sont des de résolution des équations de la mécanique des fluides en milieu poreux qui assurent une grande précision sur les résultats obtenus. Leurs temps de calcul les rendent inutilisables sur de grands domaines. On utilise alors des modèles moins précis mais aussi moins coûteux c'est le cas du modèle Orchidée et des domaines de surface continentale.

### 2.3.1 Modèle de surface continentale Orchidée

Il s'agit d'un modèle dynamique, développé à l'échelle globale simulant les processus continentaux à partir des interactions sol-végétation-atmosphère. Ses résultats sont intégrés comme condition à la limite basse du modèle général de circulation atmosphérique du Laboratoire de Météorologie Dynamique (modèle LMDZ). Il modélise les interactions entre surface continentale et atmosphère, en particulier les flux d'énergie, d'eau et de carbone.

## 3 Prédictions climatiques et modélisation hydrologiques

L'étude que nous avons menée a été réalisée sur le bassin du Little Washita. Nous avons eu accès aux données NARR (North American Real Reanalysis) de 1979 à 2014 qui se présentent sous la forme d'un maillage de largeur de maille de  $32km$  ainsi qu'aux données de l'IPSL qui se présentent elles sous la forme d'un maillage de largeur de maille  $200km$ .

Les données NARR (North American Regional Reanalysis) couvrent l'entièreté du continent nord américain. La méthode de projection pour passer de  $\mathcal{S}(\mathbb{R}^3)$  à  $\mathbb{R}^2$  est ce qu'on appelle la **projection Lambert**. Cette projection d'après le théorème de Gauss ne pas être une isométrie, cependant nous verrons que l'on pourra faire l'approximation que la géométrie obtenue est encore euclidienne pour le bassin du Washita et ses alentours (voir 4.2). Dans notre étude nous ferons comme si nous étions dans  $\mathbb{R}^2$  en considérant le maillage carré, la longueur de grille du maillage est d'une trentaines de kilomètres (voir figure 1).

Afin de déterminer l'efficacité du downscaling dans le cas où l'on ajoute du bruit à la variable que l'on cherche à prédire nous avons dégradé spatialement nos données en moyennant sur différentes échelles de grille (voir 3.2.1). Nous les avons ensuite downscalées et puis les avons injectées dans un modèle hydrologique (voir à citer) et nous avons finalement comparé les résultats à partir des méthodes expliquées dans la section 1.4.

<b>GRID DESCRIPTIONS</b>	
Regional North American Grid (Lambert Conformal) used by NAM, SREF and RAP.	
Nx	349
Ny	277
La1	1.000N
Lo1	214.500E = 145.500W
Res. & Comp. Flag	0 0 0 1 0 0 0
Lov	253.000E = 107.000W
Dx	32.46341 km
Dy	32.46341 km
Projection Flag (bit 1)	0
Scanning Mode (bits 1 2 3)	0 1 0
Latin1	50.000N
Latin2	50.000N
<b>Lat/Lon values of the corners of the grid</b>	
(1,1)	1.000N, 145.500W
(1,277)	46.635N, 148.639E
(349,277)	46.352N, 2.566W
(349,1)	0.897N, 68.318W
<b>Pole point</b>	
(I,J)	(174.507, 307.764)

The Dx, Dy grid increment (at 50 deg north) was selected so that the grid spacing would be exactly 32.000 km at 40 deg north; the intersection of 40N & 107W falls on point (174.507,108.664)

**Figure 3** – Description du maillage NARR

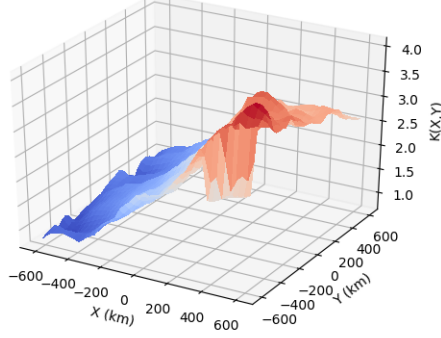
### 3.1 Analyse de la structure spatiale des données NARRs

Nous avons considéré les coordonnées du Little Washita à  $34.981^\circ N$ ,  $-97.859^\circ W$ . Nous repérons les coordonnées du point sur le maillage NARR étant le plus de ces coordonnées pour définir le point correspondant au Little Washita, en effet l'échelle du maillage NARR correspond à l'échelle du Little Washita, cette hypothèse est déjà faite dans la thèse Maquin (2016).

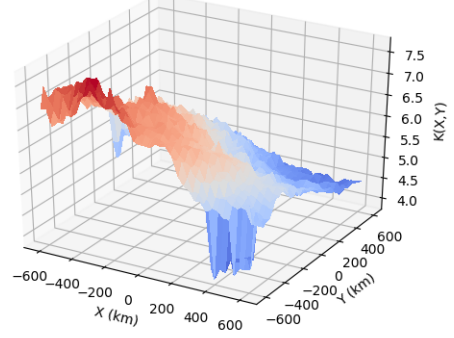
On commence par regarder les résultats autour du bassin du little washita, on fait un moyennage des précipitations et de l'évaporation sur un carré autour du little washita, le point de coordonnée (0,0) est le point correspondant au Little Washita sur ces figures.

Comme nous faisons une dégradation spatiale, il est intéressant comprendre la structure spatiale des données. Nous avons une grille spatiale pour laquelle nous avons vu que la géométrie pouvait être considérée comme euclidienne (voir 4.2). Les données NARRs que nous manipulons se présentent sous la forme d'un tenseurs  $\mathcal{T}$  de dimensions  $T \times M \times N \times v$  où  $T$  est le nombre de jours sur lesquels ont a ces observations,  $(M + 1) \times (N + 1)$  la dimension de la grille d'observation et  $v$  le nombre de variables observées. Nous faisons l'hypothèse la fonction  $S_V$  possède un noyau de covariance spatial  $K$  et nous chercherons à déterminer sa forme.

moyenne des précipitations des années 1979 aux années 2014



moyenne de l'évapotranspiration des années 1979 aux années 2014



**Définition 9.** Soit  $f : \mathbb{R}^+ \times \mathbb{R}^2 \rightarrow \mathbb{R}$  une fonction aléatoire, on définit un noyau de covariance sur cette fonction en faisant l'hypothèse que  $f$  est stationnaire en temps et espace, on a alors

$$\text{Cov}(f(t, x), f(t, x + y)) = K(y), \quad \forall t \in \mathbb{R}^+, x, y \in \mathbb{R}^2.$$

On remarque de plus que la fonction  $K$  ainsi définie est symétrique,  $K(y) = K(-y)$ ,  $\forall y \in \mathbb{R}^2$ . On peut alors simplement calculer la covariance empirique à partir de l'estimateur sans biais de la variance.

$$K(m, n) = \frac{1}{T(M - m + 1)(N - n + 1) - 1} \sum_{t=1}^T \sum_{i=m}^M \sum_{j=n}^N (\mathcal{T}_{t,i,j} - \bar{\mathcal{T}}_{1,m,n})(\mathcal{T}_{t,i-m,j-n} - \bar{\mathcal{T}}_{2,m,n}),$$

où,

$$\bar{\mathcal{T}}_{1,m,n} = \frac{1}{T(M - m + 1)(N - n + 1)} \sum_{t=1}^T \sum_{i=m}^M \sum_{j=n}^N (\mathcal{T}_{t,i,j}), \quad \bar{\mathcal{T}}_{2,m,n} = \frac{1}{T(M - m + 1)(N - n + 1)} \sum_{t=1}^T \sum_{i=0}^{M-m} \sum_{j=0}^{N-n} (\mathcal{T}_{t,i,j}).$$

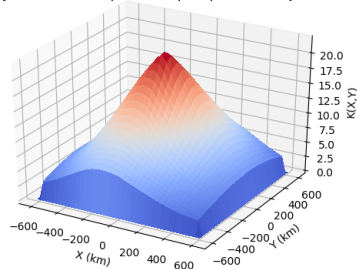
On peut donc tracer le noyau de covariance  $K(i, j)$ ,  $i \in \{-N/2, N/2\}$ ,  $j \in \{-M/2, M/2\}$ .

#### Résultats :

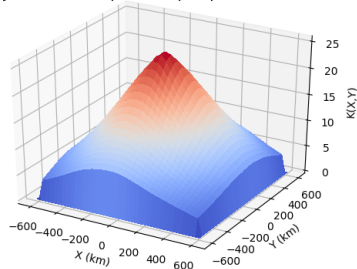
Pour obtenir nos résultats, nous avons pris en compte la saisonnalité. En effet, les résultats ont été obtenus sur chaque mois. Voici les résultats obtenus en fonction de la saisonnalité.



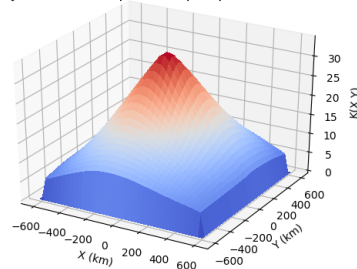
Noyau covariance spatial des précipitations en janvier



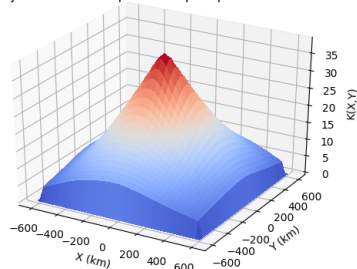
Noyau covariance spatial des précipitations en février



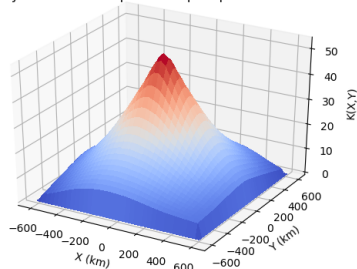
Noyau covariance spatial des précipitations en mars



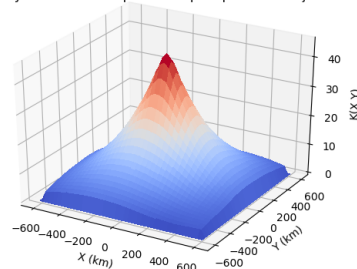
Noyau covariance spatial des précipitations en avril



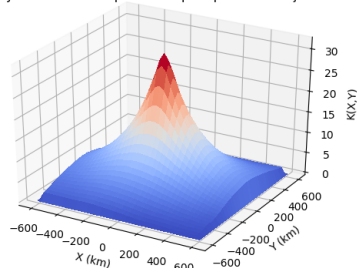
Noyau covariance spatial des précipitations en mai



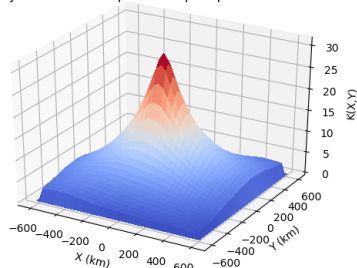
Noyau covariance spatial des précipitations en juin



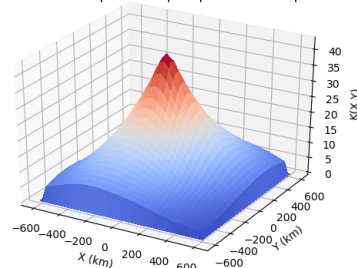
Noyau covariance spatial des précipitations en juillet



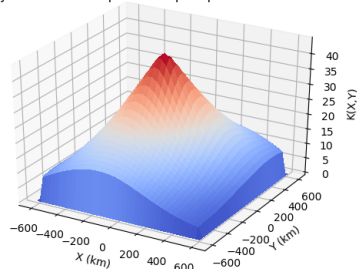
Noyau covariance spatial des précipitations en août



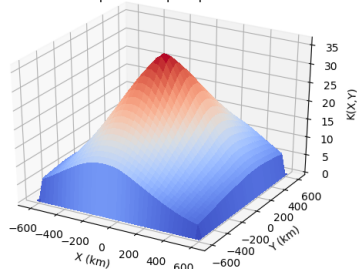
Noyau covariance spatial des précipitations en septembre



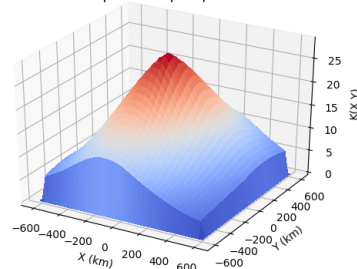
Noyau covariance spatial des précipitations en octobre



Noyau covariance spatial des précipitations en novembre

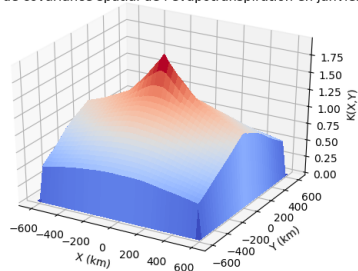


Noyau covariance spatial des précipitations en décembre

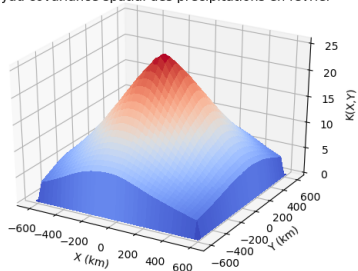


Rappelons que l'axe  $X$  correspond aux latitudes et l'axe  $Y$  aux longitudes. Nous voyons qu'en été la covariance entre deux points s'effondre relativement vite alors qu'en hiver c'est plutôt l'inverse. Il semble intéressant de remarquer les échelles d'effondrement, en effet même si l'on est à plusieurs centaines de kilomètres d'un point, on observe encore une covariance positive, on pourrait s'attendre à ce que ça s'effondre beaucoup plus vite.

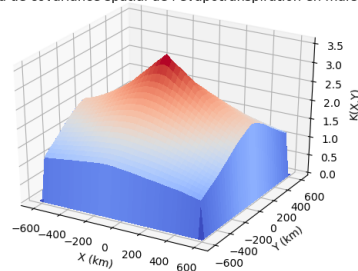
Noyau de covariance spatial de l'évapotranspiration en janvier



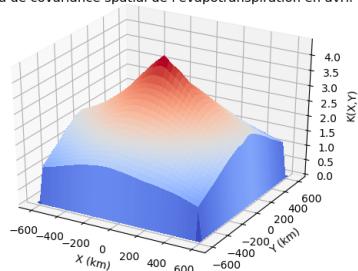
Noyau covariance spatial des précipitations en février



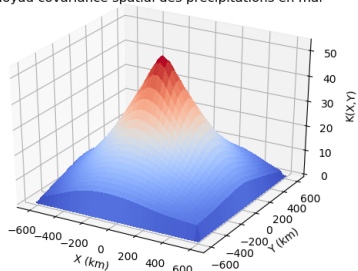
Noyau de covariance spatial de l'évapotranspiration en mars



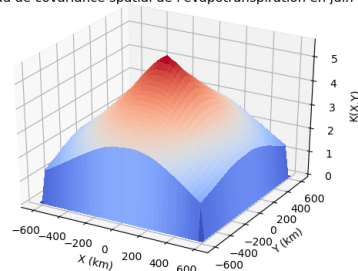
Noyau de covariance spatial de l'évapotranspiration en avril



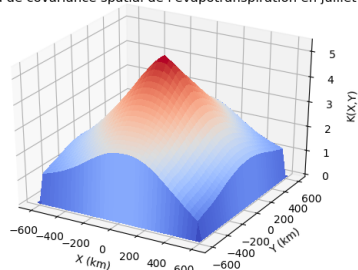
Noyau covariance spatial des précipitations en mai



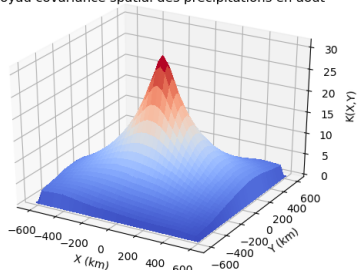
Noyau de covariance spatial de l'évapotranspiration en juin



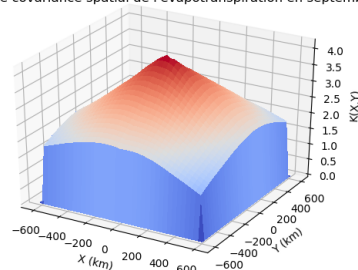
Noyau de covariance spatial de l'évapotranspiration en juillet



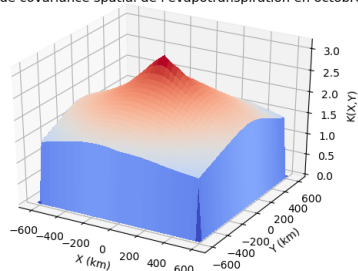
Noyau covariance spatial des précipitations en août



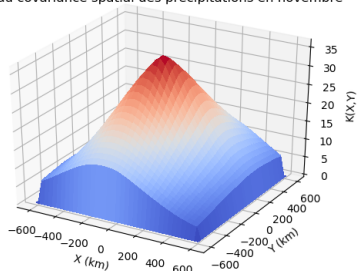
Noyau de covariance spatial de l'évapotranspiration en septembre



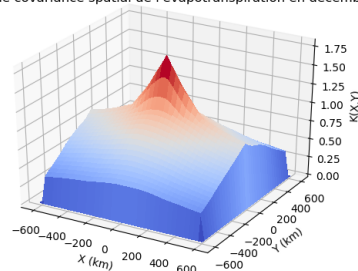
Noyau de covariance spatial de l'évapotranspiration en octobre



Noyau covariance spatial des précipitations en novembre



Noyau de covariance spatial de l'évapotranspiration en décembre



Il semble beaucoup plus difficile d'interpréter le noyau de covariance de l'évapotranspiration à partir de la saisonnalité. On voit alors que mai et août ont un noyau de covariance qui s'effondre relativement vite spatialement, que ce n'est pas le cas pour les mois de juin et juillet.

## 3.2 Présentation des résultats de prédictions climatiques

Nous allons travailler avec plusieurs séries différentes, comme nous l'avons dit précédemment nous avons travailler avec plusieurs séries temporelles.

On commence par réintroduire le tenseur  $\mathcal{T}$  de dimension  $T \times M \times N$ . Nous pouvons exprimer données NARRs par un tenseur.

$$\mathcal{T}_{t,m,n} = \mathcal{T}_V(t, lat_m, lon_n).$$

où  $t$  correspond à un temps et  $(m, n)$  une position.

### 3.2.1 La dégradation

Nous avons dégradé les valeurs par un moyennage spatial. Nous avons fait 4 dégradations différentes nous allons les numéroter ainsi  $\{1, 2, 3, 4\}$  et on appellera  $\mathcal{T}^d$  le tenseur issu de la dégradation  $d$ , défini par

$$\mathcal{T}^d_{t,m,n} = \frac{1}{(2d+1)^2} \sum_{i=-d}^d \sum_{j=-d}^d \mathcal{T}_{t,m+i,n+j} \quad \forall (m,n) \in \{d+1, \dots, M-d\} \times \{d+1, \dots, N-d\}. \quad (27)$$

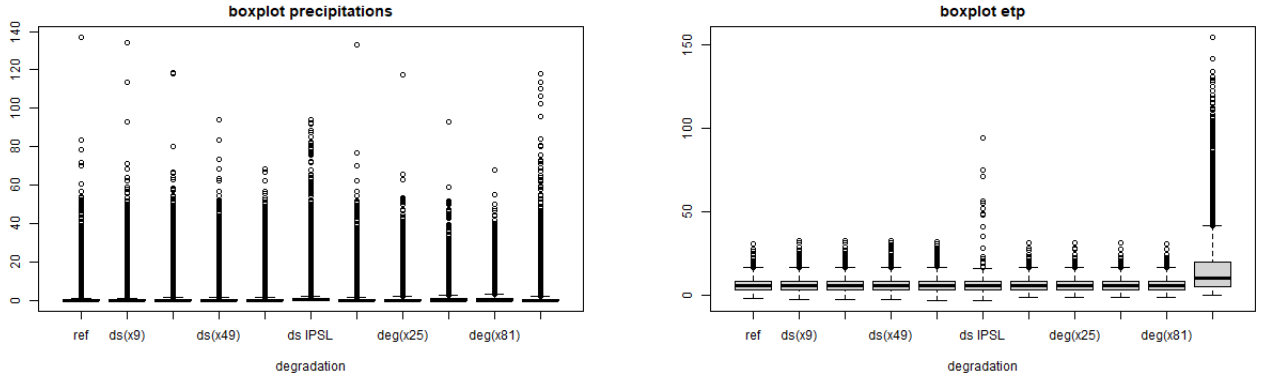
Nous avons alors étudié le point correspondant au bassin du Little Washita de coordonnées  $(a, b)$  et l'on appellera,  $(X_t^d)_{t \in \llbracket 1, T \rrbracket}$  les séries temporelles  $\mathcal{T}^d_{t,a,b}$ . Nous avons cherché à prédire la série  $(Y_t)_{t \in \llbracket 1, T \rrbracket} = \mathcal{T}^0_{t,a,b}$  à partir des séries  $(X_t^d)_{t \in \llbracket 1, T \rrbracket}$  et ainsi regarder l'impacte de la dégradation puis du downscaling sur les données.

## 3.3 Les résultats de la dégradation et du downscaling sur les données

Nous allons utiliser les outils définis dans la section 1.4 pour analyser la qualité de nos prédictions.

### 3.3.1 tracé des boxplot

Nous avons tracé les boxplot des précipitations et de l'évapotranspiration afin de voir visuellement répartitions des points.

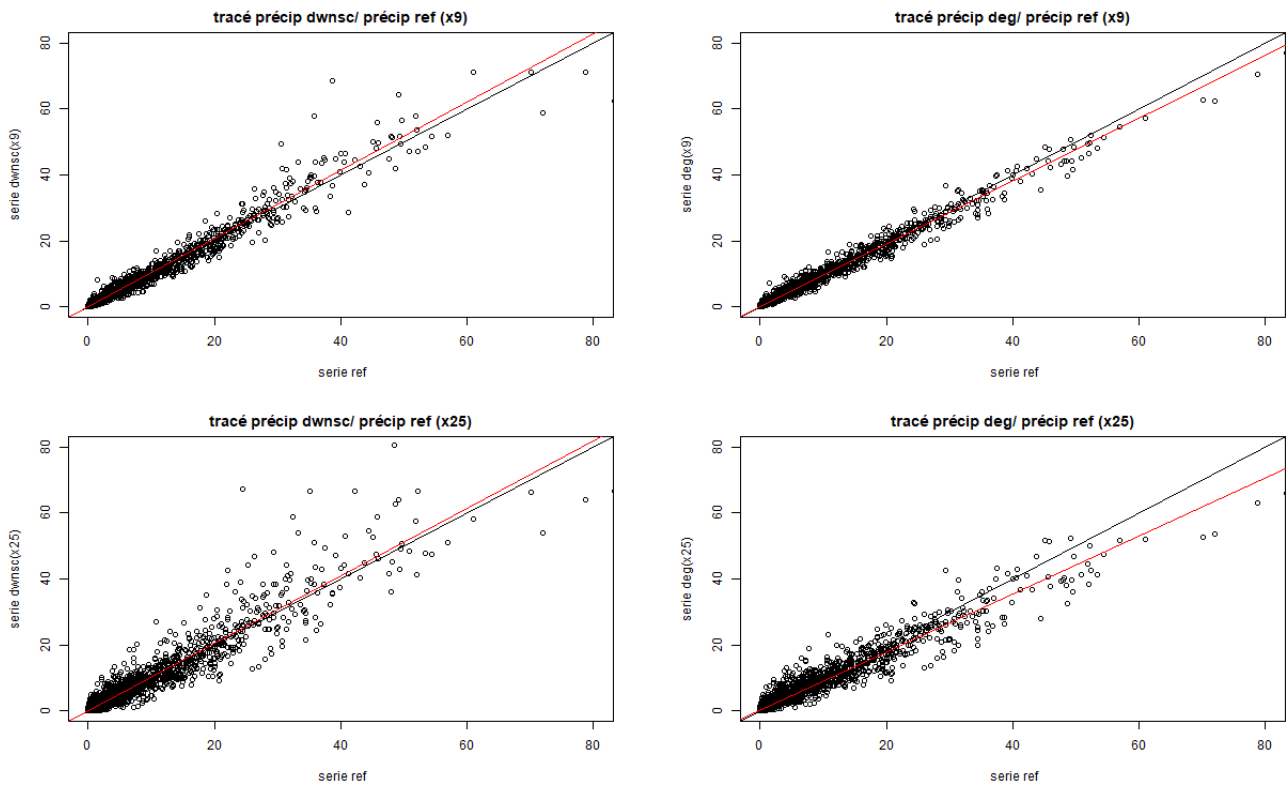


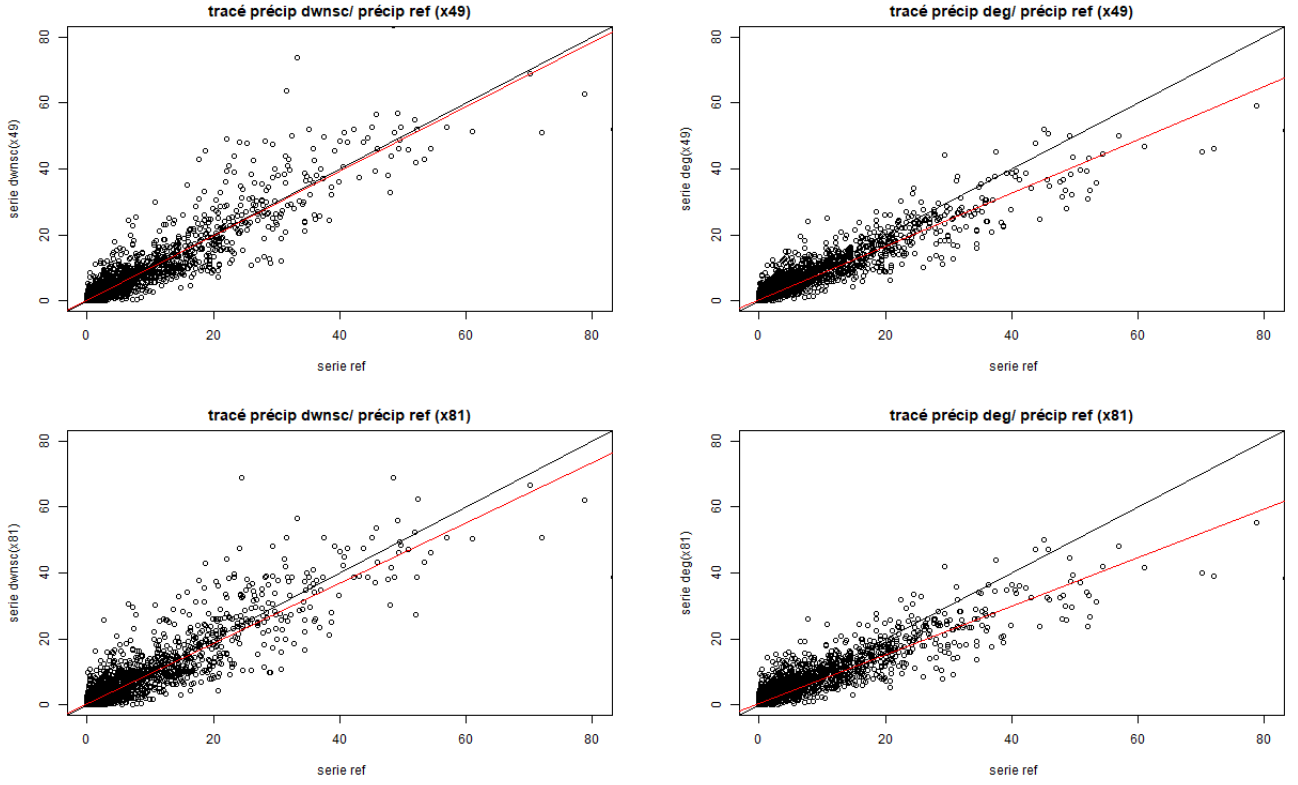
On voit que la dégradation de l'évapotranspiration n'a quasiment aucune influence sur les séries temporelles, ce qui est à prévoir d'après le noyau de covariance que nous avons obtenu précédemment. On peut aussi voir que le modèle de climat de l'IPSL prédit relativement mal l'évapotranspiration sur la région du Little Washita. On peut aussi voir la méthode de prédiction utilisée pour downscaler les séries temporelles dégradées est la méthode **CDF-t avec support égale**.

### 3.3.2 affichage des données prédites en fonction des données réelles

On commence par afficher les résultats sur un graphique en mettant les données de précipitation prédites en fonction des données de précipitation réelles, cette méthode permet de juger visuellement de l'efficacité

du downscaling sur les données. Les figures qui affichées indiquent le degré de déformation dans leur titre sous la forme  $(2d - 1)^2$  où  $d$  est le degrés que celui de l'équation (27). L'on peut voir sur chaque ligne à droite les données dégradées et à gauches les mêmes série après avoir appliquer un downscaling sur ces données.

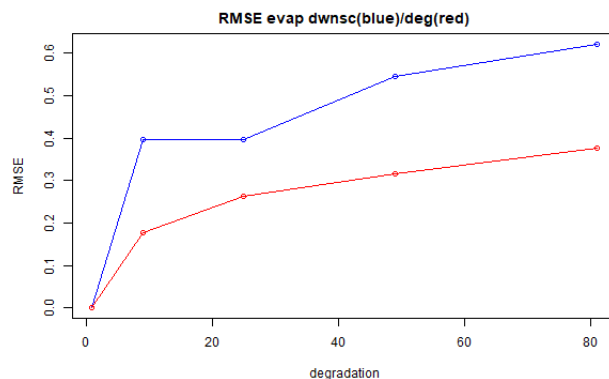
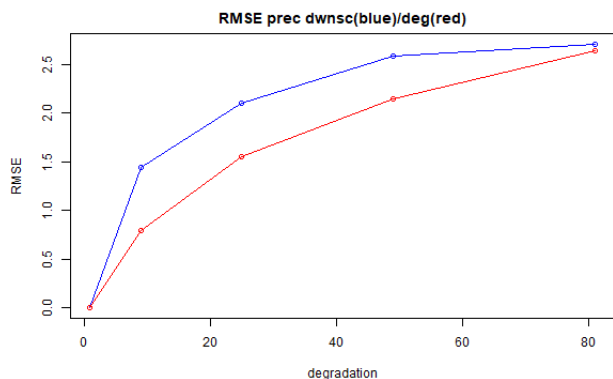




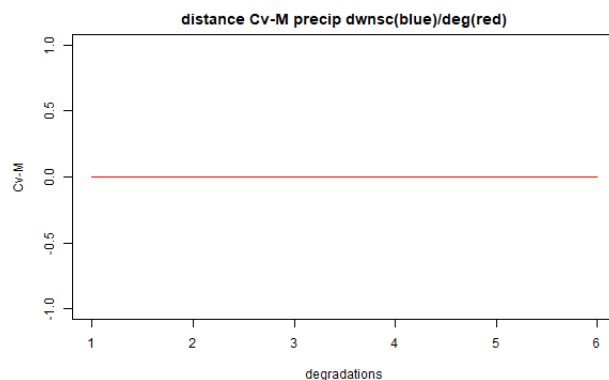
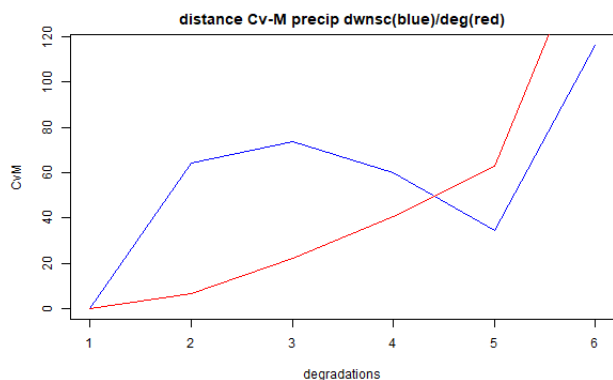
Visuellement il semble que le downscaling améliore les prédictions. Nous avons tracé pour chaque figure en noir la droite identité et en rouge les droites de régression linéaire de la forme  $y = ax + b$ , où  $y$  est la précipitation de référence et  $x$  le prédicteur de  $y$  dans notre cas c'est soit la dégradation  $\mathcal{T}_{d,a,b}$  soit son downscalé à partir de la méthode définie dans la section 1.2.2. Remarquons que la dégradation de nos données ne crée pas de changement important de l'éparpillement des points dans les graphiques, c'est en accord avec les figures obtenues dans la section ?? montrant une grande corrélation spatiale. Il semble cependant que la dégradation ait une tendance à sous-évaluer les précipitations élevées, ce qui ne semble pas aberrant si l'on considère les tailles des cumunolimbus (allant de 2 à 10km ), inférieur à la largeur de grille des données NARRs (32km).

### 3.3.3 Différence analyse Root Mean Square error et distance de Cramér-von Mises

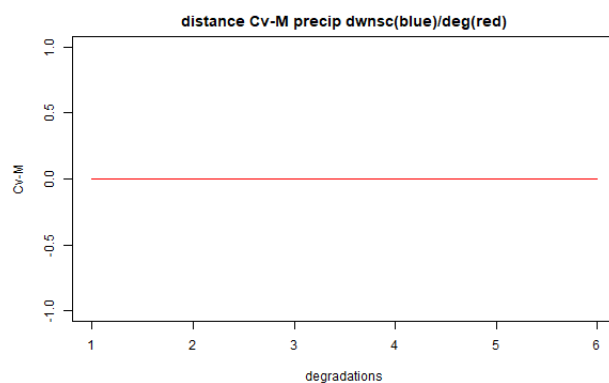
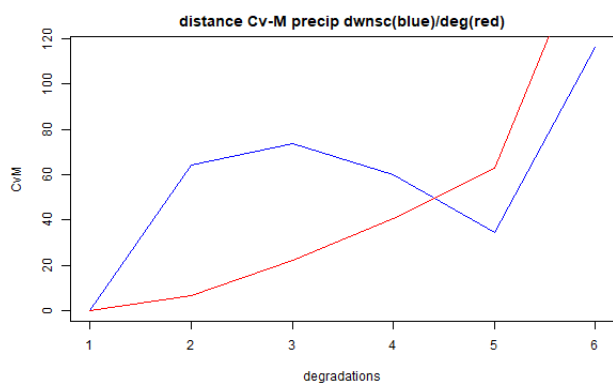
Nous avons calculer la distance quadratique moyenne entre la série de référence et les différentes séries de dégradation. Les résultats ainsi obtenues sont :



Nous voyons que les résultats sont mauvais, mais encore une fois, ce n'est pas cette distance que nous cherchons à minimiser avec le downscaling, mais une distance sur les fonctions de répartition. Regardons maintenant la distance de Cramér-Von Mises.



L'on voit qu'encore une fois la prédiction après downscaling est plus mauvaise que seulement la série dégradée, jusqu'à la dégradation  $d = 5$  et  $d = 6$  sont les données *IPSL*. L'algorithme effectue alors une grosse correction pour ce cas-ci. Cette erreur de prédiction soulève un point que nous n'avons pas encore abordé jusqu'à présent, c'est le fait que notre échantillon de précipitation possède de nombreuses valeurs à 0, ce qui signifie que la densité n'est plus une fonction mais une distribution avec un poids non nul en 0. C'est ce détail qui crée cette différence, si l'on enlève les précipitations nulle, nous obtenons :



## 4 Indexes

### 4.1 Indexe 1 : La statistique de Cramér-von Mises

Soient  $(X_i)_{i \in \llbracket 1, n \rrbracket}$  et  $(Y_i)_{i \in \llbracket 1, m \rrbracket}$  des réalisations indépendantes issues de variables aléatoires réelles  $X$  et  $Y$ . On appelle  $\mathcal{F}_n$  et  $\mathcal{G}_m$  les fonctions de répartition empiriques définies à partir de ces réalisations et  $\mathcal{H}_{m,n}$  la fonction de répartition empirique définie à partir de l'ensemble de ces réalisations  $(Z_i)_{i \in \llbracket 1, m+n \rrbracket} = X_1, \dots, X_n, Y_1, \dots, Y_m$ . Par la suite on considérera que tous les éléments sont triés dans leur ensemble ( $i \leq j \Rightarrow E_i \leq E_j$ ). Nous avons alors l'égalité suivante

$$C_{n,m} = \frac{nm}{n+m} \int_{\mathbb{R}} [\mathcal{F}_n(x) - \mathcal{G}_m(x)]^2 d\mathcal{H}_{m,n}(x) = \frac{1}{nm(m+n)} \left[ n \sum_{i=1}^n (R_{X_i} - i)^2 + m \sum_{i=1}^m (R_{Y_i} - i)^2 \right] - \frac{4nm-1}{6(m+n)}. \quad (28)$$

où  $R_{X_i}$  est le rang de  $X_i$  dans  $X_1, \dots, X_n, Y_1, \dots, Y_m$  autrement dit

$$R_{X_i} = \text{Card}(\{j \in \llbracket 1, m+n \rrbracket, Z_j \leq X_i\}).$$

Notons que cette égalité transforme un problème d'analyse en un problème de dénombrement beaucoup plus simple. On rappelle la définition de l'intégrale par rapport à une fonction.

**Définition 10.** Soient  $f$  une fonction continue par morceaux de  $\mathbb{R}$  dans  $\mathbb{R}$ , soit  $g$  une fonction continue par morceaux on définit l'intégrale en appelant  $x_{i,n} = i/n$

$$\int_{\mathbb{R}} f(x) dg(x) = \lim_{n \rightarrow \infty} \sum_{i \in \mathbb{Z}} f(x_{i,n}) (g(x_{i,n}) - g(x_{i,n-1})).$$

*Démonstration.* Commençons par montrer l'égalité

$$\frac{nm}{n+m} \int_{\mathbb{R}} [\mathcal{F}_n(x) - \mathcal{G}_m(x)]^2 d\mathcal{H}_{m,n}(x) = \frac{mn}{(m+n)^2} \sum_{i=1}^{m+n} (\mathcal{F}_n(Z_i) - \mathcal{G}_m(Z_i))^2.$$

On pose  $\delta = \inf\{|Z_i - Z_j|, Z_i \neq Z_j\}$ , quel que soit  $n \geq n_0$  tel que  $1/n_0 < \delta$  on a :

$$\sum_{i \in \mathbb{Z}} (\mathcal{F}_n(\frac{i}{n}) - \mathcal{G}_m(\frac{i}{n}))^2 \left( \mathcal{H}_{m,n}(\frac{i}{n}) - \mathcal{H}_{m,n}(\frac{i-1}{n}) \right) = \sum_{i=1}^{m+n} (\mathcal{F}_n(Z_i) - \mathcal{G}_m(Z_i))^2,$$

on obtient donc directement l'égalité voulue en passant à la limite.

Observons maintenant que  $\mathcal{F}_n(X_i) = i/n$  et  $\mathcal{G}_m(X_i) = (R_{X_i} - i)/m$  ainsi que  $\mathcal{F}_n(Y_i) = (R_{Y_i} - i)/n$  et  $\mathcal{G}_m(Y_i) = i/m$ . On peut alors réécrire  $C_{n,m}$  en séparant la somme sur les  $X_i$  et  $Y_i$

$$\begin{aligned} C_{n,m} &= \frac{mn}{(m+n)^2} \left[ \sum_{i=1}^n \left( \frac{i}{n} - \frac{R_{X_i} - i}{m} \right)^2 + \sum_{i=1}^m \left( \frac{R_{Y_i} - i}{n} - \frac{i}{m} \right)^2 \right] \\ &= \frac{mn}{(m+n)^2} \left[ \frac{1}{m^2} \sum_{i=1}^n \left( R_{X_i} - i \frac{m+n}{n} \right)^2 + \frac{1}{n^2} \sum_{i=1}^m \left( R_{Y_i} - i \frac{m+n}{m} \right)^2 \right] \end{aligned}$$

Remarquons que  $C_{n,m}$  est de la forme

$$C_{n,m} = \frac{mn}{(m+n)^2} \left[ \frac{C_1}{m^2} + \frac{C_2}{n^2} \right],$$

et que  $C_1$  et  $C_2$  sont symétriques en  $n$  et  $m$ . On définit  $\Sigma_1 = \sum_{i=1}^n R_{X_i}^2$ ,  $\Sigma_2 = \sum_{i=1}^m R_{Y_i}^2$  et  $\mathcal{S}_k = \sum_{i=1}^k i^2$  nous allons travailler sur l'expression

$$C_1 = \sum_{i=1}^n \left( R_{X_i} - i \frac{m+n}{n} \right)^2.$$

On la développe puis factorise pour obtenir

$$C_1 = \frac{m+n}{n} \sum_{i=1}^n (R_{X_i} - i)^2 - \frac{m}{n} \Sigma_1 + \frac{m(m+n)}{n^2} \mathcal{S}_n.$$

On obtient de la même manière

$$C_2 = \frac{m+n}{m} \sum_{i=1}^m (R_{Y_i} - i)^2 - \frac{n}{m} \Sigma_2 + \frac{n(m+n)}{m^2} \mathcal{S}_m.$$

D'après ce qu'on a dit précédemment on a donc :

$$C_{n,m} = \frac{1}{nm(m+n)} \left[ n \sum_{i=1}^n (R_{X_i} - i)^2 + m \sum_{i=1}^m (R_{Y_i} - i)^2 \right] - \frac{\Sigma_1 + \Sigma_2}{(m+n)^2} + \frac{\mathcal{S}_n}{n(n+m)} + \frac{\mathcal{S}_m}{m(m+n)}.$$

On remarque  $\Sigma_1 + \Sigma_2 = \mathcal{S}_{m+n}$  et que l'on a la première moitié de notre somme. Il ne reste plus qu'à développer l'expression

$$\begin{aligned} & -\frac{\mathcal{S}_{m+n}}{(m+n)^2} + \frac{\mathcal{S}_n}{n(n+m)} + \frac{\mathcal{S}_m}{m(m+n)} \\ &= -\frac{(m+n+1)(2m+2n+1)}{6(m+n)} + \frac{(n+1)(2n+1)}{6(m+n)} + \frac{(m+1)(2m+1)}{6(m+n)} = -\frac{4mn-1}{6(m+n)} \end{aligned}$$

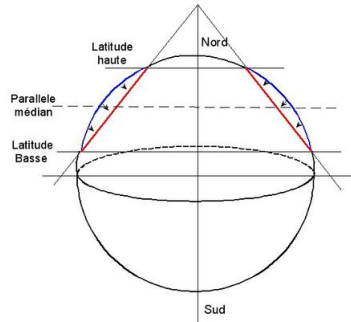
En regroupant nos deux résultats nous avons finalement :

$$C_{n,m} = \frac{1}{nm(m+n)} \left[ n \sum_{i=1}^n (R_{X_i} - i)^2 + m \sum_{i=1}^m (R_{Y_i} - i)^2 \right] - \frac{4nm-1}{6(m+n)}$$

□

## 4.2 Indexe 2 : Projection conique conforme de Lambert

La plupart des fonctions de projection de  $S(\mathbb{R}^3)$  dans  $\mathbb{R}^2$  sont des surfaces développables sur lesquelles on projette les points de la terre. Par exemple des cônes, des cylindres et des plans (projection stéréographique) sont les surfaces développables les plus connues. La **projection Lambert** est une projection conique aussi appelée la projection orthomorphique. Ses caractéristiques sont décrites dans le livre Grafarend and Krumm (2014). Elle possède la caractéristique de préserver les angles et les distances pour deux latitudes choisies, pour les données NARR les latitudes choisies sont  $33^\circ\text{N}$  et  $45^\circ\text{N}$ . De plus les lignes de latitudes égales sont des cercles et celles de longitudes égales des lignes droites. Les coordonnées que nous étudions sont entre  $33^\circ\text{N}$  et  $36^\circ\text{N}$ . On va montrer que les longueurs étudiées dans cette zone de l'espace ne souffrent que de très peu de déformation.



**Figure 4** – Projection conique conforme de Lambert



**Définition 11.** On définit les fonctions  $lat : S(\mathbb{R}^3) \rightarrow [-\pi/2, \pi/2]$  et  $lon : S(\mathbb{R}^3) \rightarrow [-\pi, \pi]$  qui associent à chaque point  $x$  de  $S(\mathbb{R}^3)$  sa latitude et sa longitude en radian.

**Définition 12.** On définit le cône convexe  $\zeta_{\theta, \theta+\epsilon}$  comme l'ensemble des droites passant les points  $x_1$  et  $x_2$  de même longitudes de latitudes égale à  $\theta$  et  $\theta + \epsilon$ . Autrement dit

$$\zeta_{\theta, \theta+\epsilon} = \{D(x_1, x_2), x_1, x_2 \in S(\mathbb{R}^3), lon(x_1) = lon(x_2), lat(x_1) = \theta, lat(x_2) = \theta + \epsilon\}.$$

où  $D(x_1, x_2)$  est la droite passant par  $x_1$  et  $x_2$ .

Il est évident que pour les lignes de latitude haute  $(\theta + \epsilon)$  et basse  $(\theta)$  les longueurs sont conservées. On définit  $\epsilon = \pi(45 - 33)/180$ .

**Proposition 3.** Pour toute courbe  $\gamma : [0, 1] \rightarrow S(\mathbb{R}^3)$  continue dont les latitudes sont comprises entre  $\theta$  et  $\theta + \epsilon$ , on a

$$\min \left( \cos(\epsilon/2), \frac{\cos(\theta + \epsilon)}{\cos(\theta)} \right) \leq \frac{\|P(\gamma)\|_{\cdot} \cdot \|\cdot\|_{\mathbb{R}^2}}{\|\gamma\|_{\cdot} \cdot \|S(\mathbb{R}^3)} \leq 1.$$

Où  $P$  est la projection de Lambert conservant les longueurs pour les latitudes  $\theta$  et  $\theta + \epsilon$ .

Ce résultat permettra de conclure que la géométrie des lieux peut être considérée comme euclidienne si  $\epsilon$  est suffisamment petit.

*Démonstration.* (Esquisse) On montrera cette inégalité pour les courbes de latitudes constantes et pour celles de longitudes constantes et la densité des fonctions de longitude ou de latitude par morceaux constantes dans l'ensemble des courbes permettra de conclure cette inégalité pour toutes les courbes.

On définit trois ensembles de courbes,  $\mathcal{C}_1$ ,  $\mathcal{C}_2$  et  $\mathcal{C}_m$  tels que

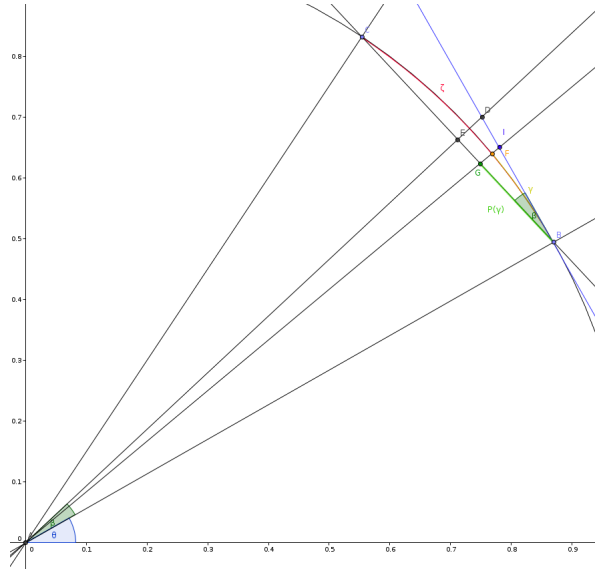
$$\mathcal{C}_1 = \{\gamma \in \mathcal{C}, lon(\gamma(t)) = c, \forall t \in [0, 1],$$

$$\mathcal{C}_2 = \{\gamma \in \mathcal{C}, lat(\gamma(t)) = c, \forall t \in [0, 1]\},$$

$$\mathcal{C}_m = \{\gamma \in \mathcal{C}, \exists t_1 < \dots < t_n, \forall i \leq n, \gamma_i : t \mapsto \gamma(t_i + t(t_{i+1} - t_i)) \in \mathcal{C}_1 \cup \mathcal{C}_2\}.$$

On commence par étudier les courbes  $\gamma$  dans  $\mathcal{C}_1$  injectives, on peut alors sans perte de généralité se placer dans le cas du cercle unité dans  $\mathbb{R}^2$  (le cercle de longitude constante). On a alors l'égalité

$$\|\gamma\|_{\cdot} \cdot \|S(\mathbb{R}^3)\| = \|\gamma\|_{\cdot} \cdot \|S(\mathbb{R}^2)\| = \int_0^1 |\gamma'(t)| dt.$$



**Figure 5** – Figure explicative de la projection lambert conforme

D'après ce schéma, on voit que pour tout point  $F$  sur le cercle la longueur de la courbe  $\gamma$  allant de  $B$  à  $F$  est majorée par la longueur du segment  $BI$ , de plus sur la figure on a  $\epsilon/2 = \beta$ . On obtient facilement l'inégalité

$$\cos(\epsilon/2) \leq \frac{\|P(\gamma)\|_{\|\cdot\|_{\mathbb{R}^2}}}{\|\gamma\|_{\|\cdot\|_{S(\mathbb{R}^3)}}},$$

Étudions maintenant les courbes  $\gamma$  dans  $\mathcal{C}_2$  de latitude égale à  $\theta + \epsilon$  et injectives. On sait que la courbe  $\gamma$  ainsi que sa projection  $P(\gamma)$  décrivent un arcs de cercle dans  $\mathbb{R}^3$ . Le rapport entre la longueur de l'arc de cercle défini par  $P(\gamma)$  et  $\gamma$  dans  $\mathbb{R}^3$  est majoré grossièrement par  $\cos(\theta + \epsilon)/\cos(\beta)$ .

$$\frac{\cos(\theta + \epsilon)}{\cos(\beta)} \leq \frac{\|P(\gamma)\|_{\|\cdot\|_{\mathbb{R}^2}}}{\|\gamma\|_{\|\cdot\|_{S(\mathbb{R}^3)}}}$$

Pour chaque courbe  $\gamma$  dans  $\mathcal{C}_m$  on a alors

$$\min\left(\frac{2 \sin(\epsilon/2)}{\epsilon}, \frac{\cos(\theta + \epsilon)}{\cos(\theta)}\right) \leq \frac{\|P(\gamma)\|_{\|\cdot\|_{\mathbb{R}^2}}}{\|\gamma\|_{\|\cdot\|_{S(\mathbb{R}^3)}}} \leq 1,$$

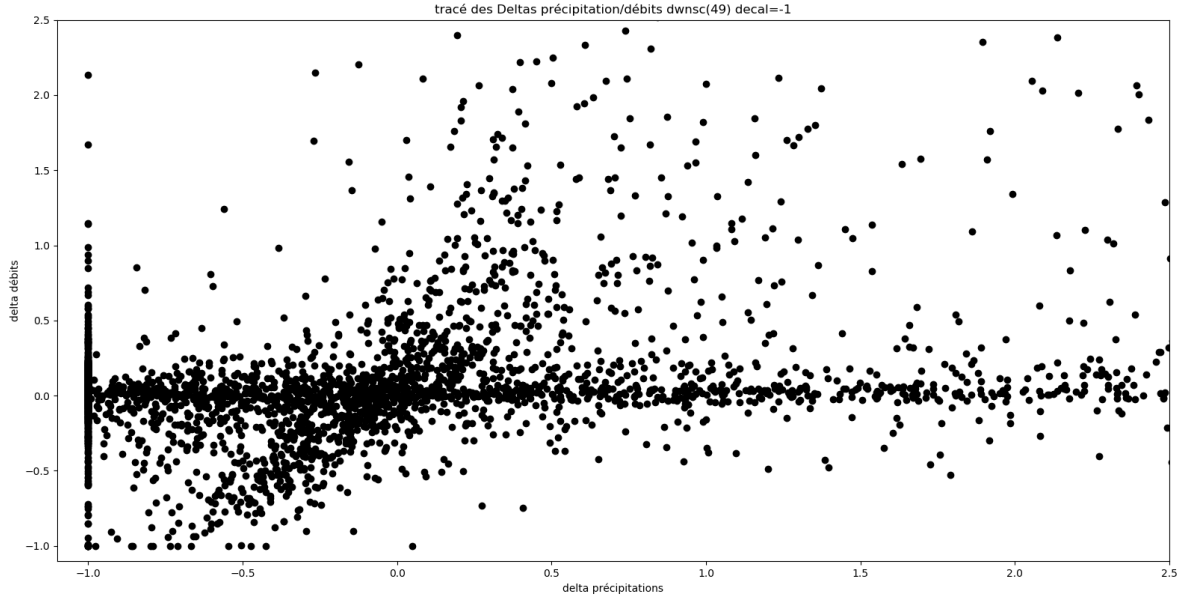
la densité de  $\mathcal{C}_m$  dans l'ensemble des courbes continues permet de conclure.  $\square$

Finalement, on peut voir que dans notre cas où  $\theta = \pi 33/180$  et  $\epsilon = \pi 12/180$ . On a que pour toute courbe  $\gamma : [0, 1] \rightarrow S(\mathbb{R}^3)$  dont la latitude est comprise entre  $\theta$  et  $\theta + \epsilon$  on a

$$0.84 \leq \frac{\|P(\gamma)\|_{\|\cdot\|_{\mathbb{R}^2}}}{\|\gamma\|_{\|\cdot\|_{S(\mathbb{R}^3)}}} \leq 1.$$

### 4.3 Indexe 3 : Classification des populations de débit

On commence par rappeler ce que sont les données  $\Delta Q$  et  $\Delta R$ , sont les différences entre les séries prédites et les séries de référence ( $R$  correspond à la pluie et  $Q$  au débit). Nous avons vu que les schémas nous incitent à considérer deux classes de points (citer la partie correspondante).



**Figure 6** – tracé des  $\Delta Q(t)$  en fonction des  $\Delta R(t-1)$

On définit alors les ensembles  $\mathcal{C}$  et  $\mathcal{I}$  tels que

$$\mathcal{C} = \{(\Delta Q(t), \Delta R(t-1)), \Delta Q(t) = f(\Delta R(t-1)) + \epsilon(t)\},$$

où  $\epsilon$  est un bruit blanc et  $f$  une fonction affine, et

$$\mathcal{I} = \{(\Delta Q(t), \Delta R(t-1)), \text{cov}(\Delta Q(t), \Delta R(t-1)) = 0\},$$

c'est à dire les lorsque les points  $\Delta Q(t)$  et  $\Delta R(t-1)$  sont indépendants.

Il paraît alors pertinent d'utiliser deux droites  $D_1$  et  $D_2$  pour classifier ces débits. Et l'on va chercher à définir les ensembles  $\mathcal{C}$  et  $\mathcal{I}$  à partir de ces deux droites. Soit  $X$  un ensemble de points  $x_1, x_2, \dots, x_n \in \mathbb{R}^2$ , on cherche deux droites  $D_1$  et  $D_2$  minimisant la valeur

$$\sum_{i=1}^n d(x_i, D_1 \cup D_2),$$

où  $d$  est la distance définie par

$$d(x, E) = \min_{e \in E} |x - e|^2.$$

On peut définir les droites de  $\mathbb{R}^2$  par un couple de points  $(u, v)$  dans  $\mathcal{U}_1(\mathbb{R}^2) \times \mathbb{R}^2$ , où  $\mathcal{U}_1(\mathbb{R}^2)$  est le cercle unité et  $u$  définit la direction de la droite et  $v$  son orientation. Alors, on peut redéfinir le problème dans ce cadre. On cherche les couples  $(u_1, v_1)$  et  $u_2, v_2$  minimisant la fonction  $F$  satisfaisant

$$F : (\mathcal{U}_1(\mathbb{R}^2) \times \mathbb{R}^2)^2 \rightarrow \mathbb{R} \\ ((u_1, v_1), (u_2, v_2)) \mapsto \sum_{i=1}^n \min(d(x_i, \mathbb{R}u_1 + v_1), d(x_i, \mathbb{R}u_2 + v_2)) \quad (29)$$

On commence par développer la distance d'un point à une droite, c'est une formule de projection classique. On appelle  $f$  la fonction  $f : (\mathcal{U}_1(\mathbb{R}^2) \times \mathbb{R}^2) \times \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f((u, v), x) = d(x, \mathbb{R}u + v)$ .

$$\begin{aligned} f_x(u, v) &= d(x, \mathbb{R}u + v) \\ &= |x - v - u\langle x - v, u \rangle|^2 \\ &= |x - v|^2 - \langle x - v, u \rangle^2 \\ &= |x|^2 - 2\langle x, v \rangle + |v|^2 - \langle x, u \rangle^2 + 2\langle x, u \rangle \langle v, u \rangle - \langle v, u \rangle^2. \end{aligned} \quad (30)$$

On cherche maintenant à calculer  $\vec{\nabla} f_x$ , le gradient étant une application linéaire, on a

$$\vec{\nabla} f_x(u, v) = \vec{\nabla} |x|^2 - \vec{\nabla} 2\langle x, v \rangle + \vec{\nabla} |v|^2 - \vec{\nabla} \langle x, u \rangle^2 + \vec{\nabla} 2\langle x, u \rangle \langle v, u \rangle - \vec{\nabla} \langle v, u \rangle^2,$$

après avoir développé puis factoriser on obtient finalement,

$$\vec{\nabla} f_x(u, v) = 2 \begin{pmatrix} (v - x)\langle u, x - v \rangle \\ v - x + u\langle u, x - v \rangle \end{pmatrix}. \quad (31)$$

On peut donc en déduire une formule pour la fonction  $F$  définie dans (29)

$$\begin{aligned} F((u_1, v_1), (u_2, v_2)) &= \sum_{i=1}^n \min(f_{x_i}(u_1, v_1), f_{x_i}(u_2, v_2)) \\ &= \sum_{i=1}^n \frac{f_{x_i}(u_1, v_1) + f_{x_i}(u_2, v_2) - |f_{x_i}(u_1, v_1) - f_{x_i}(u_2, v_2)|}{2}. \end{aligned} \quad (32)$$

On obtient finalement une expression pour  $\vec{\nabla} F$ ,

$$\vec{\nabla} F(u_1, v_1, u_2, v_2) = \left( \sum_{i=1}^n \mathcal{I}^-_{x_i}((u_1, v_1), (u_2, v_2)) \vec{\nabla} f_{x_i}(u_1, v_1) \right) \quad (33)$$

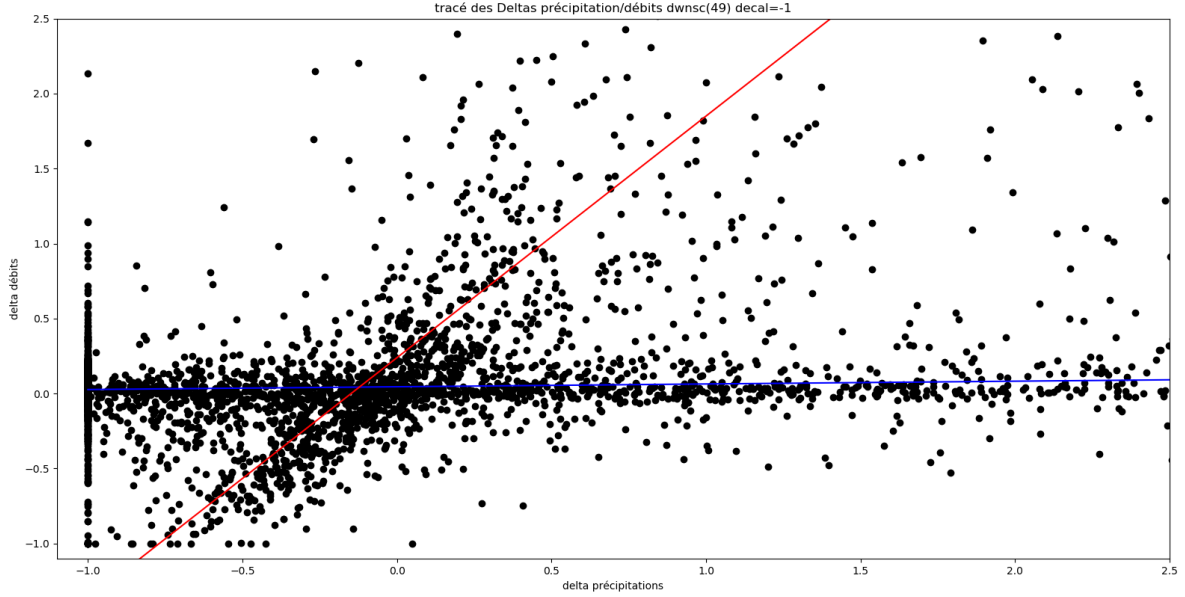
où

$$\mathcal{I}^-_{x_i}((u_1, v_1), (u_2, v_2)) = \mathbb{1}_{]-\infty, 0]}(f_{x_i}(u_1, v_1) - f_{x_i}(u_2, v_2))$$

et

$$\mathcal{I}^+_{x_i}((u_1, v_1), (u_2, v_2)) = \mathbb{1}_{[0, \infty[}(f_{x_i}(u_1, v_1) - f_{x_i}(u_2, v_2)).$$

On peut appliquer l'algorithme d'Uzawa (voir par exemple Boyd et al. (2004)) pour trouver le minimum sur  $(\mathcal{U}_1(\mathbb{R}^2) \times \mathbb{R}^2)^2$  en considérant le plongement de  $F$  dans  $(\mathbb{R}^2 \times \mathbb{R}^2)^2$  avec les contraintes  $|u_i| = 1$ . Nous avons cherché ces droites avec la fonction de minimize de "scipy.optimize" voir Jones et al. (01). Les deux droites trouvées sont alors



**Figure 7** – tracé des  $\Delta Q(t)$  en fonction des  $\Delta R(t-1)$  avec droites de classification

## Références

- Allaire, G. (2005). *Analyse numérique et optimisation : une introduction à la modélisation mathématique et à la simulation numérique*. Editions Ecole Polytechnique.
- Ayar, P. V., Vrac, M., Bastin, S., Carreau, J., Déqué, M., and Gallardo, C. (2016). Intercomparison of statistical and dynamical downscaling models under the euro-and med-cordex initiative framework : present climate evaluations. *Climate dynamics*, 46(3-4) :1301–1329.
- Boyd, S., Boyd, S. P., and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.
- Büning, H. (2002). Robustness and power of modified lepage, kolmogorov-smirnov and cramér-von mises two-sample tests. *Journal of Applied Statistics*, 29(6) :907–924.
- Cappus, P. (1960). Etude des lois de l’écoulement-application au calcul et à la prévision des débits. *La houille blanche*, pages 493–520.
- Christensen, J. H., Boberg, F., Christensen, O. B., and Lucas-Picher, P. (2008). On the need for bias correction of regional climate change projections of temperature and precipitation. *Geophysical Research Letters*, 35(20).
- De Marsily, G. (1986). Quantitative hydrogeology. Technical report, Paris School of Mines, Fontainebleau.
- Durrett, R. (2019). *Probability : theory and examples*, volume 49. Cambridge university press.
- Éthier, F. (2011). *À propos de divers tests statistiques pour l’égalité des lois*. PhD thesis, Université du Québec à Trois-Rivières.
- Fisz, M. (1963). Probability theory and mathematical statistics.
- Grafarend, E. W. and Krumm, F. W. (2014). *Map projections*. Springer.
- Hillel, D. (2003). *Introduction to environmental soil physics*. Elsevier.
- Horton, R. E. (1933). The role of infiltration in the hydrologic cycle. *Eos, Transactions American Geophysical Union*, 14(1) :446–460.
- Huang, W.-L. and Chen, S.-P. (2012). Optimal aggregate production planning with fuzzy data. *International Journal of Industrial and Manufacturing Engineering*, 6(8) :1633–1638.
- Jones, E., Oliphant, T., Peterson, P., et al. (2001–). SciPy : Open source scientific tools for Python.
- Lindgren, F., Rue, H., and Lindström, J. (2011). An explicit link between gaussian fields and gaussian markov random fields : the stochastic partial differential equation approach. *Journal of the Royal Statistical Society : Series B (Statistical Methodology)*, 73(4) :423–498.
- Maquin, M. (2016). *Développement d’un modèle hydrologique de colonne représentant l’interaction nappe-végétation-atmosphère et applications à l’échelle du bassin versant*. PhD thesis, Université Paris-Saclay (ComUE).
- Maraun, D. (2012). Nonstationarities of regional climate model biases in european seasonal mean temperature and precipitation sums. *Geophysical Research Letters*, 39(6).
- Nahar, J., Johnson, F., and Sharma, A. (2017). Assessing the extent of non-stationary biases in gcms. *Journal of Hydrology*, 549 :148–162.
- Robin, Y., Vrac, M., Naveau, P., and Yiou, P. (2019). Multivariate stochastic bias corrections with optimal transport. *Hydrology and Earth System Sciences*, 23(2) :773–786.

- Villani, C. (2003). *Topics in optimal transportation*. Number 58. American Mathematical Soc.
- Vrac, M., Drobinski, P., Merlo, A., Herrmann, M., Lavaysse, C., Li, L., and Somot, S. (2012). Dynamical and statistical downscaling of the french mediterranean climate : uncertainty assessment. *Natural Hazards and Earth System Sciences*, 12(9) :2769–2784.