

UNIVERSITY GRENOBLE ALPES



MASTER 2 SCIENCE IN INDUSTRIAL AND APPLIED  
MATHEMATICS

---

STATISTICAL METHODOLOGIES OF PRECIPITATION BIAS CORRECTION  
APPLIED TO THE OUTPUTS OF A REGIONAL CLIMATE MODEL IN THE  
ANTIZANA REGION IN ECUADOR

---

INTERNSHIP MEMORY

MARÍA BELÉN HEREDIA GUZMÁN

`Maria-Belen.Heredia-Guzman@etu.univ-grenoble-alpes.fr`

Advisor: CLEMENTINE JUNQUAS

`clementine.junquas@univ-grenoble-alpes.fr`

Co-advisor: CLEMENTINE PRIEUR

`clementine.prieur@imag.fr`

Co-advisor: THOMAS CONDOM

`thomas.condom@ird.fr`

GRENOBLE, JUIN 2017



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	In-situ data . . . . .	1
1.2	WRF simulation . . . . .	3
1.3	CHIRPS satellite product . . . . .	4
<b>2</b>	<b>The bias correction methods</b>	<b>6</b>
2.1	The Gaussian process model . . . . .	6
2.1.1	Some concepts of Kriging in Spatial Statistics . . . . .	9
2.2	Cumulative Distribution Function-transform approach . . . . .	9
2.3	Spatial CDF-t approach . . . . .	11
2.3.1	Basic time series concepts and tests . . . . .	11
2.3.2	The spatial CDF-t approach algorithm . . . . .	13
2.3.3	The final procedure . . . . .	13
2.4	Comparison criteria . . . . .	14
<b>3</b>	<b>Implementation of Gaussian process models to correct WRF precipitation biases</b>	<b>16</b>
3.1	The spatial distribution of WRF precipitation biases . . . . .	16
3.2	Selection of the deterministic function $f(x)$ . . . . .	17
3.3	Annual accumulated precipitation correction during the 2014 and 2015 periods	18
3.4	The daily precipitation correction . . . . .	19
<b>4</b>	<b>Implementation of Cumulative Distribution Function-transform spatial approach</b>	<b>22</b>
4.1	Time series analysis . . . . .	22
4.2	Spatial CDF-t approach results . . . . .	24
4.2.1	Spatial correction evaluation . . . . .	25
<b>5</b>	<b>Intercomparison between the CDF-t spatial method and Gaussian process model</b>	<b>26</b>
5.1	Evaluation of the CDF-t in a "future" period . . . . .	26
5.2	Evaluation of spatial correction methods . . . . .	28
5.3	Precipitation gridded products . . . . .	29
<b>6</b>	<b>Conclusions and Perspectives</b>	<b>32</b>
<b>A</b>	<b>The generalized inverse <math>F^{-1}</math></b>	<b>34</b>

# List of Figures

1.1	Region under study map and INAMHI stations. Color circles indicate the total precipitation from 2014-2015 period. . . . .	3
1.2	Four nested domains of the WRF simulation. Mean daily precipitation [mm day <sup>-1</sup> ] during the 2014-2015 period. . . . .	4
1.3	CHIRPS mean daily precipitation [mm day <sup>-1</sup> ] map in the 2014-2015 period. Meteorological stations are white colored. . . . .	4
3.1	WRF precipitation biases for the <b>a)</b> 2014 and <b>b)</b> 2015 periods. WRF precipitation biases in percentage for the <b>c)</b> 2014 and <b>d)</b> 2015 periods. . . . .	17
3.2	Correlation diagram of bias with different drifts (which are: longitude, latitude and altitude). Blue points correspond to the stations belonging to the Pacific region, black points to the Andes, green points to the Amazon region. The linear regression between the bias and its external drifts is depicted in red. . . . .	18
3.3	Daily standard deviation from the GP+longitude+latitude model for <b>a)</b> 2014 and <b>b)</b> 2015 periods. . . . .	19
3.4	Daily $Q_2$ evolution from the 2014-2015 period. . . . .	20
3.5	Standard deviation daily evolution for the prediction from the 2014-2015 period calculated in three stations: blue line for a Pacific station (number 22), green line for an Amazon station (number 25) and black line for an Andes (number 26). . . . .	20
3.6	Mean of daily precipitation [mm day <sup>-1</sup> ] maps during the 2014-2015 period. . . . .	21
4.1	<b>a)</b> Box-plots of Mann-Kendall and <b>b)</b> KPSS $p\_value$ results for observed and simulated time series (Obs: Observed time series $Y_t$ , Obs_diff: $Y_t - Y_{t-1}$ , Simu: Simulated time series $X_t$ and Simu_diff: $X_t - X_{t-1}$ ). . . . .	23
4.2	Time series treatment for an Andes region station (Antizana station number 26). <b>a)</b> Original ( $X_t$ observed in red and $Y_t$ simulated in blue) time series. <b>b)</b> Differentiated series CDFs. <b>c)</b> differentiated observed and <b>d)</b> differentiated simulated time series. . . . .	23
4.3	<b>a)</b> Differentiated observed time series ACF estimated and <b>b)</b> differentiated simulated time series ACF estimated . . . . .	24
4.4	Voronoi diagram of 26 meteorological stations in the study region. . . . .	24
4.5	Mean of daily precipitation [mm day <sup>-1</sup> ] maps during the 2014-2015 period for <b>a)</b> the WRF simulation and, <b>b)</b> spatial CDF-t approach. . . . .	25
5.1	Criteria related to precipitation occurrence (rainy/no-rainy events) in a calibration-evaluation framework (calibration over 01/2015 to 06/2016 and evaluation over 07/2016-12/2016). <b>a)</b> FAR criterion (ideal 0), <b>b)</b> POD criterion (1), <b>c)</b> PODF criterion (0) and <b>d)</b> HSS criterion (1). . . . .	26

5.2	Criteria related to precipitation intensity in a calibration-evaluation framework (calibration over 01/2015 to 06/2016 and evaluation over 07/2015-12/2015). <b>a)</b> KS, <b>b)</b> RMSE, <b>c)</b> Spearman correlation and <b>d)</b> $Q_{95}$ . . . . .	27
5.3	Data CDFs of <b>a)</b> an Andes station (number 26) and <b>b)</b> an Amazon station (number 25). . . . .	27
5.4	Boxplots of criteria related to precipitation occurrence (rainy/no-rainy events) for three gridded products: WRF, Spatial CDFt and GP, using a cross-validation leave-one-out framework. <b>a)</b> FAR criterion (ideal value 0), <b>b)</b> POD criterion (1), <b>c)</b> PODF criterion (0) and <b>d)</b> HSS criterion (1). . . . .	28
5.5	Boxplots of criteria related to precipitation intensity for three gridded products: WRF, spatial CDF-t and GP using a cross-validation leave-one-out framework. <b>a)</b> KS, <b>b)</b> RMSE, <b>c)</b> Spearman correlation and <b>d)</b> $Q_{95}$ . . . . .	28
5.6	Mean of daily precipitation biases maps using as reference comparison precipitation from CHIRPS and gridded products: WRF, spatial CDF-t and GP during the 2014-2015 period. (Biases= CHIRPS-gridded product). . . . .	30
5.7	Mean of daily precipitation maps [ $\text{mm day}^{-1}$ ] during 2014-2015 period. <b>a)</b> In-situ measures, <b>b)</b> WRF, <b>c)</b> spatial CDF-t, <b>d)</b> GP, and <b>e)</b> CHIRPS. . . . .	31

# List of Tables

1.1	Description of meteorological in-situ stations. Total precipitation during the 2014-2015 period at each meteorological station and total precipitation in [mm] simulated by WRF at 1 km resolution. The stations belonging to each region are differently colored: blue stations to the Pacific region, dark black stations belonging to the Amazon region and black stations to the Andes region. . . .	2
2.1	Contingency table to evaluate the accuracy of the approaches. The value 1 codes a rainy day, and 0 codes a no-rainy day. . . . .	14
3.1	Cross-validation leave-one-out results of annual accumulative precipitation for the three Gaussian Process models proposed with three drifts: longitude, latitude and longitude and latitude. The criteria are calculated for the 2014 and 2015 periods, separately. . . . .	19
3.2	Mean $Q_2$ predictivity coefficient for two daily approaches during the 2014-2015 period. . . . .	20
3.3	Criteria calculated over the accumulated precipitation quantity during 2014 and 2015 periods obtained from the results of the corrected precipitation by GP strategy separ. variog. . . . .	21

## Abbreviations

**ACF** Auto-correlation function.

**ADF** Augmented Dickey Fuller.

**AIC** Akaike information criterion.

**AR** Autoregressive process.

**CDF** Cumulative Distribution Function.

**CDF-t** Cumulative Distribution Function transform.

**CHIRPS** Climate Hazards Group InfraRed Precipitation with Station data.

**FAR** False alarm rate.

**GP** Gaussian process.

**HSS** Heidke skill score.

**INAMHI** Institut National de Météorologie et d'Hydrologie.

**KPSS** Kwiatkowski Phillips Schmidt Shin.

**KS** Kolmogorov-Smirnov Test.

**NCAR** National Center of Atmospheric Research.

**NOAA** National Oceanic and Atmospheric Administration.

**POD** Probability of detection.

**PODF** Probability of false detection.

**Q<sub>2</sub>** Predictivity coefficient.

**RCM** Regional Climate Model.

**RMSE** Root mean square error.

**SSR** Singularity Stochastic Removal.

**WRF** Weather Research and Forecasting Model.

## Abstract

The propose of this study is to correct statistically daily precipitation biases from the regional atmospheric WRF (Weather Research and Forecasting) model outputs in the Antizana region (Equatorial Andes) during the 2014-2015 period. Therefore, two methodologies of bias correction are studied: the first one is to model the bias through a Gaussian process model and the second one is to correct the bias using a spatial and time series adaptation of the Cumulative Distribution Function transform method. Four Gaussian process models are constructed by using common external drifts for this type of studies. The adaptation proposed to the Cumulative Distribution Function transform method is to correct differentiated time series and to spatialize the approach by using Voronoi polygons. The two methodologies are compared using a cross-validation leave-one-out framework in terms of precipitation occurrence and intensity criteria. The Gaussian process model shows the best results in most part of the criteria calculated.

## Résumé

L'objectif de cette étude est de corriger statistiquement les biais des précipitations issues du modèle WRF (Weather Research and Forecasting) dans la région de l'Antizana (Andes équatoriales), pendant la période 2014-2015. Deux méthodologies de correction des biais sont étudiées: la première consiste en modéliser le biais par un modèle de processus Gaussien. La deuxième méthodologie est la correction des biais en utilisant une adaptation spatiale des séries chronologiques de la méthode "Cumulative Distribution Function transform". On construit quatre modèles de processus Gaussien en utilisant des dérives externes utilisées communément pour ce type d'études. L'adaptation proposée est de corriger les séries chronologiques différenciées et de spatialiser la méthode en utilisant des polygones de Voronoi. On utilise finalement une méthode de "cross-validation leave-one-out" pour comparer les deux méthodologies en termes d'occurrence et d'intensité des précipitations. Le modèle de processus Gaussien montre les meilleurs résultats dans la plupart des critères calculés.

## Acknowledgement

I would like to express my deep gratitude to Clémentine Junquas, Clémentine Prieur and Thomas Condom for having guided me during the M2 internship that I developed in the Institut des Géosciences de l'Environnement. Your advice has been completely useful and the experience of working with you has been really rewarding. I express also my gratitude to the Institut de Recherche pour le Développement, LMI GREATICE, SNO GLACIOCLIM, INAHMI and SENESCYT for having supported this study.



## Déclaration

Je soussignée María Belén HEREDIA GUZMÁN auteur du mémoire *Statistical Methodologies of precipitation bias correction applied to the outputs of a regional climate model in the Antizana region in Ecuador*; déclare sur l'honneur que ce mémoire est le fruit d'un travail personnel et que je n'ai ni contrefait, ni falsifié, ni copié tout ou partie de l'oeuvre d'autrui afin de la faire passer pour la mienne. Toutes les sources d'information utilisées et les citations d'auteur ont été mentionnées conformément aux usages en vigueur. Je suis consciente que le fait de ne pas citer une source ou de ne pas la citer clairement et complètement est constitutif de plagiat, et que le plagiat est considéré comme une faute grave au sein de l'Université, pouvant être sévèrement sanctionnée par la loi.

Fait à GRENOBLE, le 15/06/2017  
María Belén HEREDIA GUZMÁN

# Chapter 1

## Introduction

The Antisana glacier is located in the Equatorial Andes Cordillera around 50 km of Quito city, the capital of Ecuador. The Antisana region is characterized by an important water reserve with around 60% of the drinking water used by Quito city which has a population of 2'234.000 inhabitants (Basantes-Serrano, 2015; Hall et al., 2012).

The glaciers evolution in the Tropical Andes is determined by several factors, the most important being the precipitation variability (e.g. Favier et al., 2004; Sicart et al., 2011). Therefore, it is crucial to better understand the precipitation spatio-temporal variability in this region. The meteorological in-situ stations in the Antisana region are few due to the complexity of its topography. For this reason, Regional Climate Models (RCMs) are used to simulate the local climate with high spatio-temporal resolutions. The use of RCM is essential in the Antisana glacier region considering that its surface is approximately 16.35 km<sup>2</sup>.

In this study the Weather Research Forecasting (WRF) model is used to simulate the atmospheric regional climate, including precipitation variables. Several previous studies used the WRF model in the Andes as for example Mourre (2015), Mourre et al. (2016), Ochoa et al. (2014), Ochoa et al. (2016), among others. In Mourre (2015) and Ochoa et al. (2014), WRF simulations are compared to rainfall products derived from satellite products and in-situ stations. Nevertheless, some authors (e.g. Giovannettone and Barros, 2009; Ochoa et al., 2014; Mourre et al., 2016) have shown that the WRF model simulates precipitation biases in the Andes, in terms of intensity (precipitation amounts) and occurrence (rainy/no rainy days), because of the complex topography. For this reason, it is important to develop bias correction methods of the simulated precipitation before using it in climate impact studies (Vrac and Friederichs, 2015). Commonly corrected precipitation gridded products are also needed as external forcing data for hydrological and glaciological models to understand water resources and glaciers evolution.

In this project two precipitation bias correction methods are studied, the first one consists in modeling the bias with a **Gaussian process based metamodel**. This approach is also known as kriging in geostatistics. The second approach generalizes the quantile-quantile methodology (Déqué, 2007) and is based on the **Cumulative Distribution function transform with singularity stochastic removal approach** (hereafter CDF-t) developed by Vrac et al. (2016). These two methods are explained in more detail in Chapter II. For comparing these two methods, Climate Hazards Group InfraRed Precipitation with Station Data rainfall (CHIRPS) dataset is used.

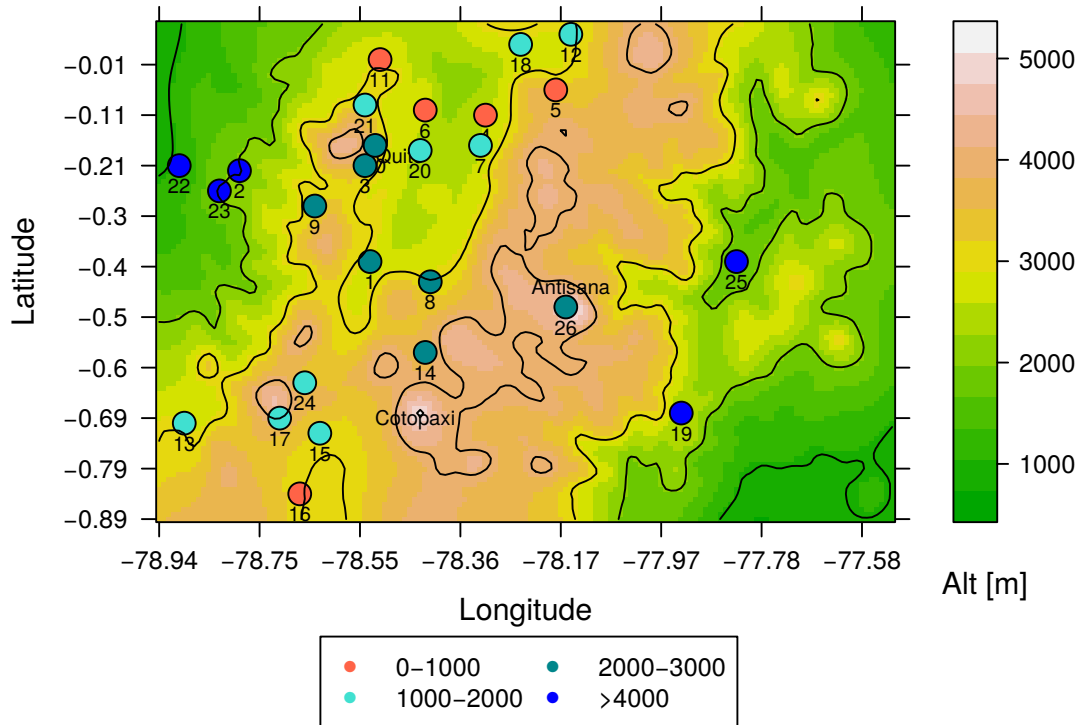
### 1.1 In-situ data

We use 26 meteorological in-situ stations, installed by the Instituto Nacional de Meteorología e Hidrología (INAMHI) in Ecuador, with a complete chronology of daily precipitation [mm

day<sup>-1</sup>] during the 2014-2015 period. Figure 1.1 shows a map with the location of the stations. The majority of them are distributed in Quito, there are two stations in the Amazon Region (stations number 19 and 25) and one station in the Antisana Region (station number 26). The study area (Figure 1.1) is divided into three regions corresponding to three regions in Ecuador: the *Pacific Region* formed by the stations 2, 22 and 23, the *Amazon Region* formed by the stations 19 and 25, and the *Andes Region* formed by the remaining ones (21 stations). There is an in-homogeneous stations distribution because the most of them are located in Andes region (80 % of the stations). Table 1.1 presents a description of the location and total precipitation during the 2014-2015 period for each meteorological station. The meteorological stations located in the Amazon Region registered the highest precipitation total values (with a total precipitation higher than 6.000 mm).

Number	Lon.	Lat.	Obs.	WRF 1km
1	-78.53	-0.39	2447	1475
2	-78.78	-0.21	8656	1568
3	-78.54	-0.20	2510	2478
4	-78.30	-0.10	769	659
5	-78.17	-0.06	835	972
6	-78.42	-0.10	831	932
7	-78.32	-0.16	1221	562
8	-78.42	-0.43	2403	1992
9	-78.63	-0.28	2886	1082
10	-78.52	-0.16	2602	2368
11	-78.51	0.00	995	623
12	-78.14	0.05	1763	1835
13	-78.89	-0.70	1694	2105
14	-78.43	-0.56	2695	753
15	-78.63	-0.72	1506	757
16	-78.66	-0.83	962	839
17	-78.70	-0.68	1203	1229
18	-78.23	0.03	1080	1327
<b>19</b>	<b>-77.93</b>	<b>-0.67</b>	<b>8276</b>	<b>2297</b>
20	-78.43	-0.18	1699	948
21	-78.54	-0.09	1824	695
22	-78.90	-0.21	8954	6892
23	-78.82	-0.25	6132	2140
24	-78.66	-0.62	1989	690
<b>25</b>	<b>-77.82</b>	<b>-0.39</b>	<b>6261</b>	<b>1186</b>
26	-78.15	-0.47	2255	2062

**Table 1.1:** Description of meteorological in-situ stations. Total precipitation during the 2014-2015 period at each meteorological station and total precipitation in [mm] simulated by WRF at 1 km resolution. The stations belonging to each region are differently colored: blue stations to the Pacific region, dark black stations belonging to the Amazon region and black stations to the Andes region.



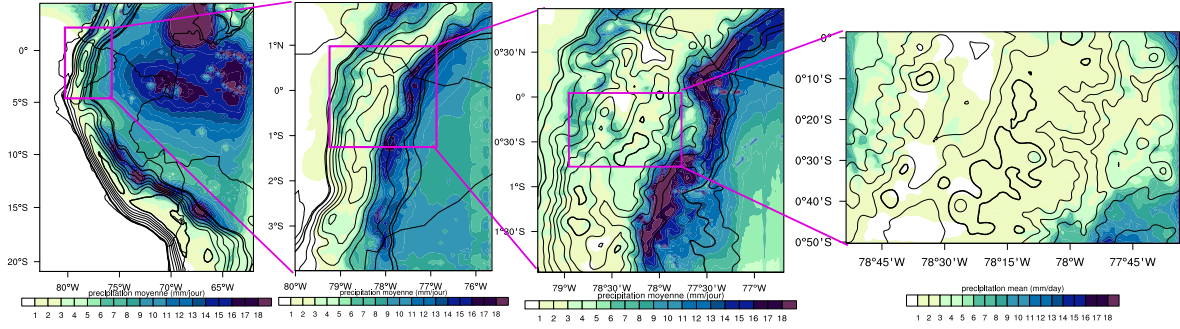
**Figure 1.1:** Region under study map and INAMHI stations. Color circles indicate the total precipitation from 2014-2015 period.

## 1.2 WRF simulation

The WRF model is a numerical weather prediction system developed since 1990 by the National Oceanic and Atmospheric Administration (NOAA) and the National Center for Atmospheric Research (NCAR). The WRF model can be used to simulate regional climate at high spatio-temporal resolution (a spatial resolution of 1 km and a time resolution of 1 hour).

The WRF model is used with nested domains, which consists into making simulations in bigger domain areas containing the region of interest and gradually reduce the resolution, in order to use the large domain simulation to force the atmospheric border conditions to simulate the smaller one. Different options of dynamical and physical parametrizations were tested in a previous study, so for this work it was chosen the WRF simulation with the parameters that have provided the better precipitation results in the Andes Region.

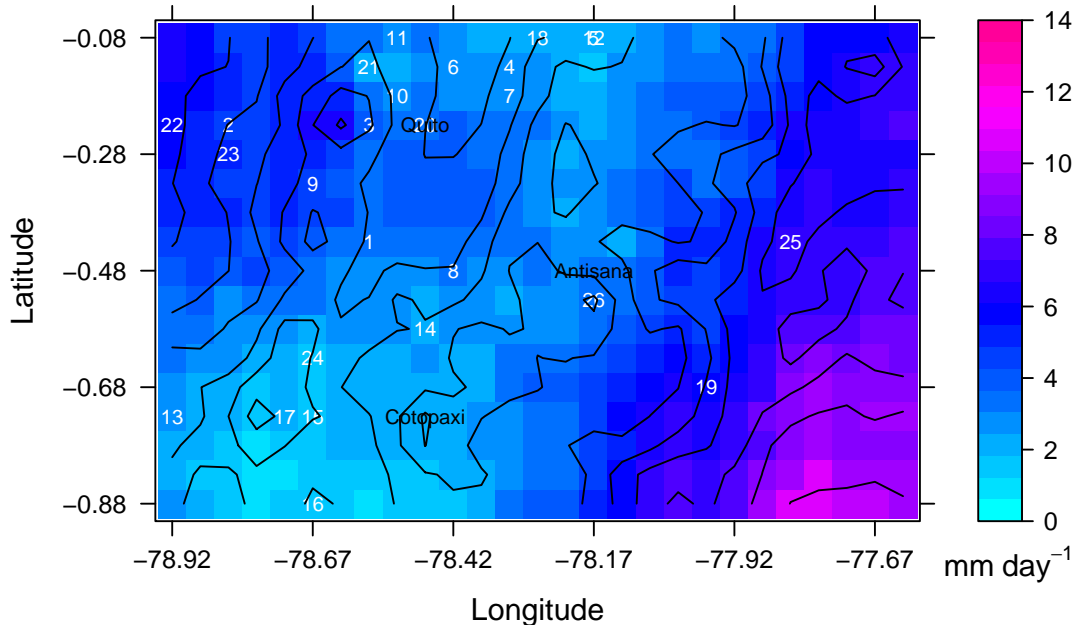
For this study, a one-way nested WRF simulation is used with four nested domains during the years 2014-2015. Figure 1.2 shows the four nested domains of the simulation. As we can see a largest domain was previously simulated to finally obtain the precipitation in the region of interest with a high spatio-temporal resolution. The largest domain has a resolution of 27 km covering an area of  $2403 \times 2943 \text{ km}^2$ . The next domain has a resolution of 9 km ( $540 \times 720 \text{ km}^2$ ). The domain of 3 km resolution covers an area of  $297 \times 333 \text{ km}^2$ , and finally, the highest resolution domain (1 km) covers an area of  $147 \times 99 \text{ km}^2$ . The simulation in the last domain (hereafter, the WRF simulation) is used as dataset for this study. The data registered in a meteorological station is associated with the closest 1 km grid of the WRF simulation. The comparison is not that precise because we are comparing one point (an in-situ station) and  $1 \text{ km}^2$  but until this moment, this is the highest spatial resolution that can be achieved with a RCM.



**Figure 1.2:** Four nested domains of the WRF simulation. Mean daily precipitation [ $\text{mm day}^{-1}$ ] during the 2014-2015 period.

### 1.3 CHIRPS satellite product

CHIRPS dataset was created by the U.S. Geological Survey and Climate Hazards Group scientists and it is a satellite image product that has been recorded since around 30 years ago. It has a spatial resolution of 5 km. One of the disadvantages of satellite products is that they are not accurate because they suffer from biases, caused by extreme precipitation which are underestimated (Climate Hazard Group). Figure 1.3 shows daily mean precipitation recorded by CHIRPS in the region under study. Due to the scarce number of station, CHIRPS product is used to evaluate the corrected gridded precipitation products (WRF, CDF-t and Gaussian Process model).



**Figure 1.3:** CHIRPS mean daily precipitation [ $\text{mm day}^{-1}$ ] map in the 2014-2015 period. Meteorological stations are white colored.

This study is organized as follows: Chapter II describes the different methodologies. Chapter III and Chapter IV present the Gaussian process model and spatial CDF-t approach results, respectively. An intercomparison between these two methods is presented in Chapter V and

finally, the conclusions and perspectives are detailed in Chapter VI.

## Chapter 2

# The bias correction methods

In this Chapter, the modeling of the bias based on Gaussian processes is described in Section 2.1, the Spatial CDF-t approach is explained in Sections 2.2 and 2.3, and finally the evaluation criteria to compare gridded precipitation products are described in Section 2.4.

### 2.1 The Gaussian process model

The Gaussian process model, known as Kriging method in spatial interpolation, takes into account the spatial statistical structure of an estimated variable; for example, in our case the variable of interest is the precipitation bias. The method of Kriging in spatial statistics was investigated by Georges Matheron, who based his investigation in the work of Daniel G. Krige in 1960. Krige used the method to estimate the distribution of gold in a region based in some samples from a few boreholes. Several studies have been developed to correct the precipitation bias based on Gaussian process models; Hanchoo Wong et al. (2012) developed a bias correction of radar rainfall based on kriging approach in Thailand, Müller and Thompson (2013) made a bias adjustment of satellite rainfall in Nepal and they used kriging to interpolate precipitation from in-situ measures and, Moure et al. (2016) made a precipitation interpolation based on kriging using as external drift the WRF simulation. The principal concepts of the Gaussian process modeling and parameters estimation as described hereafter are based on Marrel et al. (2008).

**Definition 2.1.1.** A **Gaussian Process** is a collection of random variables such that any finite number of its combination has a joint Gaussian distribution, in other words, it is a generalization of the Gaussian probability distribution (Rasmussen and Williams, 2005), and it is fully specified by its mean and covariance function.

**Definition 2.1.2.** Let  $d \in \mathbb{N}^*$ , a stochastic process  $(X_t)_{t \in \mathbb{R}^d}$  is **stationary** if

$$(X_{t_1}, \dots, X_{t_n}) \stackrel{D}{=} (X_{t_1+k}, \dots, X_{t_n+k})$$

for all  $n \in \mathbb{N}$   $t_1, \dots, t_n \in \mathbb{R}^d$ ,  $k \in \mathbb{N}$ .

Consider that  $n$  observations of a phenomenon are taken in  $\mathbb{R}^2$  (for example, the WRF bias precipitation registered in a point of the region under study). Each observation  $y(x)$  corresponds to a realization registered in a point  $x = (x_1, x_2) \in \mathbb{R}^2$ . The set of points where the observations are collected is denoted by  $x_s = (x^{(1)}, \dots, x^{(n)})$  with  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^2$  (each  $x$  corresponds to a point in the space). The set of observations of the phenomenon is denoted by  $y_s = (y^{(1)}, \dots, y^{(n)})$  with  $y^{(i)} = y(x^{(i)})$ ,  $\forall i = 1, \dots, n$ . The **Gaussian process modeling** consists into representing  $y(x)$  as a realization of a random function  $Y(x)$  that can

be decomposed into a **deterministic function**  $f(x)$  and a centered stationary **Gaussian process**  $Z(x)$ . The **Gaussian process model** is defined as:

$$Y(x) = f(x) + Z(x). \quad (2.1)$$

The function  $f(x)$  represents the tendency and commonly, it is constructed as a finite linear combination of  $k$  elementary functions ( $f_i$   $i \in \{0, \dots, k\}$ ):

$$f(x) = \sum_{j=0}^k \beta_j f_j(x) = F(x)\beta$$

where  $\beta = (\beta_0, \dots, \beta_k)$  is the regression parameter vector and  $F(x) = (f_0(x), \dots, f_k(x))$ . In  $\mathbb{R}^2$  commonly the elementary functions used are:  $f_0(x_1, x_2) = 1$ ,  $f_1(x_1, x_2) = x_1$  and  $f_2(x_1, x_2) = x_2$ . The function  $f(x)$  allows the addition of an external drift into and this is advantageous because it allows a nonstationary global modeling framework, in other words the variable  $Y$  does not need to be stationary (see Definition 2.1.2) but the variable  $Z$  is assumed to be stationary.

The Gaussian centered process  $Z(x)$  has as covariance function:

$$\text{Cov}(Z(x), Z(u)) = K(x - u) = \sigma^2 R(x - u), \quad (2.2)$$

where  $x, u \in \mathbb{R}^2$ ,  $\sigma^2$  is the variance of  $Z$  and,  $R$  is its correlation function. The process  $Z$  is stationary because it is considered that its correlation function only depends on the difference between  $x$  and  $u$ .

In this study, we use the **Matérn covariance functions** because they are stationary and commonly used in spatial statistics studies due to their flexibility (Paciorek and Schervish, 2006), and they are defined as:

$$K(x, u) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left[ \frac{\sqrt{2\nu}}{\kappa} |x - u| \right]^\nu K_\nu \left( \frac{\sqrt{2\nu}}{\kappa} |x - u| \right), \quad (2.3)$$

where  $K_\nu$  is the modified Bessel function of second kind of order  $\nu > 0$ , and  $\kappa$  is a positive parameter that represents the characteristic length scale (Rasmussen and Williams, 2005). Let's suppose  $\text{Cov}(Z(x), Z(u)) = \sigma^2 R_\theta(x - u)$  where  $\theta = (\nu, \kappa)$  is the correlation parameter vector of 2.3.

Commonly, the observations  $y^{(1)}, \dots, y^{(n)}$  are noisy. For that, an independent centered Gaussian variable  $U(x)$  with variance  $\epsilon^2 = \sigma^2 \tau$  is added to the Gaussian process model:

$$Y(x) = f(x) + Z(x) + U(x). \quad (2.4)$$

Thus, the covariance function of  $Y$  is the following:

$$\text{Cov}(Y(x), Y(u)) = \sigma^2 (R_\theta(x - u) + \tau \delta(x - u))$$

where  $\delta_v = \mathbb{1}_{\{0\}}$ .

Under the conditions of a Gaussian Process model (defined by the equation 2.4),  $Y$  follows a multinormal distribution:

$$p(Y|x_s, \beta, \sigma, \theta, \tau) = \mathcal{N}(F_s \beta, \Sigma_s), \quad (2.5)$$

where  $F_s = [F(x^{(1)})^t, \dots, F(x^{(n)})^t]^t$  and its covariance matrix  $\Sigma_s = \sigma^2 (R_\theta(x^{(i)} - x^{(j)}))_{i,j=1, \dots, n} + \tau I_n$  where  $I_n$  is a  $n$ -dimensional identity.

Let's consider a new point  $x^*$ , then the joint probability distribution of  $(Y, Y(x^*))$  is the following:



$$p(Y, Y(x^*) | x_s, x^*, \beta, \sigma, \theta, \tau) = \mathcal{N} \left[ \begin{pmatrix} F_s \\ F(x^*) \end{pmatrix} \beta, \begin{pmatrix} \Sigma_s & k(x^*) \\ k(x^*)^t & \sigma^2(1 + \tau) \end{pmatrix} \right]$$

where

$$\begin{aligned} k(x^*) &= [\text{Cov}(Y, Y(x^*))]^t \\ &= \sigma^2 [R_\theta(x^{(1)} - x^*) + \tau\delta(x^{(1)} - x^*), \dots, R_\theta(x^{(n)} - x^*) + \tau\delta(x^{(n)} - x^*)]^t. \end{aligned}$$

Then, the conditional distribution of  $Y(x^*)$  is Gaussian:

$$p(Y(x^*) | y_s, x_s, x^*, \beta, \sigma, \theta, \tau) = \mathcal{N}(E[Y(x^*) | y_s, x_s, x^*, \beta, \theta, \tau], \text{Var}[Y(x^*) | y_s, x_s, x^*, \beta, \sigma, \theta, \tau]) \quad (2.6)$$

where

$$E[Y(x^*) | y_s, x_s, x^*, \beta, \theta, \tau] = F(x^*)\beta + k(x^*)^t \Sigma_s^{-1} (y_s - F_s \beta), \quad (2.7)$$

$E[Y(x^*) | y_s, x_s, x^*, \beta, \theta, \tau]$  is the predictor of  $Y(x^*)$  and the variance is:

$$\text{Var}[Y(x^*) | y_s, x_s, x^*, \beta, \theta, \tau] = \sigma^2(1 + \tau) - k(x^*)^t \Sigma_s^{-1} k(x^*), \quad (2.8)$$

and, the variance corresponds to the mean square error of the predictor.

Finally, the estimation of parameters  $(\beta, \theta, \sigma, \tau)$  is developed in order to obtain the mean (described in 2.7) and the variance (see 2.8) of the Gaussian Process model. The parameters are estimated by using the **maximum likelihood method**. The likelihood of  $Y$  is:

$$\begin{aligned} l_Y(\beta, \theta, \sigma, \tau) &= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2} \ln(\det(R_\theta + \tau I_n)) \\ &\quad - \frac{1}{2\sigma^2} (Y - F_s \beta)^t (R_\theta + \tau I_n)^{-1} (Y - F_s \beta) \end{aligned} \quad (2.9)$$

Given the parameters  $\theta$  and  $\tau$ , the maximum likelihood estimator of  $\beta$  is:

$$\hat{\beta} = (F_s^t (R_\theta + \tau I_n)^{-1} F_s)^{-1} F_s^t (R_\theta + \tau I_n)^{-1} y_s,$$

and the maximum likelihood estimator of  $\sigma^2$  is:

$$\hat{\sigma}^2 = \frac{1}{n} (y_s - F_s \hat{\beta})^t (R_\theta + \tau I_n)^{-1} (y_s - F_s \hat{\beta}).$$

Then, taking into account the estimators of  $\beta$  and  $\sigma$ , the predictor  $\hat{Y}(x^*)$  is the following:

$$\hat{Y}(x^*) |_{y_s, x_s, x^*, \sigma, \theta, \tau} = F(x^*) \hat{\beta} + k(x^*)^t \Sigma_s^{-1} (y_s - F_s \hat{\beta})$$

and its variance is:

$$\text{Var}[\hat{Y}(x^*) | y_s, x_s, x^*, \sigma, \theta, \tau] = \sigma^2(1 + \tau) - k(x^*)^t \Sigma_s^{-1} k(x^*) + u(x^*) (F_s^t \Sigma_s^{-1} F_s)^{-1} u(x^*)^t,$$

where  $u(x^*) = F(x^*) - k(x^*)^t \Sigma_s^{-1} F_s$ .

The estimation of  $\hat{\beta}$  and  $\hat{\sigma}^2$  depends on  $\theta$  and  $\tau$  parameters. Thus, replacing the estimation of  $\beta$  and  $\sigma^2$  on 2.9, the estimation of  $\theta$  and  $\tau$  are obtained as the parameters that maximizes the following function:

$$(\hat{\theta}, \hat{\tau}) = \arg \min_{\theta, \tau} \det(R_\theta + \tau I_n)^{\frac{1}{n}} \hat{\sigma}^2. \quad (2.10)$$

Finding the parameters of 2.10 is a costly optimization problem. For this reason, there are several algorithms proposed to solve it as for example, the simplex method, Bayesian methods, among others (Marrel et al., 2008).

This study is based on Gräler et al. (2012). Gaussian process models are developed to correct daily precipitation biases and the time dimension is not directly considered in the modeling process. In other words, a Gaussian process model is constructed for each day of the 2014-2015 period to correct its precipitation bias. Separate daily variograms, as far as daily evolving ones (see hereafter) have been implemented. It is also possible to estimate

the covariance function in  $\mathbb{R}^3$  but this has not been tested during this internship (see more details in Gräler et al. (2016)). The strategies to deal with the spatial dimension are hereafter described.

### 2.1.1 Some concepts of Kriging in Spatial Statistics

The function  $K(x - u)$  defined in 2.2 is known as **variogram** in Spatial Statistics and it is usually denoted by  $\gamma(h)$  where  $h = x - u$ . The following concepts are needed to explain the strategies employed to obtain daily precipitation corrections. We decided to model a spatial variogram with possible daily evolution (see below).

- **Nugget:** It is defined as:

$$c_0 = \lim_{|h| \rightarrow 0} \gamma(h)$$

And the nugget effect ( $c_0 > 0$ ) is produced when the white noise  $U$  (defined in 2.4) introduces a discontinuity at the origin (Marrel et al., 2008).

- **Sill:** It is defined as  $\lim_{|h| \rightarrow \infty} \gamma(h)$ .
- **Range:** The distance or lag  $h$  at which  $\gamma(h)$  reaches the sill.

Then, to obtain daily precipitation two strategies are used; a separate daily variograms and a daily evolving variograms strategies based on Gräler et al. (2012). The main ideas of these two strategies are the following ones:

**Separate daily variograms:** It is based in using data of each day from bias precipitation to estimate a variogram separately to each day. One of the disadvantages of this approach is that, it does not take into account information of past days, and temporal data dependencies could be lost.

**Daily evolving variograms:** This approach is based in adding information of the previous day with a certain weight  $\lambda \in [0, 1]$  to the estimation of the new one. Hence, the estimation of a day  $D$  variogram is calculated as:

$$\begin{aligned} \text{range} &= \lambda \text{range}_D + (1 - \lambda) \text{range}_{D-1} \\ \text{nugget} &= \left( \lambda \frac{\text{nugget}_D}{\text{sill}_D} + (1 - \lambda) \frac{\text{nugget}_{D-1}}{\text{sill}_{D-1}} \text{sill}_D \right) \end{aligned}$$

$$\text{partial sill} = \text{sill}_D - \text{nugget}.$$

where  $D - 1$  is the previous day of  $D$ . Based on the study developed by Mourre et al. (2016), we use  $\lambda = 0.9$ . In our studies we have implemented both separate and evolving approaches. Additional tests for finding the best  $\lambda$  should have been done, but it is out of the scope of this work.

## 2.2 Cumulative Distribution Function-transform approach

The probabilistic approach "Cumulative Distribution Function-transform" (hereafter CDF-t) was developed for the correction of punctual daily wind speed and regional downscaling (e.g. Michelangeli et al., 2009; Vrac and Vaittinada, 2017). The CDF-t method has also been applied to correct biases of different atmospheric variables as for example; temperature, precipitation, relative humidity, among others (e.g. Colette et al., 2012; Vrac et al., 2012). Vrac et al. (2016) proposed a modification of the CDF-t method for bias correction, specifically

designed for precipitation, called "Singularity Stochastic Removal" (hereafter SSR). The motivation for developing an approach specialized for precipitation is because of its particular property in terms of a large number of zeros (no-precipitation events) in a daily time step. The principal advantage of this approach is that it is able to correct biases avoiding separating the correction in terms of occurrence (number of rainy days) and intensity of precipitation (quantity of precipitation). The SSR approach has been used to correct heat waves over France in Ouzeau et al. (2016) and also in a multivariate quantile mapping bias correction context to correct 3-hourly surface meteorological variables from the Canadian Centre for Climate Modelling and Analysis Regional Climate Model across a North American domain in Cannon (2017).

### Description of the CDF-t method

In our study, the CDF-t method aims at relating cumulative distribution functions (CDFs) of a climate variable (here the precipitation) from the WRF simulation to the CDF of this variable from the in-situ observation.

A mathematical transformation  $T$  is applied to the CDF of simulated precipitation to define a new CDF as close as possible to the CDF measured at the station.

Let  $F_{\text{simh}}$  and  $F_{\text{obsh}}$  define respectively the CDFs of two variables of interest from the WRF (subscript sim) and from a given station (subscript obs) over a historical calibration period (subscript h).

We assume that the transformation  $T$  allow us go from  $F_{\text{simh}}$  to  $F_{\text{obsh}}$ :

$$T(F_{\text{simh}}(x)) = F_{\text{obsh}}(x).$$

Replacing  $x$  by  $F_{\text{simh}}^{-1}(u)$ ,  $u \in [0, 1]$ , we obtain:

$$T(u) = F_{\text{obsh}}(F_{\text{simh}}^{-1}(u)),$$

which provides a simple definition of  $T$ . (See Appendix A for more details on the generalized inverse  $F^{-1}$ ).

When the observed and the simulated data have CDFs very different from each other, the domain of  $F_{\text{obsh}}$  can be theoretically restricted. For maximizing this domain, the simulated data  $\{x_{h,i}\}$  are normalized in order to be in  $[0, M_{ref}]$ , where  $M_{ref} = \max_{i \in \{1, \dots, n\}} y_{h,i}$  with  $\{y_{h,i}\}$  the observed data. The normalized time series for simulated data is:

$$\tilde{x}_{h,i} = x_{h,i} \frac{M_{ref}}{M_{Cmod}},$$

where  $M_{Cmod} = \max_{i \in \{1, \dots, n\}} x_{h,i}$ , it is the maximum value of the simulated time series. The same process is followed to normalize the simulated data  $\{x_{f,i}\}$  over a projection or validation time period:

$$\tilde{x}_{f,i} = x_{f,i} \frac{M_{ref}}{M_{Cmod}}.$$

Assuming that  $T$  is stationary in time, the transformation can be applied to  $F_{\text{simf}}$ , the CDF of the simulated precipitation over a validation or a future period  $f$ , to generate  $F_{\text{obsf}}$ , the CDF of the in situ precipitation for the same period  $f$ :

$$T(F_{\text{simf}}(x)) = F_{\text{obsf}}(x),$$

which is equivalent to:

$$F_{\text{obs}}(x) = F_{\text{obsh}}(F_{\text{simh}}^{-1}(F_{\text{simf}}(x))).$$

## The Singularity Stochastic Removal Approach

The approach developed by Vrac et al. (2016) is to correct daily punctual precipitation in terms of occurrence and intensity by replacing the 0 values of the observed and simulated time series into small random values uniformly distributed. First, a threshold  $th \in \mathbb{R}^+$  is chosen such that, all the positive values –of the observation and simulation– are greater than it. Then, each 0 value, in the observation and simulation, is changed by a uniformly distributed random variable  $v \sim U_{]0,th[}$ . Finally, the adjustment method CDF-t is applied and the corrected values of the simulated time series that are lower than  $th$  are set to 0.

The approach SSR was compared with three commonly used approaches of bias correction of precipitation: direct approach, threshold adaptation and positive approach. Briefly, the approach of *threshold adaptation* consists in finding a threshold  $th$  in ( $\text{mm day}^{-1}$ ) such that  $Pr(Obs = 0) = Pr(Sim \leq t)$ . Then all the values of the simulation lower than  $th$  are set to 0, and then, an adjustment method, as for example CDF-t is applied. The *positive approach* consists in correcting only the positive values of precipitation, so in this case no correction of the occurrence is done. Finally, the *direct approach* consists in applying the adjustment method in the complete time series.

Historically, the CDF-t method has been applied as a downscaling method and to correct future time series biases. But in this study, we also adapt the method to correct spatial precipitation data. The main idea is to partition the region into "neighbors sub-regions", in such a way that every sub-region contains a station. The precipitation biases in these sub-regions are "supposed" to behave similarly. To correct a simulated time series belonging to a given sub-region, a transformation  $T$  (described in the CDF-t method) is constructed by using as calibration time series observed-simulated time series at the station located in it. The strategy used to spacialize the CDF-t is detailed in the next section.

## 2.3 Spatial CDF-t approach

We propose two adaptations to the CDF-t method with SSR approach, the first adaptation is on the estimation of the CDFs based on time series, and the second one is the spatialization of the correction over the region under study using Voronoi polygons. The observed and simulated time series CDFs estimation is one of the essential steps in CDF-t method. Let  $(X_t)_{t \in \mathbb{N}^*}$  and  $(Y_t)_{t \in \mathbb{N}^*}$ , be the simulated and observed time series, respectively. A set of properties to estimate the CDFs have to be satisfied by  $(X_t)_{t \in \mathbb{N}^*}$  and  $(Y_t)_{t \in \mathbb{N}^*}$ . The  $(X_t)$ 's (similar to the  $(Y_t)$ 's) should be independent and identically distributed random variables. Tests to identify trends are carried out to obtain stationary time series and, an analysis of auto-correlation function of each time series is developed in order to have independent data. First, before proceeding, some basic time series concepts are presented, based on Rubenthaler (2017) and Jacques (2016).

### 2.3.1 Basic time series concepts and tests

**Definition 2.3.1.** A **time series** is a set of  $x_t \in \mathbb{R}$  with  $t \in \{1, \dots, n\}$  where  $t$  index represents a time unit, as for example; days, months or years. A time series  $x_t$  is a finite number of observations of a stochastic process  $(X_t)_{t \geq 0}$ . In this study, the analyzed time series are observed or simulated daily precipitation.

Until the moment, tests have been developed to prove a weaker form of stationary known as second-order stationary, that is defined as follows:

**Definition 2.3.2.** A stochastic process  $(X_t)_{t \geq 0}$  is **stationary of second order** if its expectation value is constant, i. e. it does not depend on  $t$  ( $E(X_t) = \mu \forall t$ ) and, its correlation

function, defined as  $\gamma(h) = \text{corr}(X_t, X_{t+h}) \forall t$ , only depends on the lag  $h$ , in other words it is independent of  $t$ .

**Definition 2.3.3.** Let's suppose  $X_t = m_t + s_t + \epsilon_t$ , where  $m_t$  and  $s_t$  are deterministic functions and  $\epsilon_t$  a random process with  $E(\epsilon_t) = 0$  and  $\text{cov}(X_t, X_{t+h}) = 0$  if  $h \neq 0$ .  $\epsilon_t$  is known as a white noise. The function  $m$  is the **trend** and  $s$  is the  $T$ -periodic **seasonal component** of  $X_t$ . For example, a simple trend case is when  $m$  is linear,  $m = a + bt$ .

In other words, in order to obtain stationary time series, its trend and seasonality must be removed. A commonly used method to remove a times series linear trend and seasonality is the difference method. The **difference method** consists in applying an operator  $\Delta_T : (X_t)_{t \in \mathbb{Z}} \rightarrow (X_t - X_{t-T})_{t \in \mathbb{Z}}$  over the time series  $X_t$ . Then, the time series  $\Delta_T X_t$  does not have linear trend and seasonality. Because, we have the following:

$$X_t - X_{t-T} = m_t - m_{t-T} + \epsilon_t - \epsilon_{t-T},$$

if we suppose there is a linear trend:  $m_t = at + b$ , therefore:

$$\begin{aligned} X_t - X_{t-T} &= at + b - a(t-T) - b + \epsilon_t - \epsilon_{t-T} \\ &= aT + \epsilon_t - \epsilon_{t-T} \end{aligned}$$

Some tests are made to identify that the time series under study have linear trend, in other words; the operator  $\Delta$  allow us to have time series without trend.

A time series dependence index is its empirical auto-correlation function defined as follows.

**Definition 2.3.4.** The **empirical auto-correlation function** (ACF) of a time series  $x_t$  is defined by:

$$\hat{\rho}(h) = \frac{\hat{\sigma}_n(h)}{\hat{\sigma}_n(0)}$$

where  $\hat{\sigma}(h) = \frac{1}{n-h} \sum_{t=1}^{n-h} (x_t - \bar{x}_n)(x_{t+h} - \bar{x}_n)$  and  $\bar{x}_n$  is the time series empiric mean.

**Definition 2.3.5.** A stochastic process  $X_t$  is **autoregressive of order  $p$**  ( $AR(p)$ ) if it can be represented as:

$$X_t = \epsilon_t + \sum_{j=1}^p a_j X_{t-j} \quad (2.11)$$

with  $\epsilon_t$  a white noise,  $p \in \mathbb{N}$  and  $a_p \in \mathbb{R} \setminus \{0\}$ .

For constructing an independent time series from the original one  $x_t$ , a lag  $h$  such that the estimated auto-correlation  $\hat{\rho}(h)$  is significant has to be found. Then, the sub-time series  $x'_t$  is constructed by skipping  $h$  positions in  $x_t$ . For example, if the auto-correlation is significant until  $h = 1$ , only the odd index  $t$  are kept  $x'_t = x_{2t}$ .

## Tests to identify trends and stationarity

There are tests to analyze the existence of a linear trend and stationary in a time series:

The **Mann-Kendall** non-parametric test is used to detect monotonic trends in time series, the hypothesis tested are:  $H_0$  : the data come from a population with independent realizations and are identically distributed and the alternative hypothesis  $H_1$  : the data follow a monotonic trend (Pohlert, 2016).

The **Augmented Dickey Fuller** (ADF) test shows if data is stationary, the hypothesis considered are:  $H_0$  : The data is not stationary and  $H_1$  : The data is stationary.

The **Kwiatkowski Phillips Schmidt Shin** (KPSS) test is used to analyze if data is stationary, similar to the Augmented Dickey Fuller test, but the hypothesis are:  $H_0$  : The data is stationary and  $H_1$  : The data is not stationary.

Following the spatialization strategy is explained and the complete algorithm is described.

### 2.3.2 The spatial CDF-t approach algorithm

A first approach to spatialize the CDF-t method is to divide the region using a partition based on Voronoi diagrams. The **Voronoi diagrams**, also known as Thiessen polygons, have been widely used in meteorological applications. As for example in Buytaert et al. (2006), a precipitation spatial interpolation with Thiessen polygons in the south Ecuadorian Andes is developed. In Ly et al. (2011), a spatial interpolation is carried out in the Ourthe and Ambleve catchments in Belgium. The method consists in dividing a region into polygons in such a way that every grid point is contained in the sub-region closest to a station (Aurenhammer and Klein, 2000; Barbulescu, 2016).

The following is based on Aurenhammer and Klein (2000). Let  $S$  be a set of  $n$  observations (with  $n \geq 3$ ) and  $D \subset \mathbb{R}^2$  a region. Let  $p, q \in S$ , then let  $B(p, q) = \{x \in D | d(p, x) = d(q, x)\}$ , where  $d$  is the Euclidean distance.  $B(p, q)$  is the perpendicular line through the center of the line segment from  $p$  to  $q$ .  $B(p, q)$  separates the halfplane  $D(p, q) = \{x \in D | d(p, x) < d(q, x)\}$ . The **Voronoi region** of  $p$  with respect to  $S$  is defined as:

$$\text{VR}(p, S) = \bigcap_{q \in S, q \neq p} D(p, q).$$

And finally, the Voronoi diagram of  $S$  is the following:

$$V(S) = \bigcup_{p, q \in S, p \neq q} \overline{\text{VR}(p, S)} \cap \overline{\text{VR}(q, S)}$$

The set  $S$  in our case is composed by the 26 in-situ meteorological stations.

### 2.3.3 The final procedure

The procedure to spatialize the precipitation bias correction method CDF-t with stochastic approach is now described. First, the Voronoi diagram with the meteorological stations points is layout, so  $n$  polygons are created. A polygon is represented by an observed and a simulated time serie,  $Y_t$  and  $X_t$  respectively. Then, each region grid point is assigned to a single polygon. The grid point also contains a simulated time series  $Z$ . Following, a threshold value of 1 mm day<sup>-1</sup> is applied to all data to avoid the recurrent problem of small precipitation values simulated by the WRF model. Other threshold values were tested but the best results are obtained with a threshold of 1 mm day<sup>-1</sup> in terms of precipitation occurrence, similar to the result obtained by Moure et al. (2016).

To obtain the CDFs for constructing the  $T$  transformation, new time series  $\Delta X_t$ ,  $\Delta Y_t$  and  $\Delta Z_t$  are created by applying the difference operator  $\Delta$ . Remember that,  $\Delta X_t = X_t - X_{t-1}$ , and this step is made because these time series have no trend and are stationary according to the previously tests. To avoid the problem that our time series have a large null values number, we have adapted the approach of Vrac et al. (2016). The procedure is next recalled:

1. Determine a threshold  $th > 0$  such that all the model and observational differentiated time series ( $\Delta X_t$ ,  $\Delta Y_t$  and  $\Delta Z_t$ ) in absolute value that are strictly positive and smaller than  $th$ .
2. Each null value is changed by a value  $v \sim U[-th, th]$  in the observation and simulation differentiated time series.

3. The method CDF-t is applied to the new time series to correct  $\Delta Z_t$ . The CDFs are constructed with the sub series previously defined, with the lag  $p$  obtained by the  $AR(p)$  model, but the transformation  $T$  is applied to all the time series  $\Delta Z_t$  and the  $\Delta Z'_t$  corrected time series is obtained.

Then, the new corrected precipitation variable  $Z'_t$  is calculated as:

$$Z'_t = Z_{t-1} + \Delta Z'_t$$

where  $t > 1$  and  $Z'_1 = Z_1$ . The values  $Z'_t$  that are smaller than 0 are set to 0 and finally, a threshold of 1 mm day<sup>-1</sup> is used to get the final corrected precipitation. Initially, the method CDF-t was developed to correct directly precipitation, but the main modification idea proposed is that we are now correcting the precipitation difference between two consecutive days ( $X_t - X_{t-1}$ ) in order to use data that have the mandatory properties to make a good CDFs estimation.

## 2.4 Comparison criteria

To compare the accuracy of the rainfall products created by these two methods (Gaussian process model and Spatial CDF-t approach), we use a **leave-one-out cross-validation** framework. The leave-one-out cross-validation consists of dividing a data set into a training and testing set recursively, by taking out one observation at each iteration from the data set. Then, the model is built with the remaining data. Finally, the accuracy is tested over the testing set –composed by one observation. In the particular case of this application, at each cross-validation iteration, a station is "removed" from the data set. Then, the model is built using the remaining stations and validated over the one station that was removed. We proceed recursively with all the stations.

We have computed several criteria, in terms of occurrence (number of rainy/no rainy days) and intensity of precipitation (precipitation quantity), to evaluate the approaches accuracy in a daily basis. The following criteria are commonly used in the literature as for example in Ochoa et al. (2014); Maussion et al. (2011); Murre et al. (2016); Vrac et al. (2016).

### Criteria related to the occurrence

A day is considered as a "rainy day" if its daily precipitation value is higher than 1 mm day<sup>-1</sup> –other threshold values were tested but the best performance between model and in-situ observations was obtained with 1 mm day<sup>-1</sup>. Then, a dummy variable is created for each time series to code the variable rainy/no rainy days. To calculate criteria related to occurrence of rainy day the contingency Table 2.1 is built where the value 1 codes a rainy day and the value 0 codes a no-rainy day.

		In-situ observation	
	Value	1	0
Simulation	1	A	B
	0	C	D

**Table 2.1:** Contingency table to evaluate the accuracy of the approaches. The value 1 codes a rainy day, and 0 codes a no-rainy day.

The criteria related to the occurrence calculated are:

The **false alarm rate** (FAR) is defined as the wrong number of rainy days simulated over the total number of rainy days simulated:

$$FAR = \frac{B}{A+B}.$$

The **probability of detection** (POD) is defined as the ratio between the number of rainy days simulated correctly and the total number of rainy days observed:

$$POD = \frac{A}{A+C}.$$

The **probability of false detection** (PODF) is the ratio between the number of rainy days incorrectly simulated over the number of no-rainy days of the observation:

$$PODF = \frac{B}{B+D}.$$

And finally, the **Heidke skill score** (HSS) is calculated as:

$$HSS = \frac{S - S_{ref}}{1 - S_{ref}}$$

where  $S = \frac{A+D}{n}$  and  $S_{ref} = \frac{(A+B)(A+C)+(B+D)(C+D)}{n^2}$

It could be interpreted as the simulation ability to be better or worst than a random simulation. A perfect product should have a FAR value of 0, a POD value of 1, a 0 PODF value and a HSS value of 1 (Maussion et al., 2011).

### Criteria related to the intensity

The following criteria are used to evaluate gridded products accuracy in terms of intensity: the **Kolmogorov-Smirnov Test** (KS) is a non-parametric test to compare two distributions, the maximal difference between them is calculated. The **Spearman correlation** coefficient, the **root mean square error** (RMSE), and the **mean bias** are computed. It is important also to know the percentage of data that is over the percentile 95 of the observation, in the case of a good precipitation product it should be close to 5%. And finally the predictivity coefficient  $Q_2$ , which can be interpreted as the percentage of the predictive ability of the model. The following is the definition of the criteria:

$$mean\_bias = \frac{1}{n} \sum_{i=1}^n (\hat{z}(x_i) - z(x_i)),$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{z}(x_i) - z(x_i))^2},$$

$$Q_2 = 1 - \frac{\sum_{i=1}^n (z(x_i) - \hat{z}(x_i))^2}{\sum_{i=1}^n (\bar{z} - z(x_i))^2},$$

where  $\hat{z}(x_i)$  is the prediction of the precipitation variable in a cross-validation framework at the location  $x_i$ ,  $z(x_i)$  is the observation in the point  $x_i$ ,  $\hat{\bar{z}}$  and  $\bar{z}$  are the observation and prediction expectations, respectively.

In this Chapter, we have described the bias correction methodologies applied in this study, the implementation details are presented in the Chapters III and IV and an intercomparison between methodologies is developed in Chapter V.



## Chapter 3

# Implementation of Gaussian Process Models to correct WRF precipitation biases

As mentioned before, the WRF model simulates precipitation biases. Our main objective in this chapter is to correct statistically WRF daily precipitation bias in the region under study during the 2014-2015 period. The first methodology implemented is to model WRF biases through Gaussian Process models, explained in Chapter II. This part of the study is inspired on the work of previous studies. Lichtenstern (2013) used Gaussian process models to interpolate temperature in Germany and described a didactic introduction to kriging in spatial statistics. Mourre et al. (2016) and Ochoa et al. (2014) have made an interpolation based on Gaussian process of precipitation in the Andes using as external drift WRF simulated precipitation. Finally, Gräler et al. (2012) elaborated an interpolation of the Particulate matter  $PM_{10}$  concentrations in Europe during 2009 and the strategies that we use for including the time are presented.

Traditionally, Gaussian Process modeling has been used to interpolate different atmospheric variables over a region (as for example, precipitation, temperature, among others, described in the previous studies). In this study, we model the WRF biases defined as the difference between the WRF simulation value and the observation:

$$\text{BIAS} = \text{WRF simulation} - \text{Observation}.$$

Then, we obtain a prediction of the bias in each point of the region ( $\widehat{\text{BIAS}}$ ) and we proceed to calculate the predicted precipitation ( $\widehat{\text{Precip.}}$ ) value as:

$$\widehat{\text{Precip.}} = \text{WRF simulation} - \widehat{\text{BIAS}}. \quad (3.1)$$

Before proceeding with the daily bias precipitation correction, a preliminary analysis over the annual cumulative precipitation registered (2014 and, 2015) are made to understand the accuracy of the Gaussian process models proposed during each year.

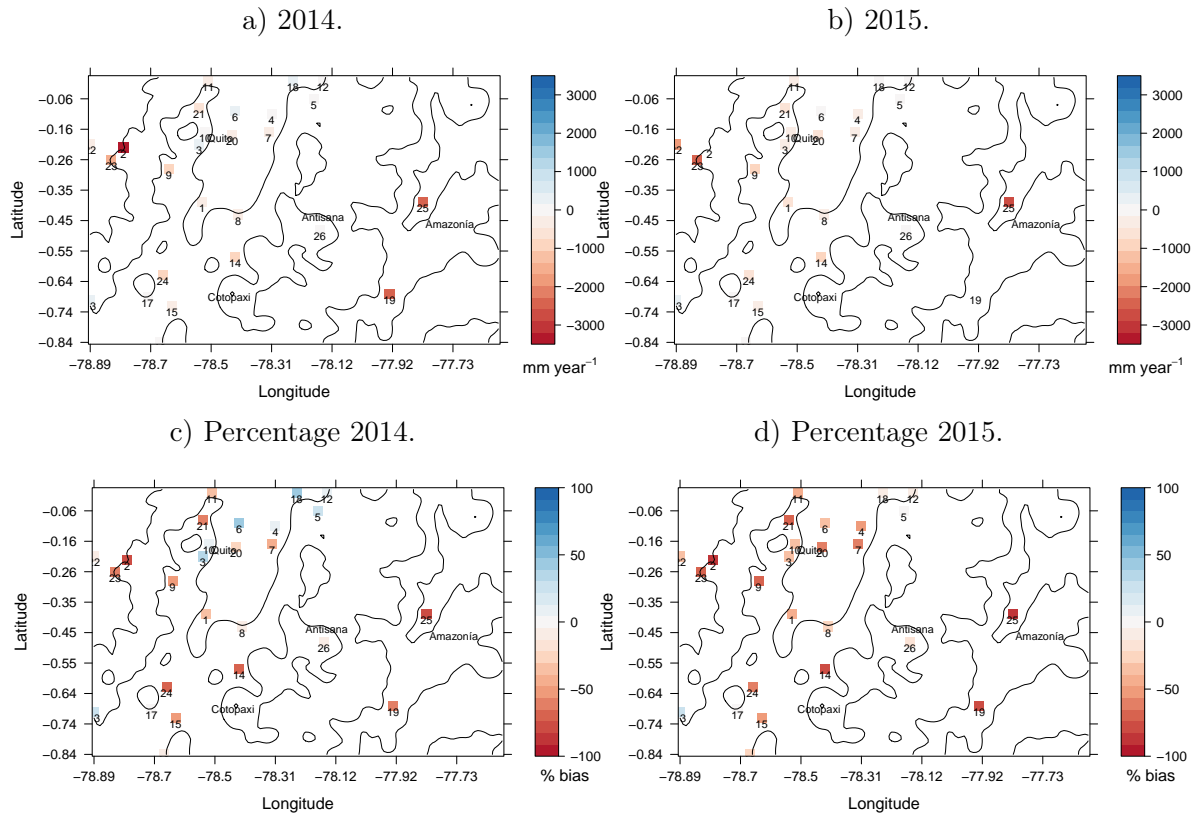
This Chapter is organized as follows: Section 1.1 presents an empirical analysis of the spatial distribution of the accumulated biases during the 2014 and 2015 periods. Section 1.2 describes the selection of the drifts ( $f(x)$  function). In Section 1.3 we study the results of an annual accumulated precipitation correction and finally, Section 1.4 presents the results of the daily bias precipitation correction for the 2014-2015 period.

### 3.1 The spatial distribution of WRF precipitation biases

The spatial bias distribution (our variable of interest  $Y$ ) during the 2014 and 2015 periods is empirically analyzed before proceeding with the Gaussian process modeling. There are tests

to analyze if the variable  $Z$  (obtained from the decomposition of  $Y$  into  $f$  and  $Z$ , described in Section 2.1) is stationary (see definition 2.1.2) that is an essential property to implement a Gaussian process model (Fuentes, 2005; Myers, 1989). But further implementation of these tests is needed, so they are out of the scope of this study.

Figures 3.1 a) and b) show the WRF simulated biases in terms of accumulated precipitation for the 2014 and 2015 years, respectively. The biases in percentage of 2014 and 2015 periods are shown in Figure 3.1 c) and d), respectively. The biases are more evident in the Amazon Region (stations 19 and 25), where an underestimation of precipitation of approximately  $3.000 \text{ mm year}^{-1}$  is simulated. The biases of 2014 and 2015 periods are slightly different because during 2014 period, there is an overestimation of the simulated precipitation in the region of Quito (North-West of the domain: stations 3, 6 and 18) in contrast to 2015, where an underestimation is displayed in all the stations.



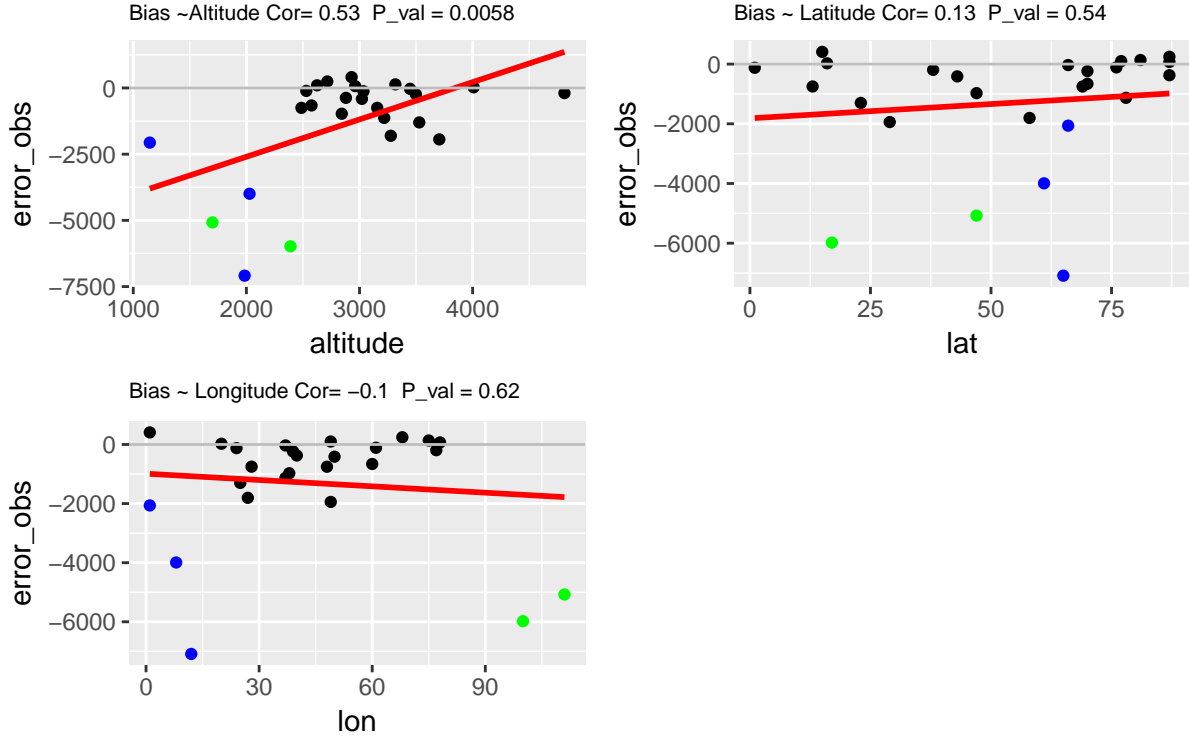
**Figure 3.1:** WRF precipitation biases for the a) 2014 and b) 2015 periods. WRF precipitation biases in percentage for the c) 2014 and d) 2015 periods.

### 3.2 Selection of the deterministic function $f(x)$

The variables longitude, latitude and altitude are commonly used as drifts (the  $f(x)$  function described in Section 2.1) in this type of studies. We follow the procedure developed by Hudson and Wackernagel (1994) to analyze their relevance as external drifts in this context. Figure 3.2 shows the linear relationship between our variable of interest (the bias) and three drifts (longitude, latitude and altitude) in a bi-annual basis, during the 2014-2015 period. Through this analysis, we have found that there is a notable difference between the bias amounts in the three geographic regions of Ecuador (Pacific, Andes and Amazon regions). Notice that, the bias values are smaller in the Andes region in contrast to the values obtained in the other two regions. This result is obtained because for this work it was chosen the WRF simulation

with the parameters that have provided the better precipitation results in the Andes Region.

Figure 3.2 shows that the linear regression between the bias and the altitude is the only significant ( $p\_value < 0.05$ ). In other words there is a linear significant relationship only between the altitude and bias.



**Figure 3.2:** Correlation diagram of bias with different drifts (which are: longitude, latitude and altitude). Blue points correspond to the stations belonging to the Pacific region, black points to the Andes, green points to the Amazon region. The linear regression between the bias and its external drifts is depicted in red.

### 3.3 Annual accumulated precipitation correction during the 2014 and 2015 periods

The accuracy of the Gaussian process (hereafter GP) models developed to correct the bias of the accumulated precipitation using three different drifts (altitude, longitude and latitude) are studied during the 2014 and 2015 periods, separately. Therefore, the criteria described in the Section 2.4 are calculated following a cross-validation leave-one-out framework over the region under study. Four corrected precipitation products are constructed by using the three drifts previously described:

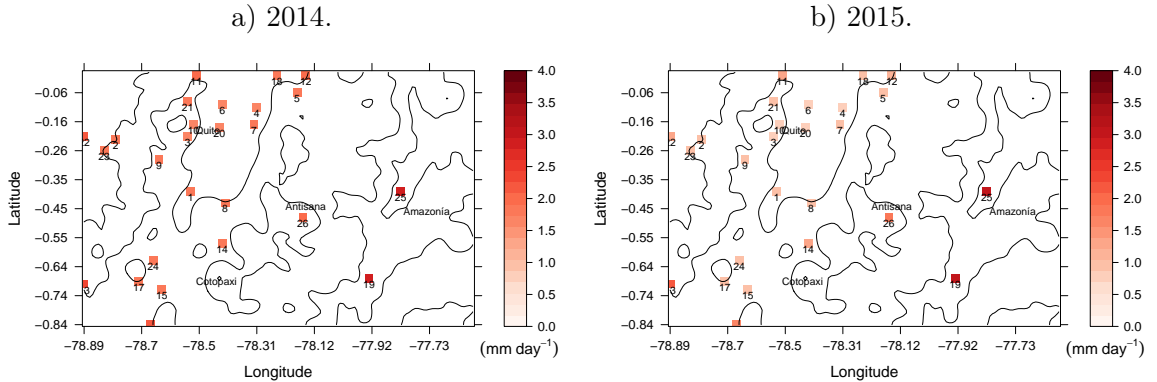
1. Gaussian process model with drift longitude, latitude and altitude (GP+longitude+latitude+alt).
2. Gaussian process model with drift longitude and latitude (GP+longitude+latitude).
3. Gaussian process model with drift longitude (GP+longitude).
4. Gaussian process model with drift altitude (GP+altitude).

Table 3.1 shows the cross-validation results over the three corrected precipitation gridded products. In the case of the criteria calculated for the WRF simulation, they are obtained directly from the registered value simulated in the grid point that corresponds to the station.

	2014				2015			
	Bias	RMSE	Correlation	Q <sub>2</sub>	Bias	RMSE	Correlation	Q <sub>2</sub>
WRF	1.77	2.89	0.59		2.19	3.67	0.65	
GP+longitude+latitude+alt.	1.71	2.31	0.71	0.44	1.56	2.14	0.80	0.65
GP+longitude+latitude	1.48	2.14	0.76	0.56	1.31	2.08	0.84	0.71
GP+longitude	1.50	2.14	0.76	0.49	1.32	2.11	0.84	0.64
GP+altitude	1.68	2.30	0.72	0.56	1.72	2.29	0.78	0.70

**Table 3.1:** Cross-validation leave-one-out results of annual accumulative precipitation for the three Gaussian Process models proposed with three drifts: longitude, latitude and longitude and latitude. The criteria are calculated for the 2014 and 2015 periods, separately.

All the three proposed GP models show better results in terms of the criteria calculated compared to the WRF simulation. But in general, the GP+longitude+latitude model obtains the best results in all the criteria (bias, RMSE, correlation and Q<sub>2</sub>). Thus, the GP model selected to correct the bias precipitation is the GP+longitude+latitude model. Figure 3.3 shows the daily standard deviation of the selected model (GP+longitude+latitude) during the 2014 and 2015 period. The Amazon stations (numbers 19 and 25) have higher variances values (3 to 4 mm day<sup>-1</sup>) than the other stations because one of the problems faced in the implementation of this methodology is the scarcity of observations and its in-homogeneous distribution. Notice that, the year 2014 (Figure 3.1 a)) presents biases values bigger than the biases values of 2015, this result is caused because the year 2015 was a dry year.



**Figure 3.3:** Daily standard deviation from the GP+longitude+latitude model for a) 2014 and b) 2015 periods.

### 3.4 The daily precipitation correction

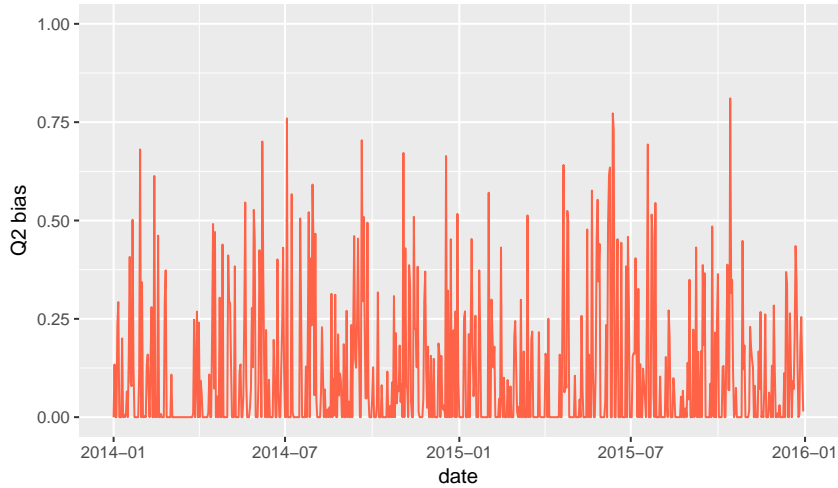
In this section we proceed with the daily correction following the two strategies described in Section 2.1.1: the separate daily variograms and the daily evolving variograms. The first strategy consists into creating a GP model for each day of the period 2014 and 2015. The second strategy, similar to the first one, consist into fitting a GP model to each day  $D$  but taking into account the parameters of the model of the previous day  $D - 1$  on the estimation.

To analyze the accuracy of the two strategies with GP+longitude+latitude using the separate variogram or the daily evolving variogram strategies, the mean predictivity coefficients Q<sub>2</sub> for each day of the 2014-2015 period are calculated and Table 3.2 shows the mean Q<sub>2</sub> coefficients calculated. The best accuracy is obtained with separate daily variograms because it obtains a mean accuracy of 11%. Figure 3.4 shows the Q<sub>2</sub> daily coefficients evolution along

the 2014-2015 period.

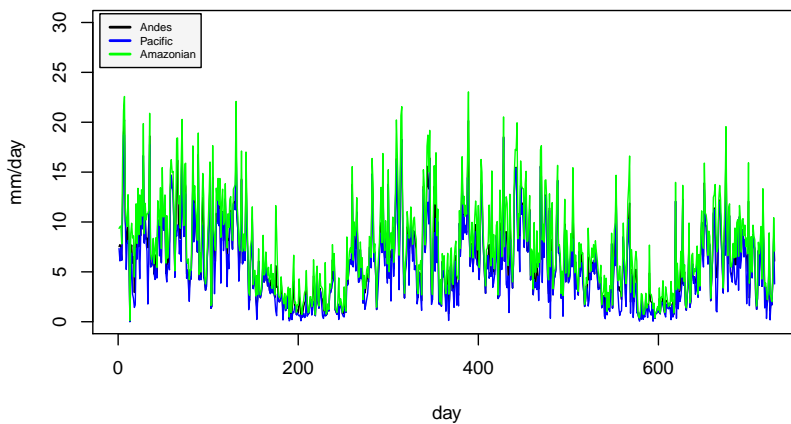
Q2	Separ. Variog.	Daily Evol. Variog.
GP+longitude+latitude	0.11	0.05

**Table 3.2:** Mean  $Q_2$  predictivity coefficient for two daily approaches during the 2014-2015 period.



**Figure 3.4:** Daily  $Q_2$  evolution from the 2014-2015 period.

Figure 3.5 shows the daily standard deviation evolution of the GP+longitude+latitude model at three stations located in each of the three geographical regions during the 2014-2015 period. The three stations show standard deviations that vary between 0 and 20 mm day<sup>-1</sup>. The time periods where the standard deviation values are higher correspond to rainy periods and, the other ones corresponds to dry periods in the region.



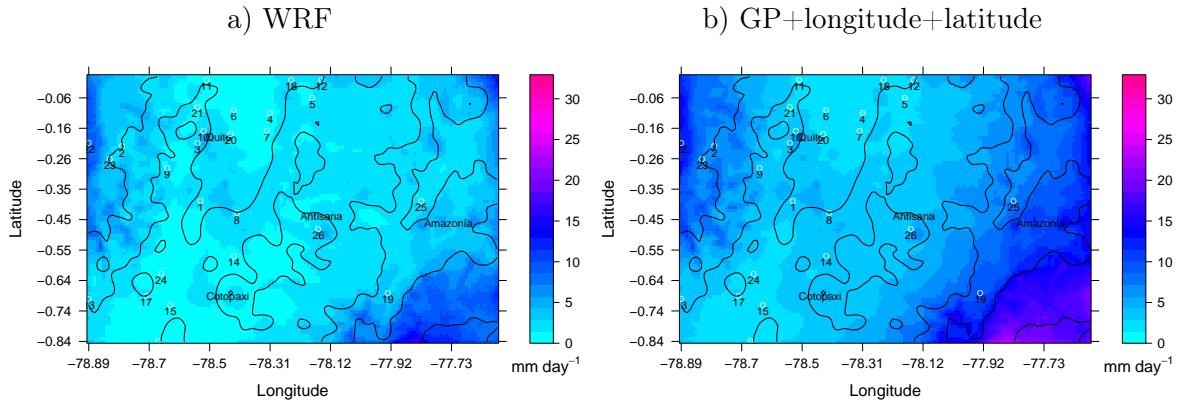
**Figure 3.5:** Standard deviation daily evolution for the prediction from the 2014-2015 period calculated in three stations: blue line for a Pacific station (number 22), green line for an Amazon station (number 25) and black line for an Andes (number 26).

Table 3.3 shows the criteria calculated over the accumulated corrected precipitation during 2014 and 2015 obtained by the GP model with the strategy separate variograms. Almost all the criteria are similar to the previous results obtained in Table 3.1, the bias criterion is slightly bigger than the annual correction, the RMSE values are smaller compared to those obtained in the annual correction and the  $Q_2$  values are similar to the previous values obtained. Thus, the method is coherent in an daily basis and an annual basis correction.

Period	BIAS	RMSE	Corr.	$Q_2$
2014	1.52	0.08	0.76	0.56
2015	1.54	0.09	0.80	0.61

**Table 3.3:** Criteria calculated over the accumulated precipitation quantity during 2014 and 2015 periods obtained from the results of the corrected precipitation by GP strategy separ. variog.

Finally, Figure 3.6 shows the mean daily precipitation during the 2014-2015 period of two gridded precipitation products: the WRF simulation and the GP+longitude+latitude model. The GP+longitude+latitude (hereafter GP) precipitation map has preserved the physical spatial patterns of the WRF simulation. A detailed comparison between the GP model and the spatial CDF-t approach is carried out in Chapter V.



**Figure 3.6:** Mean of daily precipitation [ $\text{mm day}^{-1}$ ] maps during the 2014-2015 period.

## Chapter 4

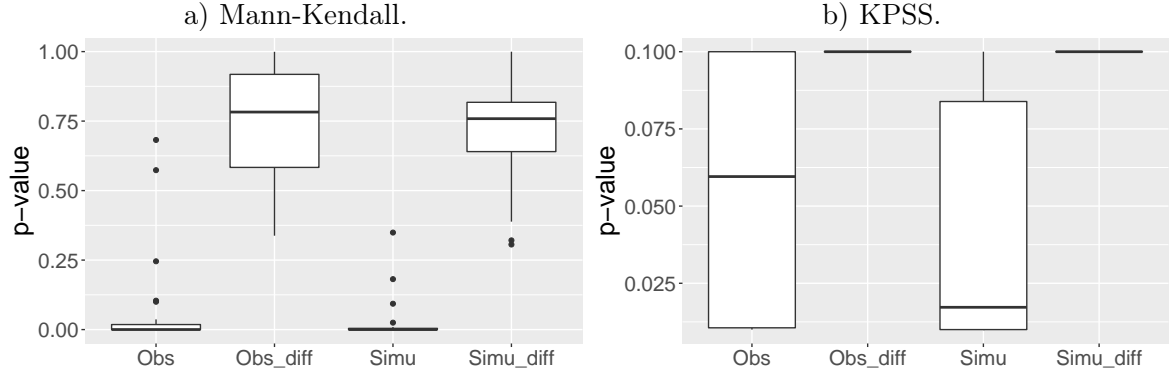
# Implementation of Cumulative Distribution Function-transform spatial approach

In previous studies, the CDF-t method has been applied as a downscaling method and to correct future time series biases (e.g. Michelangeli et al., 2009; Vrac et al., 2016). But in this study, we try to adapt the method to correct spatial precipitation data. The main idea is to partition the region into "neighbors sub-regions", in such a way that every sub-region has an associated observed-simulated time series. The precipitation biases in these sub-regions are "supposed" to have similar values. To correct a simulated time series belonging to a given sub-region, a transformation  $T$  (described in the CDF-t method) is constructed by using the observed-simulated time series CDFs associated with the sub-region and the grid point time series CDFs. In other words, instead of using a future time series as it is originally employed, we use time series belonging to an other grid of the sub-region assuming that the biases in all the sub-region grids are similar.

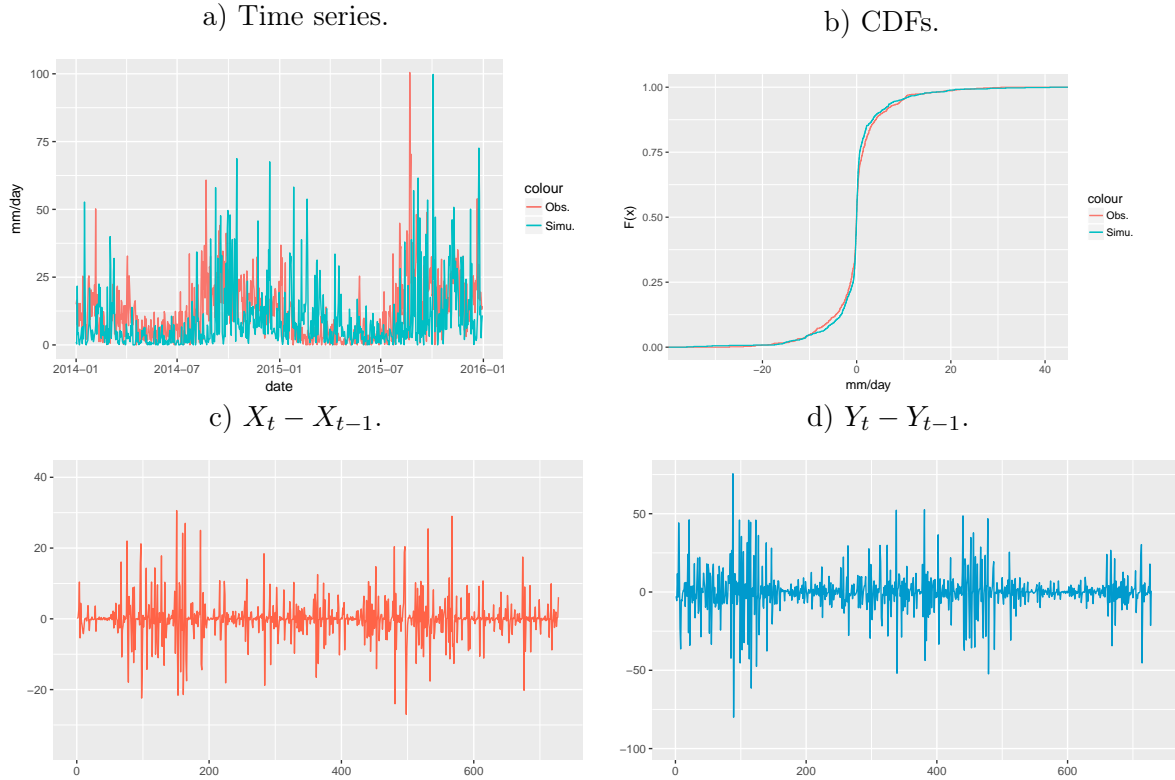
As it was mentioned before, each time series have to be stationary, independent and identically distributed random variables. Therefore, tests to identify trends are carried out to obtain stationary time series, and an analysis of auto-correlation function is developed. This Chapter is organized as follows: in Section 4.1, we analyze the stationarity and independence of the time series, and the spatial correction is studied in Section 4.2.

### 4.1 Time series analysis

The Mann-Kendall, ADF, and KPSS tests (described in Chapter II) are applied to observed and simulated time series. The ADF  $p\_value$  test is lower than 0.05 in all the time series, meaning that our data is stationary. But Mann-Kendall test results show that the hypothesis  $H_0$  is rejected, meaning that there is a monotonic trend in the data for almost all the time series except for four stations and three corresponding grid points in the simulations. The KPSS test results show that in all the time series the hypothesis  $H_0$  is rejected with a significance of 10%, in other words the data is not stationary. So, having this in mind, the operator difference  $\Delta X_t = X_t - X_{t-1}$  must be applied in the time series. We applied the operator difference to them and then, the Mann-Kendall and KPSS tests are recalculated in the differentiated time series. The two tests results are drawn in box-plots in Figure 4.1. According to these two tests, the differentiated time series ( $Y_t - Y_{t-1}$  and  $X_t - X_{t-1}$ ) are stationary and do not present monotonic trends. Figure 4.2 show the original time series, their respective differentiated series and CDFs of an Andes station (station number 26).



**Figure 4.1:** a) Box-plots of Mann-Kendall and b) KPSS  $p\_value$  results for observed and simulated time series (Obs: Observed time series  $Y_t$ , Obs\_diff:  $Y_t - Y_{t-1}$ , Simu: Simulated time series  $X_t$  and Simu\_diff:  $X_t - X_{t-1}$ ).

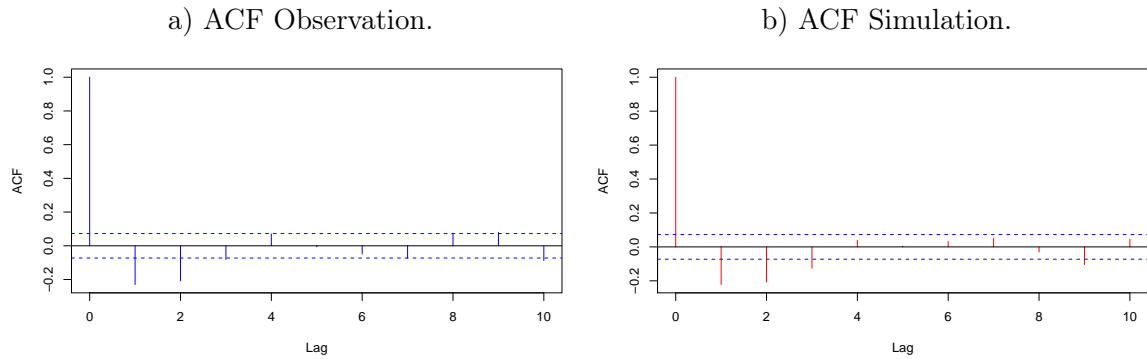


**Figure 4.2:** Time series treatment for an Andes region station (Antizana station number 26). a) Original ( $X_t$  observed in red and  $Y_t$  simulated in blue) time series. b) Differentiated series CDFs. c) differentiated observed and d) differentiated simulated time series.

The differentiated time series ( $\Delta X_t$  and  $\Delta Y_t$ ) almost accomplish the needed conditions to obtain CDFs estimations. The remaining condition to be accomplished is independence and it is obtained by modeling the  $\mathbf{AR}(p)$  process to identify a lag  $p$  until that each of the time series show a significant dependence pattern. In order to do that, models AR are fitted over the differentiated time series ( $\Delta X_t$  and  $\Delta Y_t$ ). We chose the model that minimize the Akaike information criterion (AIC). The AIC criterion is a measure of a statistical model accuracy based on a commitment between accuracy and the model parameters number used.



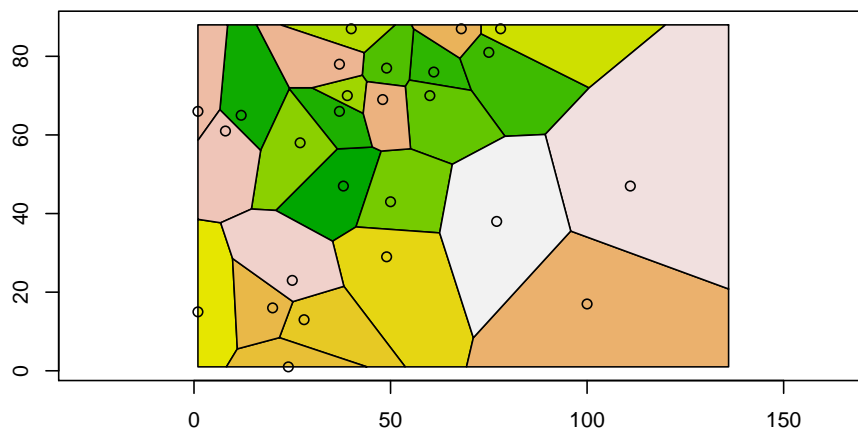
The AR models parameter estimation are obtained by using the function `auto.arima` of the R package `Forecast`. Once that the lag  $p$  in each time series has been identified, sub time series are constructed by skipping  $p$  positions in the unit time  $t$ . Figure 4.3 shows the ACF functions estimated for the differentiated observed and simulated time series for an Andes station (number 26). Data majority fits an AR(1) model (in other words, a lag of 1 is considered) and the higher lag  $p$  found is 3 for two time series in the data.



**Figure 4.3:** a) Differentiated observed time series ACF estimated and b) differentiated simulated time series ACF estimated

## 4.2 Spatial CDF-t approach results

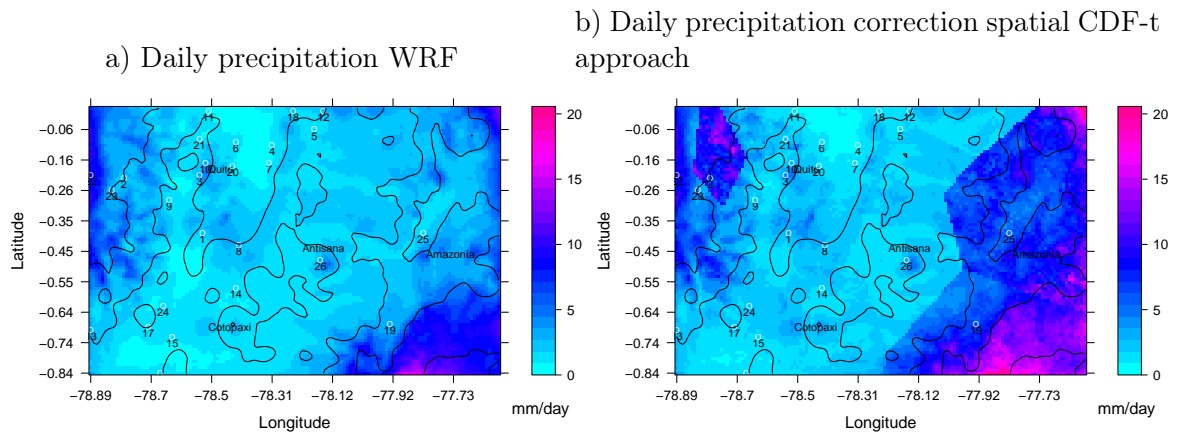
The Voronoi diagram is implemented in the region under study in Figure 4.4. Remark, as it was mentioned previously, that there is a in-homogeneous stations distribution and this fact produces that small polygons are constructed around the Andes region while the lack of stations in the Amazon region, with only two stations, produces bigger polygons. The lack of stations between the borders of the Amazon and the Andes region produce that the borders are to marked.



**Figure 4.4:** Voronoi diagram of 26 meteorological stations in the study region.

### 4.2.1 Spatial correction evaluation

The correction is applied to the region under study, and mean of daily precipitation map results are presented in Figure 4.5. The Voronoi polygon borders are completely marked in the Amazon region, due to the in-homogeneous distribution of the stations and also, the strong underestimated precipitation in this region (for example, around  $3.000 \text{ mm year}^{-1}$  at the station 25). In the Pacific Region, the border of the polygon associated to one station (station 22) is marked because it has recorded higher precipitation values. On the contrary, the polygons borders around the Andes region are not completely visible in most of the cases. As it has been seen in Chapter III, the biases values in the Andes region were quietly similar, so in this region the spatial CDF-t approach shows good results, by conserving the precipitation physical gradients well simulated by WRF. An homogeneous station distribution could increase the method accuracy by taking into account other variables in addition to geometrical properties. A deeper comparison between WRF, spatial CDF-t approach and GP is carried out in Chapter V.



**Figure 4.5:** Mean of daily precipitation [ $\text{mm day}^{-1}$ ] maps during the 2014-2015 period for **a)** the WRF simulation and, **b)** spatial CDF-t approach.

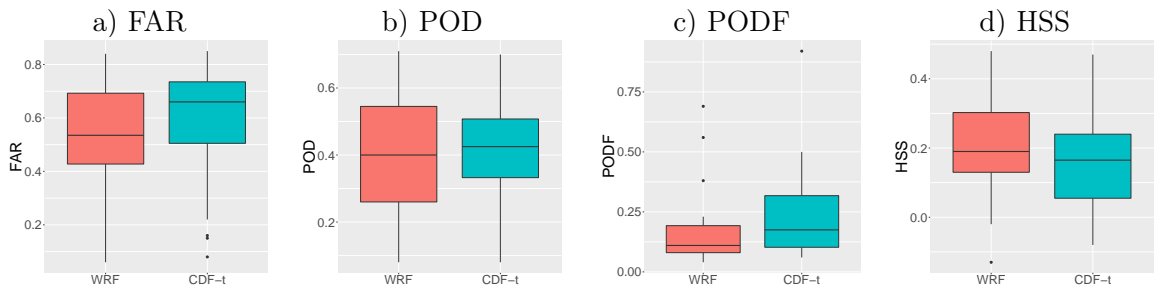
## Chapter 5

# Intercomparison between the CDF-t spatial method and Gaussian process model

After analyzing separately the implementation of the spatial CDF-t approach and the GP correction methodologies in this Chapter we proceed to make an intercomparison between these corrections and the original gridded products of the WRF simulation gridded. This Chapter is organized as follows: Section 5.1 describes a calibration-evaluation framework to analyze the accuracy of the modification proposed to the CDF-t method in Chapter IV. Section 5.2 shows a comparison in terms of criteria related to precipitation occurrence and intensity. Finally, a spatial graphical comparison between the resulting daily products and biases maps during the 2014-2015 period using the satellite product CHIRPS, is presented in Section 5.3.

### 5.1 Evaluation of the CDF-t in a "future" period

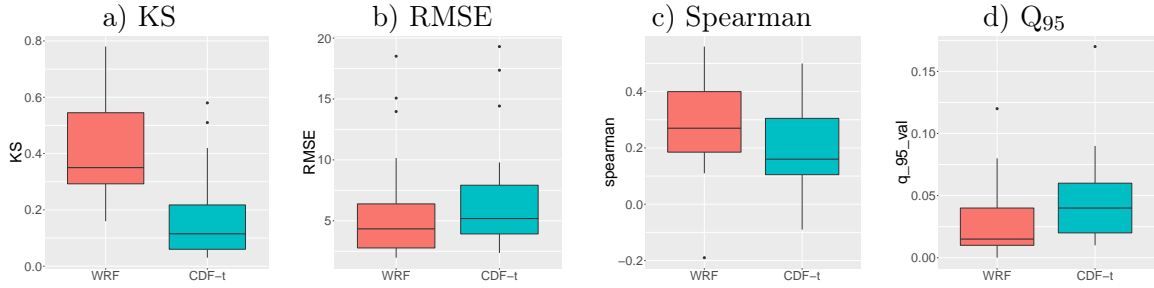
The modification of the method CDF-t method proposed in this study is evaluated in a calibration-evaluation framework for the 26 stations. To evaluate the modified method accuracy, it is calibrated over 01/2014 to 06/2015 and evaluated over 07/2014 to 12/2015. Then, the criteria related to precipitation occurrence and intensity are calculated to proceed with the accuracy evaluation.



**Figure 5.1:** Criteria related to precipitation occurrence (rainy/no-rainy events) in a calibration-evaluation framework (calibration over 01/2015 to 06/2016 and evaluation over 07/2016-12/2016). **a)** FAR criterion (ideal 0), **b)** POD criterion (1), **c)** PODF criterion (0) and **d)** HSS criterion (1).

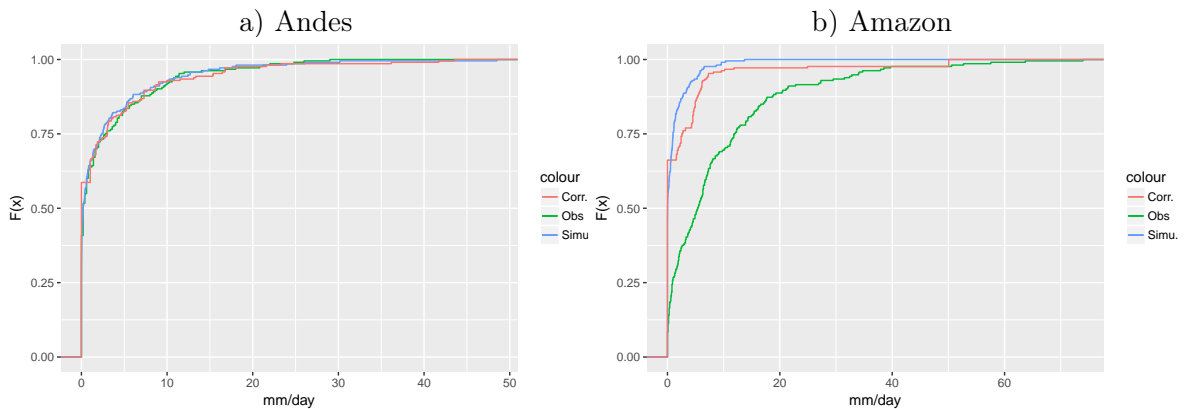
Figure 5.1 shows the results obtained in terms of criteria related to precipitation occur-

rence. The modified method results are slightly worse than the WRF original simulation in the majority of these criteria. The criteria related to the precipitation intensity are shown in Figure 5.2. A major improvement is obtained in the KS and  $Q_{95}$  criteria where almost all the stations show better results values in comparison to the WRF simulation. This is an expected result because the correction is carried out over the CDFs. The correlation criterion obtained by CDF-t of the corrected precipitation is reduced in contrast to the WRF simulation and in addition, the RMSE is slightly bigger in the corrected precipitation than in the original simulation.



**Figure 5.2:** Criteria related to precipitation intensity in a calibration-evaluation framework (calibration over 01/2015 to 06/2016 and evaluation over 07/2015-12/2015). **a)** KS, **b)** RMSE, **c)** Spearman correlation and **d)**  $Q_{95}$ .

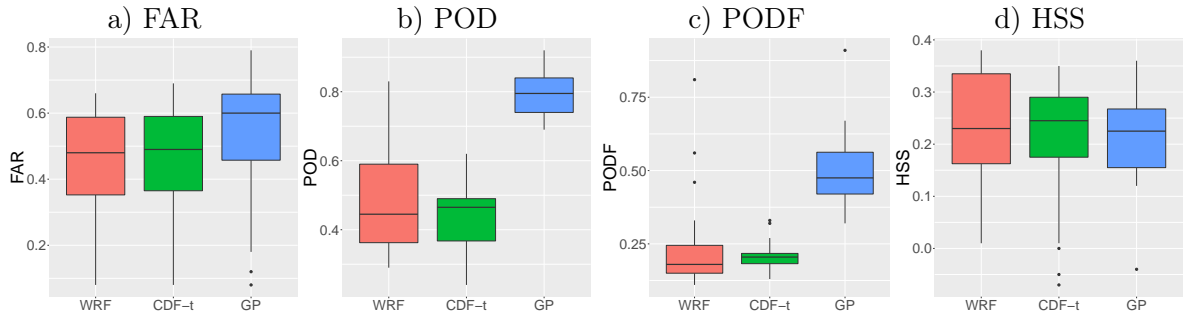
Finally, Figure 5.3 shows the simulation, observation and correction CDFs for two stations, one station located in the Andes region (number 26) and the other one, in the Amazon region (number 25). We choose this two ones because they have two different biases intensity. The simulation and observation CDFs in the Andes station are pretty similar in contrast to the CDFs in the Amazon region where there is a strong underestimation of precipitation. We can see that in both cases, the corrected precipitation CDFs (red) are closest to the observed precipitation CDFs (green) than the simulated WRF precipitation CDFs (blue), and that was the objective of the method. Notice that, the months chose for the evaluation period (07/2014 - 12/2015) belong to a sec "season" in contrast to the period of calibration that has two humid "seasons" and one sec "season", thus this could affect the accuracy of the evaluation.



**Figure 5.3:** Data CDFs of **a)** an Andes station (number 26) and **b)** an Amazon station (number 25).

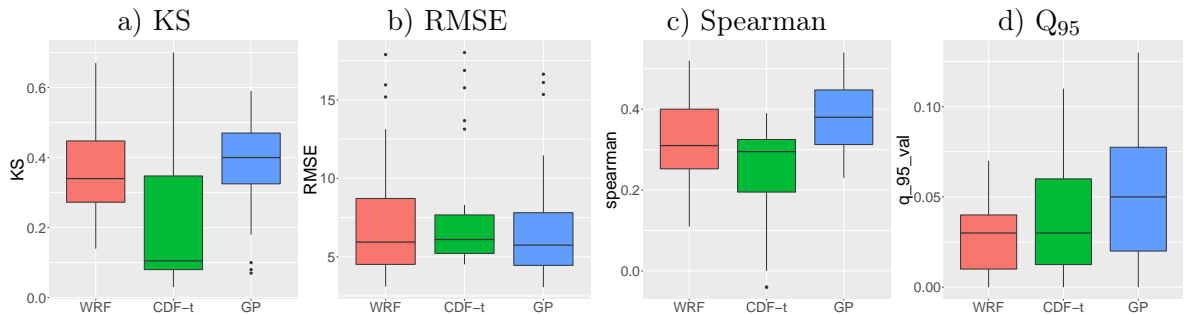
## 5.2 Evaluation of spatial correction methods

We use a cross-validation leave-one-out framework to compare the two correction (spatial CDF-t and GP) and WRF. The value used in the case of WRF is the one simulated in the corresponding grid point. The criteria related to precipitation occurrence are shown in Figure 5.4. The spatial CDF-t method shows similar results than WRF in FAR criterion, in contrast to GP where FAR results are worst than spatial CDF-t and WRF. POD criterion is highly improved by GP, in comparison to spatial CDF-t and WRF results. On the contrary, PODF results are worst in GP compared to WRF results. And finally, HSS criteria is worsened in both methods beside WRF results. In general, CDF-t method obtains, mostly, an improvement besides GP or at least, it obtains the same results as WRF. In other words CDF-t does not worsen WRF simulation in these criteria, except for HSS.



**Figure 5.4:** Boxplots of criteria related to precipitation occurrence (rainy/no-rainy events) for three gridded products: WRF, Spatial CDFt and GP, using a cross-validation leave-one-out framework. **a)** FAR criterion (ideal value 0), **b)** POD criterion (1), **c)** PODF criterion (0) and **d)** HSS criterion (1).

The results to criteria related to the precipitation occurrence are shown in Figure 5.5. The KS, RMSE and  $Q_{95}$  criteria are highly improved with the spatial CDF-t approach in contrast to GP. But on the contrary, Spearman correlation spatial CDF-t results are worst than WRF and also, than GP. Overall, GP model have shown better results, in comparison to spatial CDF-t, in terms of intensity and occurrence criteria.

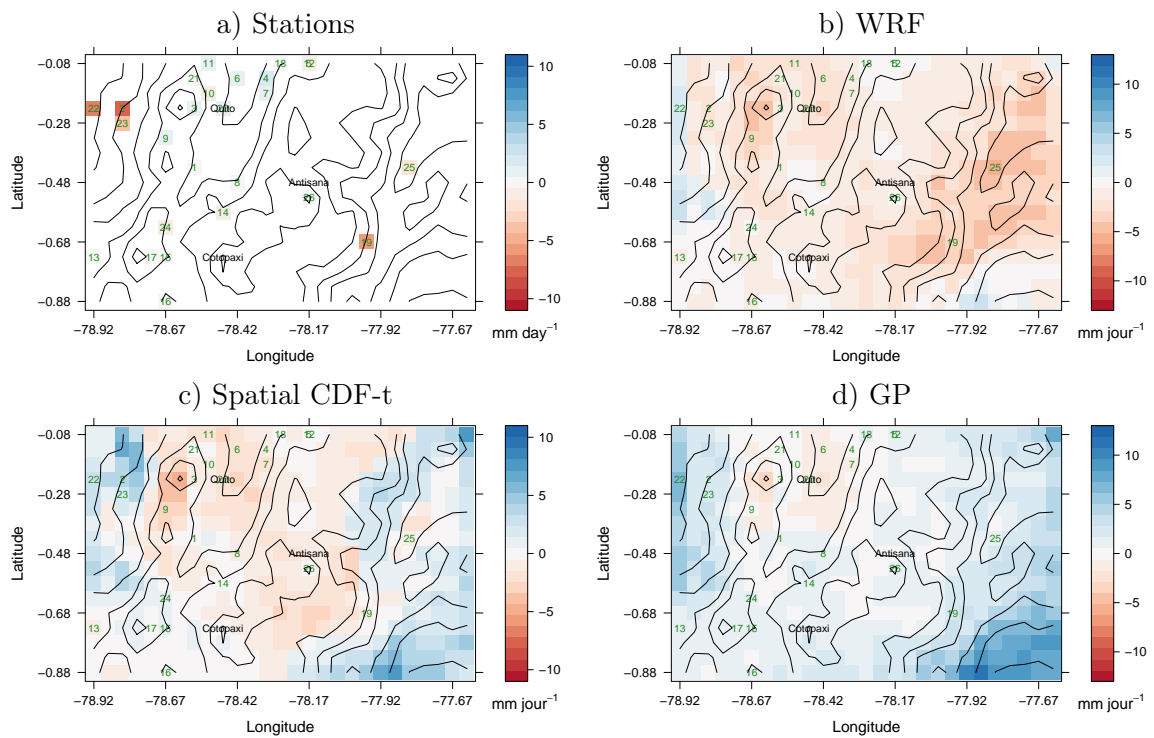


**Figure 5.5:** Boxplots of criteria related to precipitation intensity for three gridded products: WRF, spatial CDF-t and GP using a cross-validation leave-one-out framework. **a)** KS, **b)** RMSE, **c)** Spearman correlation and **d)**  $Q_{95}$ .

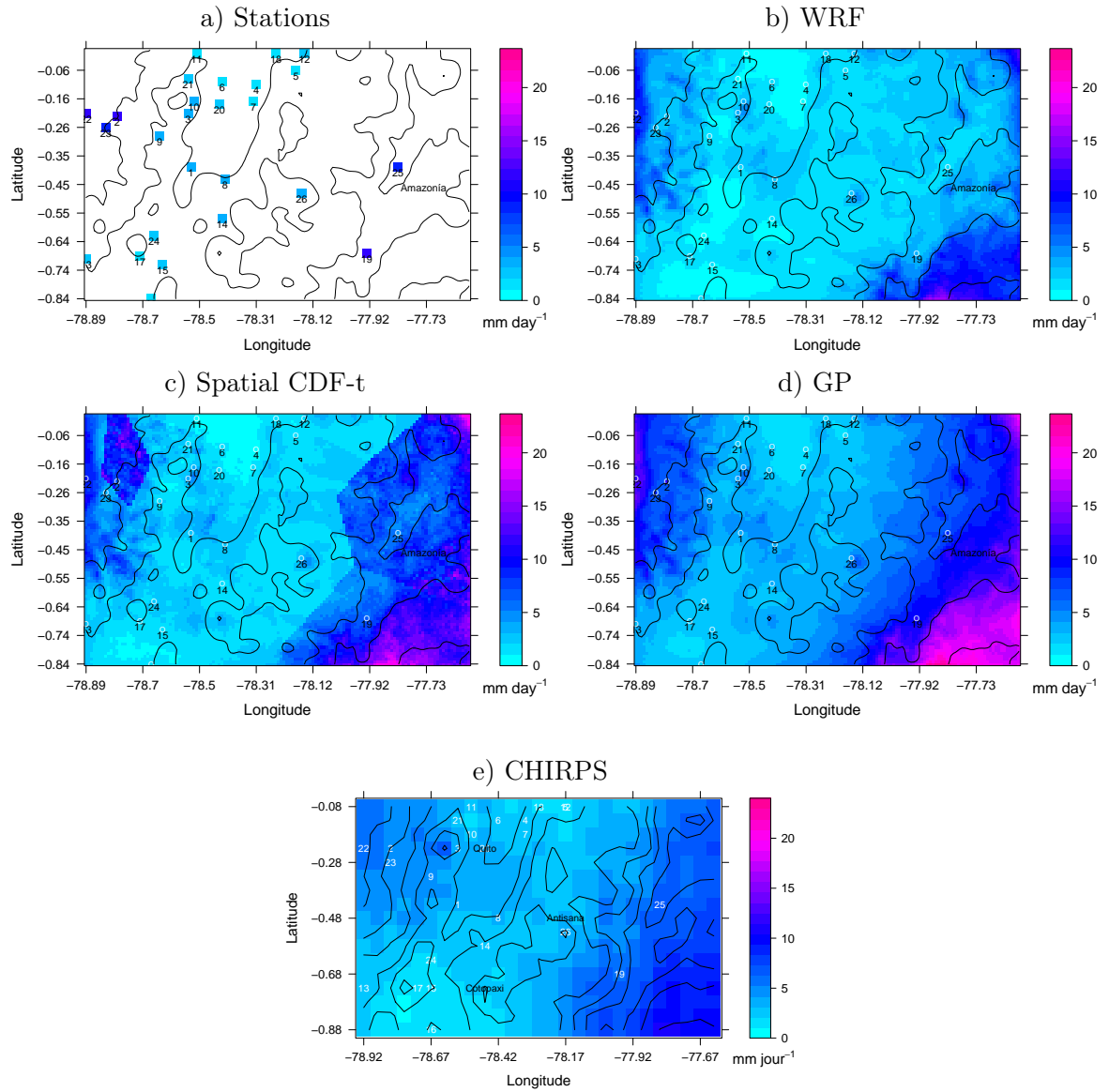
### 5.3 Precipitation gridded products

The mean of daily precipitations maps of the final products of precipitation for the 2014-2015 period (WRF, spatial CDF-t and GP) and in situ measures are presented in Figure 5.7. As it was pointed before, one of the spatial CDF-t approach disadvantages is that borders between polygons are too marked where there is a scarce number of stations (as it is the case in the Amazon region). But one of its advantages is that in sub regions where the stations are homogeneously distributed (for example, Andes region), the polygons borders are not visible and, the well simulated WRF spatial properties are preserved. On the contrary, with the GP model some precipitation spatial properties are lost, as for example, precipitation gradients in the Antisana glacier.

Finally, to compare the precipitation gridded products biases, they are compared to CHIRPS mean of daily precipitation during the 2014-2015 period. The biases precipitation maps are shown in Figure 5.6. Remember that CHIRPS, as other satellite products, tends to underestimate precipitation amounts. We also find this underestimation in Figure 5.6 a) where the difference between CHIRPS and observations is displayed. In most of the stations, CHIRPS precipitation is underestimated, except for 6 stations located in the Andes region where there is a slight overestimation. The CHIRPS biases are stronger in the Amazon and Pacific regions. The precipitation biases from WRF shows a strong underestimated precipitation in almost all the region, specially in the Amazon and to the western part of the region of Quito. The spatial CDF-t approach correction shows a decrease of this underestimation in the Amazon region, but small changes are found in the Andes region, where the biases are smaller. GP correction also shows a decrease of the WRF underestimation in almost all region, but also add a positive bias in the Amazon and the Pacific side, close to the borders of the domain, where no in-situ observations exists. Considering that the satellite product underestimates precipitation, this positive bias could be seen as most realistic that a negative one. Therefore, the GP method seems to be more realistic in most part of the Andes, except in the Amazon region and the Pacific side. However, a deeper study of the satellite biases is needed to confirm these results.



**Figure 5.6:** Mean of daily precipitation biases maps using as reference comparison precipitation from CHIRPS and gridded products: WRF, spatial CDF-t and GP during the 2014-2015 period. (Biases= CHIRPS-gridded product).



**Figure 5.7:** Mean of daily precipitation maps [ $\text{mm day}^{-1}$ ] during 2014-2015 period. a) In-situ measures, b) WRF, c) spatial CDF-t, d) GP, and e) CHIRPS.



## Chapter 6

# Conclusions and Perspectives

The aim of this study was to correct the precipitation biases of the WRF simulation in the Antisana region. Then, the final products of precipitation will be use as external forcing data for hydrological and glaciological models to understand water resources and glaciers evolution in the Andes. Therefore, two methodologies of precipitation bias correction were explored: the first one consisted into model statistically the daily WRF biases (defined as  $\text{BIAS} = \text{WRF simulation} - \text{Observation}$ ) through Gaussian Process (GP) models and, the second one was to corrected the biases by employing a spatial and time series adaptation of the CDF-t method developed by Michelangeli et al. (2009) and Vrac et al. (2016).

In first place, four GP models were proposed by using four external drifts  $f(x)$  (generally used in this type of studies: latitude, longitude and altitude) to model the annual accumulated bias during the years 2014 and 2015. The accuracy of the GP models was tested in a cross-validation leave-one-out framework. The best model was GP with drift longitude and latitude because it got the best results in the criteria calculated (Bias, RMSE, Correlation and  $Q_2$ ) during both years. Then, we explored the variance of the predictions, and the results show higher values in the Amazon Region where there is a non-homogeneous distribution and a scar number of station.

We chose the GP+longitude+latitude model to correct the daily precipitation. Therefore, we followed two strategies commonly used in the literature to obtain a daily correction: separating daily variograms and daily evolving variograms. The first strategy consists into creating a GP model for each day of the years 2014 and 2015. The second strategy, similar to the first one, consists into fitting a GP model to each day  $D$  but taking into account the parameters of the model of the previous day  $D - 1$  on the estimation. The separate daily variograms obtained the best  $Q_2$  mean result during the 2014-2015 period, thus we chose this one to correct the daily precipitation biases by the GP model.

One of the limitations of the current methods of bias correction of precipitation is the treatment of null values (no rainy days). As for example, the *threshold adaptation method* consists into finding a threshold  $t$ , such that simulation and observation have the same number of null values. Another example is the *positive approach method* where only the positive values of precipitation are corrected. Therefore, the SSR method developed by Vrac et al. (2016) allows to correct the precipitation in terms of occurrence and intensity without making subcases (dividing the correction into rainy periods and no rainy periods).

We employed the SSR method with a time series adaptation in order to obtain the CDFs estimation and a spatial adaptation to obtain the correction in the region. We proceeded as follows: first a threshold of  $1 \text{ mm day}^{-1}$  was applied to all the time series in order to solve the problem of small WRF precipitation and measurement errors in the observations (other threshold values were tested in a previous study but the value of  $1 \text{ mm day}^{-1}$  gives the best results). Then, we applied the differential operator ( $\Delta X_t : X_t - X_{t-1}$ ) over all the time series

(observed, simulated and to be corrected) in order to obtain time series without linear trend. Because the results of the Mann-Kendall test showed that these new time series do not have linear trend. Then, an  $AR(p)$  model is identified for each time series in order to recognize the lag  $p$  until that the time series are dependent, and sub series are constructed by skipping  $p$  positions. Next, we used the SSR approach over the differentiated time series ( $\Delta X_t$ ) and to recover the precipitation in a day  $t$ . We calculated the sum between the WRF simulation in the day  $t - 1$  in the corrected differentiated time series at the "period"  $t$ . Finally, a threshold of value  $1 \text{ mm day}^{-1}$  was used over the corrected time series.

The CDF-t method with the modification proposed was applied in a calibration-evaluation framework for the 26 stations. The method was calibrated over 01/2014 to 06/2015 and evaluated over 07/2014 to 12/2015. Then, the results between the WRF simulation and CDF-t modified method have shown that in terms of occurrence, the results of CDF-t were worsen in some criteria but the KS and  $Q_{95}$  criteria values were better in the CDF-t method with the modification proposed in comparison to the WRF simulation.

The two methodologies were compared in terms of precipitation occurrence (number of rainy/no-rainy days) and intensity (precipitation amounts) by applying a cross-validation leave-one-out framework. For comparing them, criteria related to the occurrence (FAR, POD, PODF and HSS) and criteria related to the intensity were calculated (mean bias, Spearman correlation, KS, RMSE and  $Q_2$  and  $Q_{95}$ ). In terms of almost all the criteria calculated, the GP model obtained the best results. The spatial CDF-t approach obtained the best results in terms of KS and  $Q_{95}$  that are criteria directly related to the CDFs. This result was expected because the spatial CDF-t approach correction is carried out over the CDFs.

Finally, the mean precipitation corrected maps obtained from spatial CDF-t and GP methodologies were compared with the satellite product CHIRPS. Considering that the satellite product underestimates precipitation, the GP method seemed to be more realistic in most part of the Andes, except in the Amazon region and the Pacific side. However, a deeper study of the satellite biases is needed to confirm these results.

There is still work to be done in the methodologies here presented to increase its accuracy. Thus, the perspectives of this study are the following ones: to include in the GP modeling the time dimension by proceeding as in Gräler et al. (2016). A second one is to analyze deeply the implementation of stationary tests for a GP model.

For the spatial CDF-t approach, other spatialization strategies should be implemented that include not only geometrical properties, as it was the case of the Voronoi polygons. For example, it could improve the results of the spatial CDF-t to construct clusters of the region under study. One alternative is to use the Functional Clustering Method as in Antoniadis et al. (2012) where a curve-based clustering is used to reduce the data dimension for constructing a metamodel for West African monsoon. The Functional Clustering method has the advantage of taking into account time-point correlations of time series spatial data (Antoniadis et al., 2012). Therefore, during this study a WRF simulation of 10 years (2005-2015 period) was made to apply the Functional Clustering Method over the region under study.

As mentioned before, the gridded precipitation corrected products are used as external forcing data for hydrological and glaciological models. Thus, an evaluation of the results of the models by using these two products could show the weakness or strength of them.

# Appendix A

## The generalized inverse $F^{-1}$

The following is obtained from Bercu and Chafai (2007). Let  $F : \mathbb{R} \rightarrow [0, 1]$  be a cumulative distribution function.  $F$  is increasing and right-continuous. The **generalized inverse function** of  $F$ , also called quantile function, denoted by  $F^{-1}$  is defined as:

$$F^{-1}(u) = \inf\{x \in \mathbb{R} \text{ such that } F(x) \geq u\}, \quad \forall 0 < u \leq 1.$$

**Theorem 1.** (*The inverse method*) If  $\mu$  is a probability distribution on  $\mathbb{R}$  with cumulative distribution function  $F$ , and if  $U$  is a random variable with uniform distribution over  $[0, 1]$ , then the random variable  $F^{-1}(U)$  is distributed as  $\mu$ .

*Proof.* Let's prove that for all  $x \in \mathbb{R}$  and  $0 < u \leq 1$ :

$$u \leq F(x) \iff F^{-1}(u) \leq x.$$

If  $u \leq F(x)$ , then  $x \in \{t \in \mathbb{R} | F(t) \geq u\}$  and thus  $x \geq F^{-1}(u)$  by definition of  $F^{-1}(u)$ . Now, let's suppose that  $F^{-1}(u) \leq x$ . As  $F$  is an increasing function, we have  $F(F^{-1}(u)) \leq F(x)$ . As  $F$  is right continuous, we have  $u \leq F(F^{-1}(u))$  which implies that  $u \leq F(x)$ . Therefore,

$$\mathbb{P}(F^{-1}(U) \leq x) = \mathbb{P}(U \leq F(x)) = F(x) \quad \forall x \in \mathbb{R},$$

The random variable  $F^{-1}(U)$  has the same cumulative distribution function  $F$  as  $X$ , therefore, it has the same probability distribution as  $X$ .  $\square$

**Theorem 2.** (*Continue distribution functions*) If a random variable  $X$  has a distribution function  $F$ , then the next properties are equivalent:

1.  $F$  is continuous over  $\mathbb{R}$ .
2.  $F(X)$  follows a uniform distribution over  $[0, 1]$ .
3.  $F(\mathbb{R}) = [0, 1]$ .

*Proof.* If  $F$  is a continuous function, the function  $F$  is not necessarily invertible. However,  $F$  is left-and-right continuous,  $u \leq F(F^{-1}(u)) \leq u$  for all  $0 \leq u \leq 1$ , therefore  $F(F^{-1}(u)) = u$ . Then, for all  $0 \leq u \leq 1$ :

$$\mathbb{P}(F(X) \leq u) = \mathbb{P}(X \leq F^{-1}(u)) = F(F^{-1}(u)) = u$$

and  $F(X)$  follows an uniform distribution over  $[0, 1]$ .  $\square$

# References

- Anestis Antoniadis, Céline Helbert, Clémentine Prieur, and Laurence Viry. Spatio-temporal metamodeling for west african monsoon. *Environmetrics*, 23(1):24–36, 2012.
- Franz Aurenhammer and Rolf Klein. Voronoi diagrams. *Handbook of computational geometry*, 5:201–290, 2000.
- Alina Barbulescu. *Studies on Time Series Applications in Environmental Sciences*. Springer Publishing Company, Incorporated, 1st edition, 2016. ISBN 3319304348, 9783319304342.
- Rubén Basantes-Serrano. *Evolution of glaciers in the Ecuadorian Andes since the 1950s and its contribution to the study of the climate change in the inner tropics*. Theses, Université Grenoble Alpes, July 2015. URL <https://tel.archives-ouvertes.fr/tel-01219778>.
- Bernard Bercu and Djalil Chafai. *Modélisation stochastique et simulation - Cours et applications*. Mathématiques appliquées pour le Master - Sciences Sup. Dunod, October 2007. URL <https://hal.archives-ouvertes.fr/hal-00669263>.
- Wouter Buytaert, Rolando Celleri, Patrick Willems, Bert De Bièvre, and Guido Wyseure. Spatial and temporal rainfall variability in mountainous areas: A case study from the south ecuadorian andes. *Journal of Hydrology*, 329(3-4):413–421, 2006. ISSN 0022-1694. doi: <https://doi.org/10.1016/j.jhydrol.2006.02.031>. URL <http://www.sciencedirect.com/science/article/pii/S0022169406001144>.
- Alex J. Cannon. Multivariate quantile mapping bias correction: an n-dimensional probability density function transform for climate model simulations of multiple variables. *Climate Dynamics*, pages 1–19, 2017. ISSN 1432-0894. doi: 10.1007/s00382-017-3580-6. URL <http://dx.doi.org/10.1007/s00382-017-3580-6>.
- Climate Hazard Group. Chirps satellite product. <http://chg.geog.ucsb.edu/data/chirps/>. [Online; accessed 27-May-2017].
- A. Colette, R. Vautard, and M. Vrac. Regional climate downscaling with prior statistical correction of the global climate forcing. *Geophysical Research Letters*, 39(13), 2012. ISSN 1944-8007. doi: 10.1029/2012GL052258. URL <http://dx.doi.org/10.1029/2012GL052258>. L13707.
- M. Déqué. Frequency of precipitation and temperature extremes over France in an anthropogenic scenario: Model results and statistical correction according to observed values. *Global and Planetary Change*, 57:16–26, May 2007. doi: 10.1016/j.gloplacha.2006.11.030.
- V. Favier, P. Wagnon, J. P. Chazarin, L. Maisincho, and A. Coudrain. One-year measurements of surface heat budget on the ablation zone of antizana glacier 15, ecuadorian andes. *Journal of Geophysical Research*, 2004.

- Montserrat Fuentes. A formal test for nonstationarity of spatial stochastic processes. *Journal of Multivariate Analysis*, 96(1):30 – 54, 2005. ISSN 0047-259X. doi: <http://dx.doi.org/10.1016/j.jmva.2004.09.003>. URL <http://www.sciencedirect.com/science/article/pii/S0047259X04001757>.
- Jason P. Giovannettone and Ana P. Barros. Probing regional orographic controls of precipitation and cloudiness in the central andes using satellite data. *Journal of Hydrometeorology*, 10(1):167–182, 2009. doi: 10.1175/2008JHM973.1. URL <http://dx.doi.org/10.1175/2008JHM973.1>.
- Benedikt Gräler, Mirjam Rehr, Lydia Gerharz, and Edzer Pebesma. Spatio-temporal analysis and interpolation of pm10 measurements in europe for 2009. *ETC/ACM Technical Paper*, 8:1–29, 2012.
- Benedikt Gräler, Edzer Pebesma, and Gerard Heuvelink. Spatio-temporal interpolation using gstat. *The R Journal*, 8:204–218, 2016. URL <https://journal.r-project.org/archive/2016-1/na-pebesma-heuvelink.pdf>.
- M Hall, P. Mothes, J. Aguilar, J. Bustillos, P. Ramón, J. P. Eissen, M Monzier, C. Robin, Egred. J., A. Militzer, and H. Yepes. Los peligros volcánicos asociados con el Antisana, 2012.
- Ratchawatch Hanchoo Wong, Uruya Weesakul, and Siriluk Chumchean. Bias correction of radar rainfall estimates based on a geostatistical technique. *ScienceAsia*, 38:373–385, 2012.
- Gordon Hudson and Hans Wackernagel. Mapping temperature using kriging with external drift: Theory and an example from scotland. *International Journal of Climatology*, 14(1): 77–91, 1994. ISSN 1097-0088. doi: 10.1002/joc.3370140107. URL <http://dx.doi.org/10.1002/joc.3370140107>.
- Julien Jacques. *Introduction aux Séries Temporelles (avec R)*. Université de Lyon 2, 2016.
- A. Lichtenstern. *Kriging methods in spatial statistics*. PhD thesis, Bachelor’s thesis, Technische Universität München, August 2013.
- S. Ly, C. Charles, and A. Degré. Geostatistical interpolation of daily rainfall at catchment scale: the use of several variogram models in the ourthe and ambleve catchments, belgium. *Hydrology and Earth System Sciences*, 15(7):2259–2274, 2011. doi: 10.5194/hess-15-2259-2011. URL <http://www.hydrol-earth-syst-sci.net/15/2259/2011/>.
- A. Marrel, B. Iooss, F. Van Dorpe, and E. Volkova. An efficient methodology for modeling complex computer codes with Gaussian processes. *Computational Statistics & Data Analysis*, 52(10):4731–4744, 2008.
- F. Maussion, D. Scherer, R. Finkelburg, J. Richters, W. Yang, and T. Yao. WRF simulation of a precipitation event over the tibetan plateau, China an assessment using remote sensing and ground observations. *Hydrology and Earth System Sciences*, 15(6):1795–1817, 2011. doi: 10.5194/hess-15-1795-2011. URL <http://www.hydrol-earth-syst-sci.net/15/1795/2011/>.
- P. A. Michelangeli, M. Vrac, and H. Loukos. Probabilistic downscaling approaches: Application to wind cumulative distribution functions. *Geophysical Research Letters*, 36(11), 2009. doi: 10.1029/2009GL038401.

- L. Mourre. *Précipitations dans les Andes tropicales: analyse spatio-temporelle, intercomparaison de forçages et impacts dans un modèle glacio-hydrologique. Cas du Rio Santa au Pérou*. PhD thesis, Université de Grenoble, 2015.
- L. Mourre, T. Condom, C. Junquas, T. Lebel, J. E. Sicart, R. Figueroa, and A. Cochachin. Spatio-temporal assessment of wrf, trmm and in situ precipitation data in a tropical mountain environment (cordillera blanca, peru). *Hydrology and Earth System Sciences*, 20(1):125–141, 2016. doi: 10.5194/hess-20-125-2016. URL <http://www.hydrolog-earth-syst-sci.net/20/125/2016/>.
- Marc F Müller and Sally E Thompson. Bias adjustment of satellite rainfall data through stochastic modeling: Methods development and application to nepal. *Advances in Water Resources*, 60:121–134, 2013.
- D. E. Myers. To be or not to be... stationary? that is the question. *Mathematical Geology*, 21(3):347–362, 1989. ISSN 1573-8868. doi: 10.1007/BF00893695. URL <http://dx.doi.org/10.1007/BF00893695>.
- A. Ochoa, L. Pineda, P. Crespo, and P. Willems. Evaluation of TRMM 3B42 precipitation estimates and WRF retrospective precipitation simulation over the Pacific-Andean region of Ecuador and Peru. *Hydrol. Earth Syst. Sci.*, 18:3179–3193, 2014. doi: 10.5194/hess-18-3179-2014.
- A. Ochoa, L. Campozano, E. Sánchez, R. Gualán, and E. Samaniego. Evaluation of downscaled estimates of monthly temperature and precipitation for a Southern Ecuador case study. *International Journal of Climatology*, 36:1244–1255, 2016. doi: 10.5194/hess-18-3179-2014.
- G. Ouzeau, J.-M. Soubeyroux, M. Schneider, R. Vautard, and S. Planton. Heat waves analysis over France in present and future climate: Application of a new method on the EURO-CORDEX ensemble. *Climate Services*, 4:1 – 12, 2016. ISSN 2405-8807. doi: <https://doi.org/10.1016/j.cliser.2016.09.002>. URL <http://www.sciencedirect.com/science/article/pii/S2405880716300309>.
- Christopher J. Paciorek and Mark J. Schervish. Spatial modelling using a new class of nonstationary covariance functions. *Environmetrics*, 17(5):483–506, 2006. ISSN 1099-095X. doi: 10.1002/env.785. URL <http://dx.doi.org/10.1002/env.785>.
- Thorsten Pohlert. Non-Parametric Trend Tests and Change-Point Detection. *R*, 116, 2016. URL <https://CRAN.R-project.org/package=trend>.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005. ISBN 026218253X.
- Sylvain Rubenthalès. *Séries Chronologiques (avec R)*. Université Nice Sophia Antipolis, 2017.
- J. Sicart, R. Hock, P. Ribstein, M. Litt, and E. Ramirez. Analysis of seasonal variations in mass balance and meltwater discharge of the Tropical Zongo Glacier by application of a distributed energy balance model. *J. Geophys. Res.*, 116, 2011. doi: 10.1029/2010JD015105.
- M. Vrac and P. Friederichs. Multivariate-Intervariable, Spatial, and Temporal-Bias Correction. *Journal of Climate*, 28(1):218–237, 2015. doi: 10.1175/JCLI-D-14-00059.1.

- M. Vrac and P. Vaittinada. Influence of bias correcting predictors on statistical downscaling models. *Journal of Applied Meteorology and Climatology*, 56(1):5–26, 2017. doi: 10.1175/JAMC-D-16-0079.1. URL <http://dx.doi.org/10.1175/JAMC-D-16-0079.1>.
- M. Vrac, P. Drobinski, A. Merlo, M. Herrmann, C. Lavaysse, L. Li, and S. Somot. Dynamical and statistical downscaling of the french mediterranean climate: uncertainty assessment. *Natural Hazards and Earth System Sciences*, 12(9):2769–2784, 2012. doi: 10.5194/nhess-12-2769-2012. URL <http://www.nat-hazards-earth-syst-sci.net/12/2769/2012/>.
- M. Vrac, T. Noel, and R. Vautard. Bias correction of precipitation through Singularity Stochastic Removal: Because occurrences matter. *Journal of Geophysical Research: Atmospheres*, 2016.