

The University of Chicago
Center for Integrating Statistical and Environmental Science
www.stat.uchicago.edu/~cises



Chicago, Illinois USA

TECHNICAL REPORT NO. 36

**STOCHASTIC DOWNSCALING OF PRECIPITATION:
FROM DRY EVENTS TO HEAVY RAINFALLS**

Mathieu Vrac and Philippe Naveau

June 2006



Although the research described in this article has been funded wholly or in part by the United States Environmental Protection Agency through STAR Cooperative Agreement #R-82940201 to The University of Chicago, it has not been subjected to the Agency's required peer and policy review and therefore does not necessarily reflect the views of the Agency, and no official endorsement should be inferred.

Stochastic downscaling of precipitation:

From dry events to heavy rainfalls

M. VRAC^{a,b}

P. NAVEAU^b

^a Center for Integrating Statistical and Environmental Science, The University of Chicago,
Chicago, IL 60637, USA.

^b Laboratoire des Sciences du Climat et de l'Environnement, IPSL-CNRS, Saclay, France

1 Abstract

2 Downscaling precipitation is a difficult challenge for the climate community. We propose
3 and study a new stochastic weather typing approach to perform such a task. In addition to
4 providing accurate small and medium precipitation, our procedure possesses built-in features
5 that allow us to model adequately extreme precipitation distributions.

6 First, we propose a new distribution for local precipitation via a probability mixture
7 model of Gamma and Generalized Pareto (GP) distributions. The latter one stems from Ex-
8 treme Value Theory (EVT). The performance of this mixture is tested on real and simulated
9 data, and also compared to classical rainfall densities.

10 Then, our downscaling method, extending the recently developed nonhomogeneous stochas-
11 tic weather typing approach, is presented. It can be summarized as a three-step program.
12 First, regional weather *precipitation patterns* are constructed through a hierarchical ascend-
13 ing clustering method. Second, daily transitions among our precipitation patterns are rep-
14 resented by a nonhomogeneous Markov model influenced by large-scale atmospheric vari-
15 ables like NCEP reanalyses. Third, conditionally on these regional patterns, precipitation
16 occurrence and intensity distributions are modeled as statistical mixtures. Precipitation
17 amplitudes are assumed to follow our mixture of Gamma and GP densities.

18 The proposed downscaling approach is applied to 37 weather stations in Illinois (USA) and
19 compared to various possible parameterizations and to a direct modeling. Model selection
20 procedures show that choosing one GP distribution shape parameter per pattern for all
21 stations provides the best rainfall representation amongst all tested models. This work

- 22 highlights the importance of EVT distributions to improve the modeling and downscaling of
- 23 local extreme precipitations.

1 Introduction

In recent decades, the accuracy of general circulation models (GCM) to simulate the large-scale behavior of the atmosphere has greatly improved. Still, such models have difficulties capturing small-scale intermittent processes, e.g. local precipitation. To better understand and represent these sub-grid scale meteorological characteristics, Regional Climate Models (RCM) offer an elegant way to integrate local processes through physical and dynamical equations. However, they can be extremely computer-intensive and their spatial resolution - generally from 5 to 50 km - does not always provide the required information needed in impact studies. Again, local precipitation can be considered as the archetypical example of such limitations. While advances in computer sciences may give the necessary computer power to resolve these smaller scales in the future, practitioners (flood planners, insurance companies, etc) need to make decisions locally with the current information today.

In order to link our large scale knowledge supplied by today's GCM, RCM and reanalysis outputs with measurements recorded at weather stations, statistical downscaling techniques offer a computationally attractive and ready-to-use route. This statistical approach consists of inferring significant relationships among large, regional and local scale variables. How to estimate, apply and test such relationships in order to have accurate representations of local features constitutes the so-called group of statistical downscaling questions. Three categories of methods are usually given to answer such questions: transfer functions, stochastic weather generators and weather typing methods. The first category is a direct approach. The relationships between large-scale variables and location-specific values are directly estimated via

45 either parametric, nonparametric, linear or nonlinear methods such as the analog method
46 [e.g. Barnett and Preisendorfer, 1978; Zorita and von Storch, 1998], multiple linear regressions
47 [e.g. Wigley *et al.*, 1990; Huth, 2002], kriging [e.g. Biau *et al.*, 1999] and neural networks
48 [e.g. Snell *et al.*, 2000; Cannon and Whitfield, 2002]. The second category focuses on weather
49 generators in which GCM outputs drive stochastic models of precipitation [e.g. Wilks, 1999;
50 Wilks and Wilby, 1999]. They are particularly of interest to assess local climate change [e.g.
51 Semenov and Barrow, 1997; Semenov *et al.*, 1998]. The weather typing approach, the third
52 and last category, encapsulates a wide range of methods that have in common an algorithmic
53 step in which recurrent *large-scale* and/or *regional atmospheric patterns* are identified. These
54 patterns are usually obtained from clustering and classification algorithms applied to geopo-
55 tential height, pressure or other meaningful atmospheric variables over a large spatial area.
56 These clustering and classification algorithms can be of different types: CART [Classification
57 and Regression Trees, see Breiman *et al.*, 1984; Schnur and Lettenmaier, 1998], “K-means”
58 methods [e.g. Huth, 2001; Yiou and Nogaj, 2004], hierarchical clustering approaches [e.g.
59 Davis *et al.*, 1993; Bunkers *et al.*, 1996], fuzzy-rules-based procedures [e.g. Pongracz *et al.*,
60 2001], neural networks [e.g. Bardossy *et al.*, 1994] or mixture of copula functions [Vrac *et*
61 *al.*, 2005]. Introducing such an intermediate layer (the weather patterns) in a downscaling
62 procedure provides a strong modeling flexibility. For example, linking directly the relation-
63 ships between large-scale atmospheric variables and precipitation recorded at a few weather
64 stations may be too complex in most inhabited regions. In comparison, it may be easier
65 and more efficient to first model the dependences between large-scale data and weather pat-
66 terns, the latter representing the recurrent atmospheric structures corresponding to a kind

67 of summary of the large scale. Then, we can focus on the coupling between weather patterns
68 and local measurements. Obviously, such a strategy will only be successful if the weather
69 patterns are carefully chosen; i.e., if they capture relevant recurrent summary information.
70 From a probabilistic point of view, the coupling step of a weather typing approach can be
71 viewed as deriving the following conditional probability density function (pdf)

$$f_{\mathbf{R}_t|\mathbf{S}_t} \tag{1}$$

72 which corresponds to the probability of observing local rainfall intensities, say \mathbf{R}_t , given the
73 current weather state, say \mathbf{S}_t , at time t . In addition to providing a simple mathematical
74 framework that can easily integrate various uncertainties, this probabilistic definition of
75 statistical downscaling is wide enough to cover many case studies. In this work, to get more
76 realistic precipitation variability than with a model only conditional on weather patterns,
77 the pdf (1) is also defined conditionally on a vector of large-scale atmospheric variables, say
78 \mathbf{X}_t , at time t :

$$f_{\mathbf{R}_t|\mathbf{X}_t,\mathbf{S}_t}. \tag{2}$$

79 In this paper, our main application is to downscale precipitation over the region of Illinois
80 (USA). Consequently, we would like to address the following questions: how to find adequate
81 regional weather patterns for \mathbf{S}_t ? How to model the coupling between large atmospheric
82 variables \mathbf{X}_t and \mathbf{S}_t ? What is an appropriate form for the conditional density defined by
83 (2)? The last question is the central one for the practitioner.

84 To our knowledge, none of the statistical downscaling methods discussed previously in this
85 section has been developed to address the issue of modeling both common and extreme values.

86 Nevertheless, although, for example, hydrologists and flood planners are interested in mean
87 precipitation, they also have a particular interest in modeling extreme local precipitation
88 because of its human, economical and hydrological impacts where large scale information
89 may help at modeling such extreme events. Past studies [Katz et al. 2002, Naveau et al.
90 2005] have illustrated how Extreme Value Theory (EVT), a statistical theory developed over
91 the past 80 years, provides the mathematical foundation for appropriately modeling extreme
92 precipitation. Hence, another important objective in this paper is to integrate EVT models
93 within a weather typing approach, i.e., throughout the density (2). To perform such a task,
94 we extend the original work on the nonhomogeneous stochastic weather typing approach by
95 Vrac *et al.* [2006].

96 The paper is organized as follows. In the first part of Section 2, we recall three clas-
97 sical distribution candidates that have been proposed to fit rainfall and we also introduce
98 a mixture model inspired by Frigessi *et al.* [2003]. A comparison and a discussion about
99 the performance of these four distributions is undertaken. In Section 3 the full data sets
100 are presented. Regional precipitation-related patterns are obtained by applying a hierarchi-
101 cal ascending clustering (HAC) algorithm to observed precipitation. Then, our statistical
102 downscaling model is explained. Section 4 contains results about our application and many
103 different diagnostics are computed to assess the quality of the models and to select the
104 most appropriate one. All along this section, instead of “pure” GCM outputs as large-
105 scale atmospheric variables, we take advantage of reanalysis data from the National Centers
106 for Environmental Prediction (NCEP). Indeed, not only are NCEP reanalyses constrained
107 GCM outputs, but also, using NCEP is necessary to assess our daily downscaling method in

108 a present climate, before fitting the method to (pure) GCM outputs to project local change
109 in precipitation. Hence, because the motivation is driven by the scale transformation of
110 large-scale atmospheric variables (GCM outputs or reanalysis data), working on reanalyses
111 is a first essential step. Lastly, in Section 5, we conclude and give some future research
112 directions.

113 2 Modeling rainfall locally

114 There exists a wide range of distribution families to statistically model rainfall intensities.
115 For example, [Katz, 1977; Wilks, 1999; Bellone *et al.*, 2000; Vrac *et al.*, 2006; Wilks, 2006]
116 argued that most of the precipitation variability can be approximated by a Gamma distribu-
117 tion. However, it is also well known [e.g. Katz *et al.*, 2002] that the tail of this distribution
118 can be too light to capture heavy rainfall intensities. This leads to the underestimation of
119 return levels and other quantities linked to high percentiles of precipitation amounts. Con-
120 sequently, the societal and economical impacts associated with heavy rains (e.g., floods) can
121 be miscalculated. To solve this issue, an increasingly popular approach in hydrology [Katz
122 *et al.*, 2002] is to disregard small precipitation values and to focus only on the largest rain-
123 fall amounts. The advantage of this strategy is that an elegant mathematical framework
124 called *Extreme Value theory* (EVT) developed in 1928 [Fisher and Tippett, 1928] and reg-
125 ularly updated during the last decades [e.g., Coles, 2001] dictates the distribution of heavy
126 precipitation. More specifically, EVT states that rainfall exceedances, i.e. amounts of rain
127 greater than a given threshold u , can be approximated by a Generalized Pareto Distribution

128 (GPD) if the threshold and the number of observations are large enough. In other words,
 129 the probability that the rainfall amount, say R , is greater than r given that $R > u$ is given
 130 by

$$P(R > r | R > u) = \left(1 + \xi \frac{r - u}{\sigma}\right)_+^{-1/\xi}, \quad (3)$$

131 where $a_+ = \max(a, 0)$ and $\sigma > 0$ represents the scale parameter. The shape parameter ξ
 132 describes the GPD tail behavior. If ξ is negative, the upper tail is bounded. If ξ is zero, this
 133 corresponds to the case of an exponential distribution (all moments are finite). If ξ is positive,
 134 the upper tail is still unbounded but higher moments eventually become infinite. These three
 135 cases are termed “bounded”, “light-tailed”, and “heavy-tailed”, respectively. The flexibility
 136 of the GPD to describe three different types of tail behavior makes it a universal tool for
 137 modeling exceedances. Although this GPD approach has been very successful to model heavy
 138 rains, it has the important drawback of overlooking small precipitation. Recently, Wilson
 139 and Toumi [2005] proposed a new probability distribution for heavy rainfall by invoking a
 140 simplified water balance equation. They claimed that the stretched exponential distribution
 141 tail defined by

$$P(R > r) = \exp \left[- \left(\frac{r}{\psi} \right)^\nu \right], \quad (4)$$

142 where $\psi > 0$ and $\nu > 0$ correspond to the scale and shape parameter. The latter should
 143 be equal to $\nu = 2/3$. This was justified by physical arguments that take into account of
 144 the distributions probabilities of quantities like the upward wind velocity \mathbf{w} (although the
 145 distribution of \mathbf{w} is much more unknown than the distribution of R). Note also that, although
 146 the parameter ν is expected to be equal to $2/3$ in theory, Wilson and Toumi did not say

147 that in practice this parameter has to be equal to $2/3$. Indeed, they estimated the shape
 148 parameter from different weather station precipitation measurements over the world. They
 149 found that, in practical applications, the estimated shape parameter is usually different from
 150 the $2/3$ constant. Despite its drawbacks, such a type of model is promising because it tries
 151 to combine probabilistic reasoning with physical arguments. But still, it is not designed
 152 to model small precipitation amounts. For their main example, Wilson and Toumi [2005]
 153 estimated the parameter (ψ, ν) in (4) for “heavy precipitation defined as daily totals with
 154 probability less than 5%”. Hence, one may wonder how to deal with the remaining 95%
 155 and what is the justification for working with 5% of the data and not 10%, 3% or any small
 156 percentages (this later problem also exists with a classical EVT approach). Because our
 157 final objective is to downscale the *full* range of precipitation values and because we do not
 158 want to choose an arbitrarily preset threshold (or percentage), we follow a different direction
 159 and opt for the method proposed by Frigessi *et al.* [2003]. These authors introduced the
 160 following mixture model

$$h_{\boldsymbol{\beta}}(r) = c(\boldsymbol{\beta}) \times [(1 - w_{m,\tau}(r)) \times f_{\beta_0}(r) + w_{m,\tau}(r) \times g_{\xi,\sigma}(r)] \quad (5)$$

161 where $c(\boldsymbol{\beta})$ is a normalizing constant, $\boldsymbol{\beta} = (m, \tau, \beta_0, \xi, \sigma)$ encapsulates the vector of unknown
 162 parameters, f_{β_0} corresponds to a light-tailed density with parameters β_0 , the function $g_{\xi,\sigma}$
 163 represents the GPD density that can be obtained from deriving the tail defined by (3) and
 164 $w_{m,\tau}(\cdot)$ is a weight function that depends on two parameters

$$w_{m,\tau}(r) = \frac{1}{2} + \frac{1}{\pi} \arctan\left(\frac{r - m}{\tau}\right). \quad (6)$$

165 Note that this weight function is non-decreasing, takes values in $(0, 1]$ and tends to 1 as r
 166 goes to ∞ ; i.e., heavy rains are represented by the GPD density $g_{\xi,\sigma}(r)$ in the mixture $h_{\beta}(r)$
 167 for large r . Conversely, small precipitation values are mostly captured by the light-tailed
 168 density $f_{\beta_0}(r)$. Hence, the idea behind equations (5) and (6) is rather simple: the mixing
 169 function $w_{m,\tau}(r)$ provides a smooth transition from a light-tailed density (small and medium
 170 precipitation) to the GPD density (heavy rainfalls). The parameters m and τ in $w_{m,\tau}(r)$
 171 correspond to the location and the speed of the transition from f_{β_0} to $g_{\xi,\sigma}$ in (5), respectively.
 172 In 2003, Frigessi *et al.* applied their model to Danish fire loss data and opted for a Weibull
 173 distribution as a light-tailed density in (5). In the context of precipitation modeling, past
 174 works [Bellone *et al.*, 2000; Vrac *et al.*, 2006; Wilks, 2006] indicate that a Gamma density,
 175 i.e.

$$f_{\beta_0}(x) = \frac{1}{\lambda^{\gamma}\Gamma(\gamma)}x^{\gamma-1}\exp(-x/\lambda), \text{ with } \beta_0 = (\gamma, \lambda), \quad (7)$$

176 should fit appropriately the bulk of the precipitation values (heavy rains excluded). This
 177 hypothesis could be challenged if the variable of interest was different, e.g. temperature. In
 178 addition, one may be puzzled by the “absence” of a threshold in Equation (5). Indeed, the
 179 threshold u in Equation (3) is forced to be equal to zero in (5). But introducing the weight
 180 function $w_{m,\tau}(r)$ and fixing the GPD threshold to zero brings two important benefits. First,
 181 the difficult threshold selection problem is replaced by a simpler unsupervised estimation
 182 procedure, i.e. finding m and τ in $w_{m,\tau}(r)$ from the data. This strategy is particularly
 183 relevant to large data sets analysis because it would be very time-consuming to find an
 184 adequate threshold for a large number of weather stations. Second, allowing for non-zero
 185 thresholds in (5) would impose an unwelcome discontinuity in $h_{\beta}(r)$. From a physical point

186 of view, such a discontinuity represents an unrealistic feature in precipitation.

187 In summary, we have four candidates for modeling local rainfall distribution:

- 188 • the Gamma density that works well for the main rainfall range but not for large values,
- 189 • the recently introduced stretched-exponential distribution function defined by (4), con-
190 structed on a physical foundation but only designed for heavy rainfall and not for small
191 precipitation values,
- 192 • the GPD function that works for extreme precipitation but not for small values, that is
193 mathematically sound and universal, in the sense that it can also fit temperature, winds
194 extremes, etc,
- 195 • and our new mixture model defined by (5) and (7) that combines the advantages of
196 the Gamma and GPD densities, and consequently can fit small and heavy rainfall.

197 To compare the performances of these four distributions, we implement the following pro-
198 cedure. We simulate 100 samples of 1000 iid realizations of each density with: $\lambda = 1$ and
199 $\gamma = 0.25$ for the Gamma distribution (see Equation (7)), $u = 0$, $\xi = 0.3$, $\sigma = 0.1$ for the GPD
200 (see Equation (3)), $u = 0$, $m = 1$, $\tau = 0.1$ for the mixture of the two previous distributions
201 (see equation (5)), and $\nu = 2/3$ and $\psi = 1$ for the stretched exponential (see Equation (4)),
202 respectively. Such parameter values were chosen because they correspond to reasonable esti-
203 mates for precipitation data. In particular, $\nu = 2/3$ is recommended by Wilson and Toumi.
204 As a second step, we fit each distribution to each of the four simulated samples by using the
205 maximum likelihood approach to compute the “optimal” parameters for each distribution.

206 To be consistent with Wilson and Toumi’s paper, the parameter ν in (4) is not considered
207 as a constant, i.e. we assume that this shape parameter has to be estimated. This has also
208 the advantage that we don’t penalize the stretched exponential distribution with respect to
209 the other distributions we test and for which the shape parameter is also not fixed but esti-
210 mated. The last step is to compare the qualities of the fit with respect to the given density.
211 Classically, one can compute the Akaike information criterion [AIC, Akaike, 1974], defined
212 by $-2\log(L) + 2p$, and the Bayesian information criterion [BIC, Schwarz, 1978], defined by
213 $-2\log(L) + p\log(n)$, where L is the likelihood of the model fitted to the data, p is the number
214 of parameters, and n is the number of data. Minimizing AIC and BIC helps to select the
215 model with a good fit to the data (i.e. high likelihood) while penalizing a model with too
216 many parameters. The BIC tends to add a larger parameter cost than the AIC. For our
217 simulations, the frequencies of selection of the four candidate distributions by the AIC and
218 BIC values are summarized in Table 1. As expected, the best AIC and BIC (in bold) are
219 majoritarily obtained along the diagonal of the table, i.e. the simulated samples are best
220 fitted by the density from which they were generated. We can remark that about one time
221 every third, the BIC indicates a Gamma fit when the true density is a mixture, i.e. the BIC
222 penalizes too much. In comparison, the AIC largely selects the correct distribution for all
223 four cases.

224 Hence, for these simulations, the AIC appears to perform reasonably and will be used in
225 the subsequent analyses. Still, we can not solely rely on these two criteria to discriminate
226 among models. In particular, these criteria may not be well adapted for extreme values.
227 Concerning the fit quality of the largest values, Figure 1 displays four quantile-quantile type

228 plots (QQplots). The \circ , \times , $+$, and \diamond signs correspond to the analytically fitted Gamma,
229 mixture, GP and stretched exponential densities, respectively. The $y = x$ black line rep-
230 resents the “true” distribution that can either be a Gamma (left-upper panel), a mixture
231 (right-upper panel), a GP (left-lower panel) and a stretched exponential (right-lower panel)
232 density. This graph mainly tells us that the mixture distribution (\times signs) appears to pro-
233 vide a very good fit in all cases. As expected, a Gamma fit (\circ signs) does not work very well
234 when the true trail is heavy. The stretched exponential (\diamond signs) is somehow limited because
235 it only provides a good fit when the true tail is stretched exponential. The worst case is
236 the GPD ($+$ signs), but this is expected because the threshold u was set to zero and it is
237 well known that the GPD only works well for very large values. An alternative would be to
238 select a high threshold, but then the main part of rainfall can not be statistically modeled
239 (and consequently, be compared with the other densities). Still, it is very interesting to see
240 that, despite of also having a GPD threshold set to zero, the mixture density provides very
241 good results. This reveals that the weight function $w_{m,\tau}$ in (6) can bring enough flexibility
242 even if the mixture threshold is equal to zero. One may argue that the mixture density has
243 too many parameters, but the AIC and BIC summarized in Table 1 do not show much cases
244 of over-fitting. Even more importantly, figure 1 shows that the other three classical distri-
245 butions for rainfall (Gamma, stretched exponential and GPD) do not offer the necessary
246 latitude to model the full spectrum of precipitation distribution.

247 Although the scope of this small simulation study is very limited and a more thorough
248 investigation would be welcome to review the arguments and problems related to local rainfall
249 distributions, Table 1 and Figure 1 strongly suggest that our mixture model could provide

250 a competitive probabilistic foundation. Consequently, this model will be used in the rest of
251 this paper. Concerning the choice between the AIC and BIC, only the AIC will be presented
252 in the remainder of this paper. In most cases, the BIC provides similar results and does not
253 change the meaning of the main findings that will be presented in Section 3.

254 With respect to real data, our goal is to analyze daily observations that were recorded
255 at 37 weather stations in Illinois (USA) from 1980 to 1999. Those stations correspond
256 to the complete dataset of precipitation provided for Illinois by the co-op observational
257 program. The stations are found to be uniformly distributed over Illinois. To reduce seasonal
258 influences, we only consider three winter months, December, January and February (DJF).
259 To illustrate the fit between our mixture model and real rainfall observations and also to
260 show the difference of fit to the data between a Gamma distribution and our mixture, we
261 select one station (Aledo) and apply a maximum likelihood estimation procedure to derive
262 the parameters of each distribution. Figure 2 shows the resulting quantile-quantile plots.
263 The upper panel displays the fit obtained using a Gamma distribution, while the lower panel
264 shows the result for our mixture distribution. As already seen in our simulation study, this
265 latter model provides a gain at capturing extreme values behavior. At this stage, one could
266 be satisfied by this type of station-per-station analysis. But from a statistical and physical
267 point of views, we prefer to go a step further in our statistical analysis by relating local
268 precipitation with large scale variables through an extension of our mixture model. This is
269 the object of the following section.

3 Our downscaling procedure

To develop a statistical model capable of downscaling precipitation, we need large-scale atmospheric variables and local observed precipitation measurements. The latter are provided here by daily observations described in Section 2. Large-scale atmospheric variables are given by NCEP reanalysis data - with a $2.5^\circ \times 2.5^\circ$ spatial resolution and at 850 mb. Three NCEP variables are considered in our analysis: geopotential height denoted Z_{850} , specific humidity, Q_{850} , and dew point temperature depression ΔT_{d850} defined as $T_{850} - T_{d850}$, where T_{850} and T_{d850} are the temperature and dew point temperature at 850 mb, respectively.

3.1 Modeling regional-scale precipitation patterns

Classically, weather typing methods are based on *circulation*-related patterns. A number of studies (e.g. Mamassis and Koutsoyiannis, 1996) showed that, according to the studied region, large-scale atmospheric patterns can be efficient to explain and characterize local precipitation variability. However, to better represent precipitation behaviors, we follow the approach of Vrac *et al.* [2006]. Instead of defining upper-air circulation patterns, these authors recently constructed *precipitation*-related patterns, directly obtained from a subset of observed local precipitations, and showed that, for Illinois, these patterns are more efficient than classical upper-air circulation patterns to characterize and simulate local precipitation. These precipitation patterns were derived from a hierarchical ascending clustering (HAC) algorithm with Ward criterion [Ward, 1963], applied to the observed precipitation of the 1980-1999 winter months (DJF). Instead of the common Euclidean distance, a special metric

290 tailored to precipitation was developed to take account of the spatio-temporal rain features.
 291 The details of this clustering algorithm can be found in Vrac *et al.* [2006]. Figure 4 shows the
 292 four precipitation patterns over the region of Illinois. It is clear that pattern 1 represents the
 293 smallest rainfall intensities whereas pattern 4 corresponds to the most intense precipitation.
 294 Patterns 2 and 3 show moderate precipitation, with opposite South/North and North/South
 295 gradients respectively. The North/South gradient (drier in the north and wetter in the south)
 296 that is also perceptible in pattern 4, is a classical recurrent feature of winter precipitation in
 297 Illinois.

298 **3.2 Relating regional precipitation patterns with large-scale NCEP** 299 **outputs**

300 At this stage, precipitation-related structures S_t have been derived (see Figure 4) and repre-
 301 sent the regional scale. How to link them to the larger scale (the NCEP reanalysis) and how
 302 to connect them to the smaller scale (the weather stations) are the two remaining questions
 303 we have to address in this paper. In this section, we focus on answering the first one. To
 304 perform this task, we model the day-to-day probability transitions from the given weather
 305 state at day t , say S_t , to the state of the following day, S_{t+1} as a function of the current
 306 large atmospheric variables, say \mathbf{X}_t , from the NCEP reanalysis. More precisely, a nonhomo-
 307 geneous Markov model [e.g., Bellone *et al.*, 2000] is fitted to our NCEP data and our states
 308 by applying the following temporal dependence structure

$$P(S_t = s | S_{t-1} = s', \mathbf{X}_t) \propto \gamma_{s's} \exp \left[-\frac{1}{2} (\mathbf{X}_t - \mu_{s's}) \Sigma^{-1} (\mathbf{X}_t - \mu_{s's})' \right]. \quad (8)$$

309 where the symbol \propto means “proportional to” and where $\gamma_{s's}$ is the baseline transition proba-
310 bility from pattern s' to pattern s , corresponding to the observed transition probability from
311 s' to s , i.e. the proportion of transitions from s' to s over the total number of transitions.
312 In the above formula, we can recognize a weight represented by the exponential term that
313 is proportional to a normal density whose mean $\mu_{s's}$ and variance matrix Σ are directly
314 representing the influence of the large atmospheric variable \mathbf{X}_t . Eq. (8) comes from Bayes’s
315 theorem, saying that:

$$\begin{aligned}
P(S_t = s | S_{t-1} = s', \mathbf{X}_t) &= \frac{P(S_t = s | S_{t-1} = s') P(\mathbf{X}_t | S_t = s, S_{t-1} = s')}{P(\mathbf{X}_t | S_{t-1} = s')} \\
&= \frac{\gamma_{s's} P(\mathbf{X}_t | S_t = s, S_{t-1} = s')}{\sum_k \gamma_{ks} P(\mathbf{X}_t | S_t = s, S_{t-1} = k)} \tag{9}
\end{aligned}$$

316 By assuming in Eq. (9) that \mathbf{X}_t is multivariate normal, Eq. (8) is easily derived. In Eq. (8),
317 $\mu_{s's}$ corresponds to the mean vector of the atmospheric variables at time t when transition-
318 ing from $S_{t-1} = s'$ to $S_t = s$. The four precipitation patterns defined in section 3.1 imply a
319 reasonable number of 16 possible transition. Hence the 16 $\mu_{s's}$ and $\gamma_{s's}$ to be computed can
320 be estimated very fast. As for Σ , it is the variance-covariance matrix for the whole dataset
321 of large-scale atmospheric data (centered around their mean). Indeed, as in Charles *et al.*
322 (1999), Bellone *et al.* (2000) or Vrac *et al.* (2006), for stability reasons, a single covariance
323 matrix is preferred over one matrix per transition. In contrast to the exponential part
324 of Eq. (8), the baseline transition probability $\gamma_{s's}$ in (8) is time invariant and corresponds
325 to the transition probabilities that one would have if large scale features did not bring any
326 information. This case corresponds to the homogeneous Markov model. Hence, allowing
327 a non-homogeneity in our Markov modeling brings the necessary flexibility to mathemati-

328 cally integrate large-scale information at the intermediate level of the regional precipitation
 329 patterns.

330 **3.3 Linking regional precipitation patterns to local precipitation**

331 In order to implement an efficient downscaling precipitation scheme, we also need to model
 332 accurately the distributional properties of precipitation at the smallest scale, i.e. the ones
 333 recorded at rain gauges.

334 We now assume that, given the current weather state s , all the rainfall intensities for
 335 station i follow the density $h_{\beta_{si}}$ given in (5) with state- and site-specific parameters. This
 336 gives us the last ingredient to determine our main density defined by (2): the probability of
 337 observing local rainfall intensities at day t , say $\mathbf{R}_t = (R_{t,1}, \dots, R_{t,N})$, given the current weather
 338 state, say $S_t = s$, and large-scale atmospheric variables, say \mathbf{X}_t . To compute $f_{\mathbf{R}_t|\mathbf{X}_t, S_t}$, we
 339 follow Bellone *et al.* [2000] who considered that each rain gauge is spatially independent
 340 given the state S_t . Mathematically, this assumption translates into the following equality

$$f_{\mathbf{R}_t|\mathbf{X}_t, S_t}(r_{t1}, \dots, r_{tN}) = \prod_{i=1}^N f_{R_{ti}|\mathbf{X}_t, S_t}(r_{ti}) \quad (10)$$

341 To give an explicit form for the density $f_{R_{ti}|\mathbf{X}_t, S_t}$, we take advantage of Vrac *et al.* [2006]
 342 who suggested the following form

$$f_{R_{ti}|\mathbf{X}_t, S_t=s}(r_{ti}) = [p(\mathbf{X}_t; \boldsymbol{\alpha}_{si})h_{\beta_{si}}(r_{ti})]^{\mathbb{1}_{\{r_{ti}>0\}}} \times [1 - p(\mathbf{X}_t; \boldsymbol{\alpha}_{si})]^{\mathbb{1}_{\{r_{ti}=0\}}} \quad (11)$$

343 where $h_{\beta_{si}}$ is given by (5), $\mathbb{1}_{\{a\}} = 1$ if a is true and 0 if false, and $p(\mathbf{X}_t; \boldsymbol{\alpha}_{si})$ represents
 344 the probability of rain occurrence for weather station i in state s . Equation (11) may look
 345 complex at first sight. Basically, it is composed of three elements:

346 (a) the indicator function $\mathbb{1}_{\{r_{ti}=0\}}$ is necessary to take into account that the rain gauge i
347 can record no precipitation during day t ,

348 (b) $1 - p(\mathbf{X}_t; \boldsymbol{\alpha}_{si})$ provides the probability of such a dry day and it depends on the atmo-
349 spheric variables \mathbf{X}_t through a logistic regression with parameters $\boldsymbol{\alpha}_{si}$, as suggested
350 by Jeffries and Pfeiffer [2000]:

$$p(\mathbf{X}_t; \boldsymbol{\alpha}_{si}) = P(R_{ti} > 0 | S_t = s, \mathbf{X}_t) = \frac{\exp(\mathbf{X}_t' \boldsymbol{\alpha}_{si})}{1 + \exp(\mathbf{X}_t' \boldsymbol{\alpha}_{si})} \quad (12)$$

351 (c) the density $h_{\beta_{si}}(r_{ti})$ corresponds to positive rainfall values.

352 Combining equations (8), (5), (10) and (11) constitutes the main components of our stochas-
353 tic weather typing approach. It integrates three scales (small, regional and large) through the
354 variables \mathbf{R}_t , S_t and \mathbf{X}_t . In addition, the full spectrum of precipitation values (dry events,
355 medium precipitation, heavy rainfall) is modeled.

356 4 A case study: Precipitation in Illinois, USA

357 As previously mentioned, Figure 4 displays our four selected regional precipitation patterns
358 over the region of Illinois. From these four patterns, the nonhomogeneous Markov model
359 is parameterized, and the parameters of the conditional distributions of precipitation are
360 estimated by Maximum Likelihood Estimation (MLE), given each observed (i.e. pre-defined)
361 pattern. In the following simulation process, the precipitation patterns are stochastically
362 simulated, for each t , according to the parameterized NMM, influenced by the large-scale
363 atmospheric variables. In other words, in the simulation step, we do not use the patterns

364 defined previously by HAC but we generate new ones according to \mathbf{X}_t and our model.
365 Conditionally on the four patterns, equations (10) and (11) offer a wide range of modeling
366 possibilities. For example, one may wonder if it is better to have a unique GPD shape
367 parameter ξ for all precipitation patterns and at all rain gauges or if a better statistical
368 fit can be obtained by allowing this shape parameter to vary from station to station, while
369 taking into account the risk of over-parametrization. Before presenting the seven different
370 models that we have tested and compared, we note that the parameter τ in Eq. (6) cannot
371 be null. For this reason, from the limit of Eq. (6) when τ goes to 0, we extend Eq. (6) to

$$w_{m,0}(r) = \begin{cases} 0, & \text{if } r < m \\ 0.5, & \text{if } r = m \\ 1, & \text{if } r > m \end{cases} \quad (13)$$

372 for $\tau = 0$, whenever we do not wish to estimate τ and we think that the transition from
373 the Gamma to the GPD distribution is very fast in the mixture defined by (10). Our seven
374 models are the following ones:

375 (0) Gamma and GPD mixtures whose parameters vary with location and precipitation
376 pattern,

377 (i) only Gamma distributions (no GPD in the model) whose parameters vary with location
378 and precipitation pattern,

379 (ii) Gamma and GPD mixtures with one ξ parameter per pattern (i.e. given the weather
380 pattern, the weather stations have the same ξ),

381 (iii) same as (ii) with τ set to be equal to 0,

382 (iv) Gamma and GPD mixtures with one common ξ for all stations and all patterns,

383 (v) same as (iv) with τ set to be equal to 0.

384 (iii)* same as model (iii) - one ξ parameter per pattern with $\tau = 0$ - except that only Gamma
385 distributions are used in pattern 1. Indeed, since this pattern corresponds to small or
386 null intensities of rainfall, a modelling of the extreme events could have no sense here.

387 From a statistical point of view, the GPD shape parameters are very difficult to estimate
388 (wide confidence intervals). Hence, diminishing the number of ξ parameters to estimate like
389 in model (iii) reduces the overall variability. In addition, interpreting four ξ parameters (one
390 per pattern, see models (ii) & (iii)) instead of 37×4 is much easier for the hydrologist.
391 Besides these two general guidelines, we need a more objective “measure” to compare our
392 seven models. As in Section 2, we opt for minimizing the classical AIC criterion (similar
393 results are obtained with the BIC).

394 Our seven models’ differences primarily focus on the degree of flexibility allowed for ξ
395 and τ . Concerning the other parameters (σ, m, \dots), we allow them to vary across stations
396 and across patterns because they mainly represent local variability.

397 For each model, we estimate its parameters by implementing a maximum likelihood
398 estimation method. To illustrate the quality and drawbacks of our approach, we will comment
399 on five example stations in this section: Aledo (North-West of Illinois), Aurora (North-East),
400 Fairfield (South-East), Sparta (South-West), and Windsor (center-East of Illinois). This
401 subset was picked because we believe that it represents a large range of cases and space
402 limitations make it impossible to provide plots and tables for all 37 stations.

403 Concerning the large-scale atmospheric variables \mathbf{X}_t , we assume that only the NCEP
404 grid-cells over Illinois have the potential to influence local precipitation and transition prob-
405 abilities. Consequently, we only work with the six grid-cells that cover Illinois. According to
406 the studied region, it is possible that taking more NCEP grid-cells into account could improve
407 the modeling and the simulation process. A few attempts have been made to enlarge the
408 NCEP area influencing local precipitation and patterns transitions. The associated results,
409 not presented here, did not show any clear improvement for the Illinois region, compared to
410 the results obtained from the six grid-cells. Moreover, the more grid-cells we work on, the
411 more parameters we have (with a risk of over-parameterization). Hence from a computa-
412 tional point of view, it is better to restrict the large-scale influence to a reasonable number
413 of NCEP grid-cells over Illinois. Based on these two considerations, we then limit the appli-
414 cation presented here to the six NCEP grid-cells over Illinois to influence local precipitation
415 and patterns transitions.

416 Instead of working directly with the *raw* variables, Z_{850} , Q_{850} , and ΔT_{a850} - corresponding
417 to $6 \times 3 = 18$ variables - we perform a Singular Value Decomposition [Von Storch and Zwier,
418 1999; Vrac *et al.*, 2006; Wilks, 2006]. This has the advantage of reducing significantly the
419 dimensionality of the NCEP data, while keeping the main part of information brought by the
420 reanalysis. The SVD operation gives us the following summary: the SVD explains 93.6%,
421 98.6%, and 97.5% of the correlation for Z_{850} , Q_{850} , and ΔT_{a850} respectively.

422 A central theme in this paper is how to capture the full range of precipitation, extremes
423 included. To determine if the addition of a GPD to a Gamma density is worthwhile, Figure
424 5 displays QQplots (empirical quantiles versus modeled quantiles) for the Sparta station for

425 two precipitation patterns (see the left and right panels) and in two models: (0) & (i), see
426 the lower and upper panels, respectively. In contrast to histograms, the QQplots are, by
427 design, capable of representing the quality of the estimated fit at the end of the distribution
428 tail, i.e. they can show the capacity of our mixture model to represent extreme precipitation.

429 Figure 5 indicates that a fitted Gamma has the tendency to either underestimate (5.a)
430 or overestimate (5.b) the largest precipitation for this station, respectively to the precipi-
431 tation patterns. Fig. 5.a and 5.c show that, for pattern 2, our mixture can model heavier
432 rainfall than the gamma distribution alone (i.e. characterizes stronger intensities for this
433 pattern/station). To explain how the Gamma model can overestimate large precipitation in
434 Fig. 5.b, we have to keep in mind that the whole rainfall range is fitted and the Gamma
435 distribution does not have a shape parameter for the tail of the distribution. In the presence
436 of a heavy tail, it is not clear how the estimation procedure is going to compensate the facts
437 that the gamma distribution is not heavy tailed and that the whole distribution has to be fit-
438 ted. Either the Gamma scale parameter can be largely overestimated (by the largest values)
439 or underestimated (depending on the spread and the size of the sample). Applying a robust
440 estimator to find the Gamma scale parameter should remove the problem of overestimation,
441 but then heavy tailed values will even be more disregarded. Consequently, a possible solution
442 is to allow a distribution (like the GPD) with a shape parameter. More generally, Fig. 5
443 clearly indicates that integrating a GPD improves the fit of “large” rainfalls for this station,
444 as the closer the estimated quantiles are to the empirical quantiles the better. Of course, this
445 does not mean that this is true for all stations and all patterns. Instead, this shows that our
446 mixture defined by (5) provides the necessary modeling flexibility to describe heavy-tailed

447 behaviors when needed. If no heavy rainfalls are observed at a given station, the estimated
448 weight defined by (6) should take small values to favor the Gamma distribution, i.e., m large
449 for this station.

450 Concerning the model selection, Table 2 compares models (0) and (i) with respect to the
451 Akaike Information Criterion (AIC) for our five selected stations and for each precipitation
452 pattern. Because the BIC values gave us equivalent results, they are not provided in this
453 table, illustrating that the optimal choice between model (0) and model (i) varies greatly
454 across stations and across patterns. For example, introducing a GPD seems to be a good
455 choice for Sparta, while a simpler Gamma model appears to be sufficient for Aurora.

456 Table 3 contain the AIC values obtained for the seven models. The bold values correspond
457 to the optimal criterion of each row. Taking model (iii)* ($\tau = 0$, a Gamma distribution for
458 pattern 1 and one ξ parameter per pattern for patterns 2-4) provides the best AIC for Sparta,
459 while setting one overall ξ parameter gives the best AIC for the four other stations. For any
460 of the five stations, we can remark that setting $\tau = 0$ in model (ii) - i.e. going from model (ii)
461 to model (iii) - brings an improvement of the AIC. This means that restricting the number
462 of ξ parameters generally provides better criteria. Models (iii)* and (iv) seem to be the most
463 competitive ones in general (i.e. for most of the stations separately), while the preferred
464 model tends to be (iii)* for the set of the five selected weather stations altogether (last row
465 of Table 3). Consequently, model (iii)*, i.e. pattern 1 associated to Gamma distributions
466 and patterns 2-4 to mixtures with one ξ parameter per pattern with the constant $\tau = 0$, is
467 chosen as the most efficient model, as it provides the best overall criterion for the set of these
468 five stations. Hence, this model can well represent both common and extreme precipitation

469 values with an acceptable number of parameters and has the overall preference.

470 Table 4 shows the values of the ξ parameters and the values of the m parameters (when
471 applicable) for the five example stations for model (iii)*. The three ξ parameters are clearly
472 positive. These positive values indicate that the heavy tail component in our mixture pdf
473 is essential to model heavy rainfalls for precipitation patterns 2 to 4, while the Gamma
474 distributions (with light tails) are sufficient in pattern 1 corresponding to small precipitation
475 events. Unsurprisingly, the m parameters tend to increase from pattern 2 (with the smallest
476 rainfall intensities among patterns 2-4) to pattern 4 (with the strongest rainfalls among all
477 patterns).

478 To visually evaluate the fit between our model (iii)* and the observed precipitation, a
479 QQplot is plotted for the Aledo station in Fig. 6. The agreement between observed and
480 theoretical quantiles (even for high quantiles) is clearly good. Fig. 6 has to be compared to
481 Fig. 2. This allows us to conclude that, not only the AIC is better for model(iii)* than for a
482 “no pattern” modeling, but also that model(iii)* improves the QQplot.

483 Besides heavy rainfalls, an important characteristic of precipitation modeling is the rep-
484 resentation of the so-called wet and dry spell periods, fundamental quantities in agriculture.
485 Note that none of the following results concerning wet and dry spells and local precipitation
486 probabilities, presented and shown from Fig. 7, depends on the Gamma or mixture models.
487 Indeed, they are only related to the nonhomogeneity introduced in the Markov model (8) -
488 that characterizes pattern transitions - and to the probabilities of local rain occurrence mod-
489 eled as logistic regressions (see (11) and (12)). So, the following results are directly derived
490 from the model developed by Vrac *et al.* [2006] and allow us to compare some precipita-

491 tion appearance characteristics obtained from the “four precipitation patterns” and those
492 obtained from the alternative “no pattern” approach.

493 In this context, we have noticed that the four precipitation patterns have to be included
494 in order to obtain adequate wet and dry spell probabilities. For example, Fig. 7 shows such
495 probabilities (in log-scale) at two stations, respectively Fairfield and Windsor. Upper panels
496 (a) and (b) display these probabilities when the four precipitation patterns are included in
497 our analysis. In contrast, lower panels (c) and (d) show the results when no patterns are
498 introduced. From these graphs, one can see that the “no pattern” option is not completely
499 satisfying, it tends to underestimate the probabilities for long spells, above all for dry spells.

500 **5 Conclusion**

501 We presented here a nonhomogenous stochastic weather typing method to downscale the full
502 spectrum of precipitation distributional behaviors. Our downscaling technique is based on a
503 nonhomogeneous Markov model that characterizes the transitions amongst different precip-
504 itation patterns obtained from a hierarchical ascending clustering algorithm. Conditionally
505 on these precipitation patterns, the precipitation distribution is modeled by a mixture model
506 that integrates heavy rainfalls, medium precipitation and no rain occurrences, and that de-
507 pends on large-scale features given from a SVD applied to NCEP reanalysis.

508 After applying our approach to the region of Illinois, it appears that a specific subclass of
509 our model (the one with Gamma distributions for pattern 1 and mixture models with a single
510 GPD shape parameter per pattern for patterns 2-4) produces the best fit with respect to

557 A.J. Cannon and P.H. Whitfield. Downscaling recent streamflow conditions in british
558 columbia, canada using ensemble neural network models. *Journal of Hydrology*, 259:136–
559 151, March 2002.

560 S.P. Charles, B.C. Bates, and J.P. Hughes. A spatio-temporal model for downscaling pre-
561 cipitation occurrence and amounts. *Journal of Geophysical Research*, 104:31,657 – 31,669,
562 1999.

563 S. Coles. *An Introduction to Statistical Modeling of Extreme Values*. Springer-Verlag,
564 London, 2001.

565 R.E. Davis, R. Dolan, and G. Demme. Synoptic climatology of atlantic coast northeasters.
566 *Intl. J. Climatol.*, 13:171–189, 1993.

567 R.A. Fisher and L.H.C. Tippett. Limiting forms of the frequency distribution of the largest
568 or smallest member of a sample. *Proceedings of the Cambridge Philosophical Society*, 24:180–
569 190, 1928.

570 A. Frigessi, O. Haug, and H. Rue. A dynamic mixture model for unsupervised tail estimation
571 without threshold selection. *Extremes*, 5:219–235, 2003.

572 R. Huth. Disaggregating climatic trends by classification of circulation patterns. *Interna-*
573 *tional Journal of Climatology*, 21:135–153, 2001.

574 R. Huth. Statistical downscaling of daily temperature in central europe. *Journal of Climate*,
575 15:1731–1742, 2002.

511 the AIC criterion for this region. In terms of extreme precipitation, this model corresponds
512 to a very fast transition from the Gamma distribution to the GPD for patterns 2-4. It is
513 also worthwhile to highlight that introducing four precipitation patterns produces better
514 precipitation characteristics than a direct “no pattern” approach does.

515 As possible improvements, spatial dependence modeling could be introduced in this model
516 to better represent the correlation between stations. In that context, Bayesian hierarchical
517 methods could provide an additional flexibility. A possible application of our downscaling
518 procedure could be the projection of future local precipitation based on large-scale climate
519 change simulated by GCMs. While the estimation step requires both present large- and
520 local-scale data, the local projection of future climate scenarios can be done by using only
521 the GCM outputs describing future time periods. Based on the NMM previously fitted, the
522 future large-scale outputs are first used to influence the simulation of future precipitation
523 patterns through Eq. (8). No local precipitation is needed for this step, since it is obviously
524 not even available. Conditionally on the generated future patterns, probabilities of local
525 rainfall events can be computed - influenced by the large-scale GCM outputs - through Eq.
526 (12) for rain appearances and through Eq. (11) for intensities. These local projections would
527 then allow economic impact studies of extreme precipitation.

528

529 **Acknowledgments**

530 Although this research has been funded in part by the United States Environmental
531 Protection Agency through STAR Cooperative Agreement # R-82940201 to the University
532 of Chicago, it has not been subjected to the Agency’s required peer and policy review and

533 therefore does not necessarily reflect the views of the Agency, and no official endorsement
534 should be inferred. P. Naveau's research work is supported by the european E2-C2 grant,
535 the National Science Foundation (grant: NSF-GMC (ATM-0327936)) and by The Weather
536 and Climate Impact Assessment Science Initiative at the National Center for Atmospheric
537 Research (NCAR).

References

- 538
- 539 H. Akaike. A new look at the statistical model identification. *IEEE Transactions on*
540 *Automatic Control*, 19:716–723, 1974.
- 541 A. Bárdossy, H. Muster, L. Duckstein, and I. Bogardi. Automatic classification of circu-
542 lation patterns for stochastic precipitation modelling. *Stochastic and Statistical Methods*
543 *in Hydrology and Environmental Emgineering*, 1. Extreme Values: Floods and Droughts,
544 1994.
- 545 T. Barnett and R. Preisendorfer. Multifield analog prediction of short-term climate fluctu-
546 ations using a climate state vector. *Journal of Atmospheric Science*, 35:1771–1787, 1978.
- 547 E. Bellone, J.P. Hughes, and J. P. Guttorp. A hidden markov model for downscaling
548 synoptic atmospheric patterns to precipitation amounts. *Climate Research*, 15:1–12, 2000.
- 549 G. Biau, E. Zorita, H. von Storch, and H. Wackernagel. Estimation of precipitation by
550 kriging in the eof space of thesea level pressure field. *Journal of Climate*, 12:1070–1085,
551 1999.
- 552 L. Breiman, J. Friedman R., Olshen, and C. Stone. *Classification And Regression Trees*
553 *(CART)*. Chapman & Hall, New York, London, 1984.
- 554 M. J. Bunkers, J. R. Miller, and A. T. DeGaetand. Definition of climate regions in the
555 northern plains using an objective cluster modification technique. *Journal of Climate*,
556 9:130–146, 1996.

576 N. Jeffries and R. Pfeiffer. A mixture model for the probability distribution of rain rate.
577 *Environmetrics*, 12:1–10, 2000.

578 R. Katz, M.B. Parlange, and P. Naveau. Statistics of extremes in hydrology. *Advances in*
579 *Water Resources*, 25:1287–1304, 2002.

580 R.W. Katz. Precipitation as a chain-dependent process. *Journal of Applied Meteorology*,
581 16:671–676, 1977.

582 N. Mamassis and D. Koutsoyiannis. Influence of atmospheric circulation types on space-
583 time distribution of intense rainfall. *J. Geophys. Res.*, 101 (D21):26,267–26,276, 1996.

584 P. Naveau, M. Nogaj, C. Ammann, P. Yiou, D. Cooley, and V. Jomelli. Statistical methods
585 for the analysis of climate extremes. *C.R. Geoscience (In press)*, 2005.

586 R. Pongracz, J. Bartholy, and I. Bogardi. Fuzzy rule-based prediction of monthly precipi-
587 tation. *Phys. Chem. Earth*, 9:663–667, 2001.

588 R. Schnur and D. Lettenmaier. A case study of statistical downscaling in australia using
589 weather classification by recursive partitioning. *Journal of Hydrology*, 212-213:362–379,
590 1998.

591 G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6:461–464,
592 1978.

593 M.A. Semenov and E.M. Barrow. Use of a stochastic weather generator in the development
594 of climate change scenarios. *Climate Research*, 35:397–414, 1997.

595 M.A. Semenov, R.J. Brooks, E.M. Barrow, and C.W. Richardson. Comparison of the
596 WGEN and the LARS-WG stochastic weather generators in diverse climates. *Climate*
597 *Research*, 10:95–107, 1998.

598 S.E. Snell, S. Gopal, and R.K. Kaufmann. Spatial interpolation of surface air temperatures
599 using artificial neural networks: Evaluating their use for downscaling gcms. *Journal of*
600 *Climate*, 13:886–895, 2000.

601 H. Von Storch and F.W. Zwiers. *Statistical Analysis in Climate Research*. Cambridge
602 University Press, Cambridge, 1999.

603 M. Vrac, A. Chédin, and E. Diday. Clustering a global field of atmospheric profiles by mix-
604 ture decomposition of copulas. *Journal of Atmospheric and Oceanic Technology*, 22:1445–
605 1459, 2005.

606 M. Vrac, M. Stein, and K. Hayhoe. Statistical downscaling of precipitation through a
607 nonhomogeneous stochastic weather typing approach. *Climate Research (In press)*, 2006.

608 J.H. Ward. Hierarchical grouping to optimize an objective function. *Journal of the American*
609 *Statistical Association*, 58:236–244, 1963.

610 T.M. Wigley, P. Jones, K. Briffa, and G. Smith. Obtaining subgrid scale information from
611 coarse resolution general circulation output. *J. Geophys. Res.*, 95:1943–1953, 1990.

612 D. Wilks. *Statistical methods in the atmospheric sciences (second edition)*. Elsevier, Oxford,
613 2006.

614 D.L. Wilks and R.L. Wilby. The weather generation game: a review of stochastic weather
615 models. *Progress in Physical Geography*, 23:329–357, 1999.

616 D.S. Wilks. Multisite downscaling of daily precipitation with a stochastic weather generator.
617 *Climate Research*, 11:125–136, 1999.

618 P.S. Wilson and R. Toumi. A fundamental probability distribution for heavy rainfall.
619 *Geophys. Res. Lett.*, 32:L14812, doi:10.1029/2005GL022465, 2005.

620 P. Yiou and N. Nogaj. Extreme climatic events and weather regimes over the north atlantic:
621 When and where? *Geophys. Res. Lett.*, 31, 2004.

622 E. Zorita and H. von Storch. The analog method as a simple statistical downscaling tech-
623 nique: Comparison with more complicated methods. *Journal of Climate*, 12:2474–2489,
624 1998.

List of Tables

625			
626	1	Frequencies of selections of the four candidate distributions by the Akaike	
627		Information Criterion (AIC) and Bayesian information criterion (BIC) values	
628		obtained from 100 samples of 1000 simulated data (for each given density).	
629		The bold fonts correspond to the highest frequencies with respect to the AIC	
630		and the BIC.	38
631	2	Akaike Information Criterion (AIC) values obtained pattern by pattern for five	
632		weather stations. The bold values correspond to the optimal criteria either	
633		for model (0) or (i)	39
634	3	Akaike Information Criterion (AIC) values obtained for our five selected weather	
635		stations and for our seven models. The bold values correspond to the optimal	
636		criterion per row. Below each model's name, the number p of parameters for	
637		n stations is provided.	40
638	4	Values of the ξ and m parameters for the five example stations for model (iii)*.	
639		Non-applicable (NA) is indicated for pattern 1, since this pattern is associated	
640		to Gamma distributions in this model.	41

List of Figures

641			
642	1	Quantiles-quantiles plots (i.e. theoretical vs. fitted quantiles). The o, ×, +,	
643		and ◊ signs correspond to the QQplots from the Gamma, mixture, GP and	
644		stretched exponential densities, respectively. Each distribution is analytically	
645		fitted by a Gamma (left-upper panel), our mixture (right-upper panel), a GP	
646		(left-lower panel) and a stretched exponential (right-lower panel) density. The	
647		99% quantile is indicated for each fitted distribution. These graphes mainly	
648		tell us that the mixture distribution (× signs) appears to provide a very good	
649		fit in all cases.	42
650	2	QQplot for Aledo with (a) Gamma distribution, and (b) our mixture.	43
651	3	Schematic graph explaining the main components of our downscaling scheme.	44
652	4	Four station-based precipitation patterns over Illinois derived by the Vrac <i>et</i>	
653		<i>al.</i> (2006) HAC method, with area proportional to mean rainfall for each cluster.	45
654	5	QQplots of precipitation patterns 2 and 3 for station “Sparta”, for function	
655		h_β in (11) as a Gamma distribution in (a) and (b) and h_β as a mixture (5) in	
656		(c) and (d). Units are cm.	46
657	6	QQplot for Aledo with four patterns and model (iii)*, i.e., Gamma distribu-	
658		tions for pattern 1 and mixtures for patterns 2-4 with one ξ per pattern and	
659		$\tau = 0$	47

660 7 Wet and dry spells probabilities (in log-scale) obtained for Fairfield and Wind-
661 sor. Upper panels (a) and (b): the “4 patterns” approach; lower panels (c)
662 and (d): the “no pattern” approach. 48

True density	Fitted density			
	Gamma	GP	Mixture	Stretched
Gamma	AIC= 90 BIC= 100	AIC= 0 BIC= 0	AIC= 10 BIC= 0	AIC= 0 BIC= 0
Generalized- Pareto	AIC= 0 BIC= 0	AIC= 86 BIC= 96	AIC= 11 BIC= 1	AIC= 3 BIC= 3
Mixture: GP + Gamma	AIC= 7 BIC= 36	AIC= 0 BIC= 0	AIC= 93 BIC= 64	AIC= 0 BIC= 0
Stretched exponential	AIC= 3 BIC= 3	AIC= 0 BIC= 0	AIC= 10 BIC= 0	AIC= 87 BIC= 97

Table 1: Frequencies of selections of the four candidate distributions by the Akaike Information Criterion (AIC) and Bayesian information criterion (BIC) values obtained from 100 samples of 1000 simulated data (for each given density). The bold fonts correspond to the highest frequencies with respect to the AIC and the BIC.

Station	Model	Pattern 1	Pattern 2	Pattern 3	Pattern 4
Aledo	(0)	-351.15	-486.98	-162.21	86.72
	(i)	-349.15	-493.05	-163.29	84.86
Aurora	(0)	-948.62	-663.43	-228.28	272.63
	(i)	-954.48	-670.41	-235.40	265.43
Fairfield	(0)	-367.72	-513.05	57.15	499.93
	(i)	-375.99	-282.21	97.42	741.63
Sparta	(0)	-131.34	-488.52	-128.03	613.23
	(i)	-129.61	-466.06	-123.57	766.54
Windsor	(0)	-632.25	-982.22	-321.01	441.08
	(i)	-613.92	-985.26	-325.78	579.47

Table 2: Akaike Information Criterion (AIC) values obtained pattern by pattern for five weather stations. The bold values correspond to the optimal criteria either for model (0) or (i)

Station	Model (0) $p = 24n$	Model (i) $p = 8n$	Model (ii) $p = 20n + 4$	Model (iii) $p = 16n + 4$	Model (iv) $p = 20n + 1$	Model (v) $p = 16n + 1$	Model (iii)* $p = 12n + 5$
Aledo	AIC=-796.52	AIC=-816.58	AIC=-795.76	AIC=-809.79	AIC=- 819.46	AIC=-816.18	AIC=-816.79
Aurora	AIC=-1137.47	AIC=-1149.99	AIC=-1256.53	AIC=-1293.89	AIC=- 1358.48	AIC=-1152.51	AIC=-1299.89
Fairfield	AIC=14.36	AIC=103.07	AIC=22.45	AIC=22.37	AIC=- 76.81	AIC=-10.21	AIC=16.37
Sparta	AIC=277.10	AIC=372.92	AIC=235.65	AIC=228.35	AIC=231.91	AIC=251.44	AIC= 222.35
Windsor	AIC=-1014.80	AIC=-920.68	AIC=-1016.25	AIC=-1017.59	AIC=- 1069.99	AIC=-1028.91	AIC=-1023.59
All five stations	AIC=-4433.18	AIC=-4422.27	AIC=-4479.50	AIC=-4515.13	AIC=-4425.06	AIC=-4423.78	AIC=- 4553.13

Table 3: Akaike Information Criterion (AIC) values obtained for our five selected weather stations and for our seven models.

The bold values correspond to the optimal criterion per row. Below each model's name, the number p of parameters for n stations is provided.

	Pattern 1	Pattern 2	Pattern 3	Pattern 4
ξ	NA	0.3	0.13	0.26
m for Aledo	NA	0.73	0.81	1.06
m for Aurora	NA	0.28	0.48	1.38
m for Fairfield	NA	1.61	1.24	1.84
m for Sparta	NA	0.46	1.01	1.83
m for Windsor	NA	0.56	0.81	0.96

Table 4: Values of the ξ and m parameters for the five example stations for model (iii)*. Non-applicable (NA) is indicated for pattern 1, since this pattern is associated to Gamma distributions in this model.

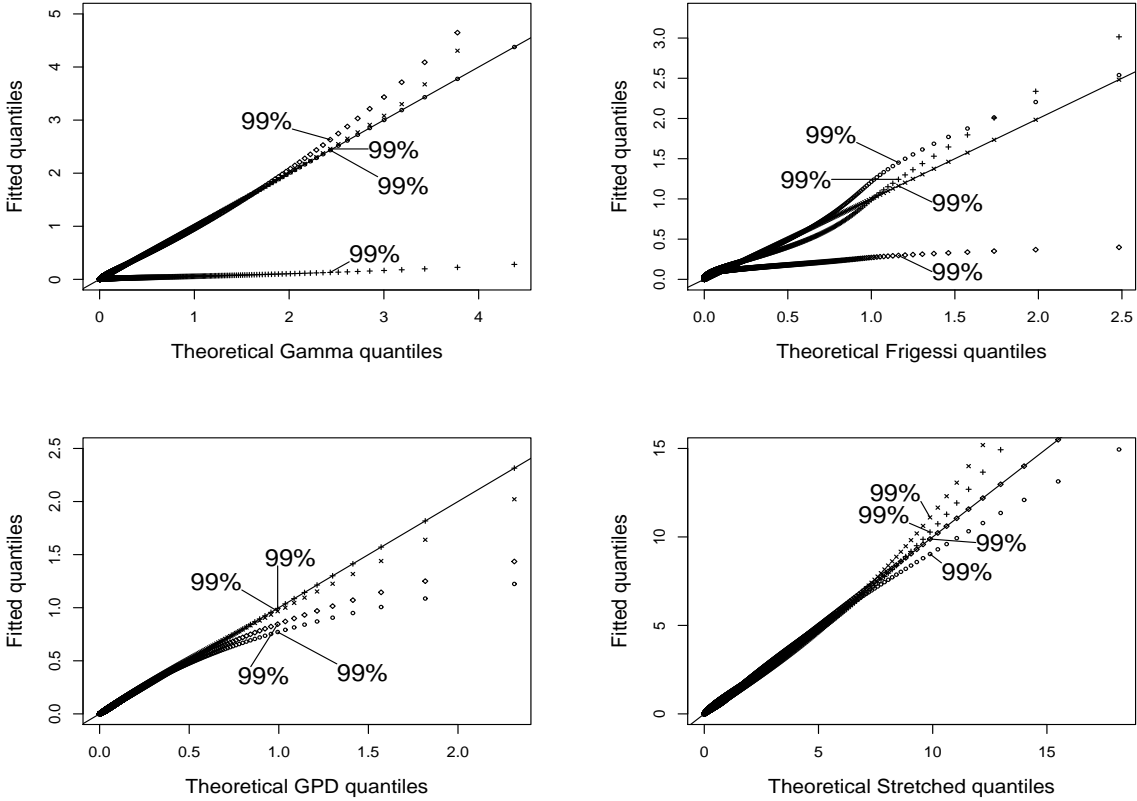
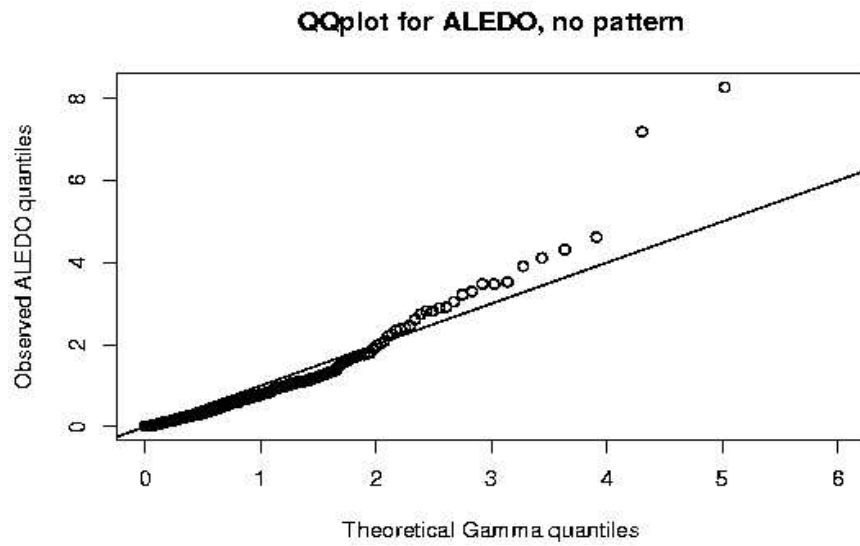
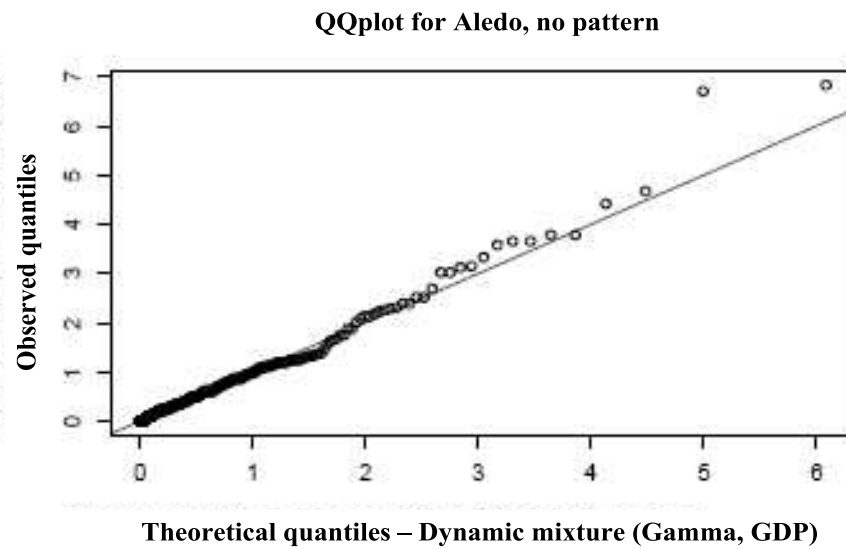


Figure 1: Quantiles-quantiles plots (i.e. theoretical vs. fitted quantiles). The o , \times , $+$, and \diamond signs correspond to the QQplots from the Gamma, mixture, GP and stretched exponential densities, respectively. Each distribution is analytically fitted by a Gamma (left-upper panel), our mixture (right-upper panel), a GP (left-lower panel) and a stretched exponential (right-lower panel) density. The 99% quantile is indicated for each fitted distribution. These graphs mainly tell us that the mixture distribution (\times signs) appears to provide a very good fit in all cases.



(a)



(b)

Figure 2: QQplot for Aledo with (a) Gamma distribution, and (b) our mixture.

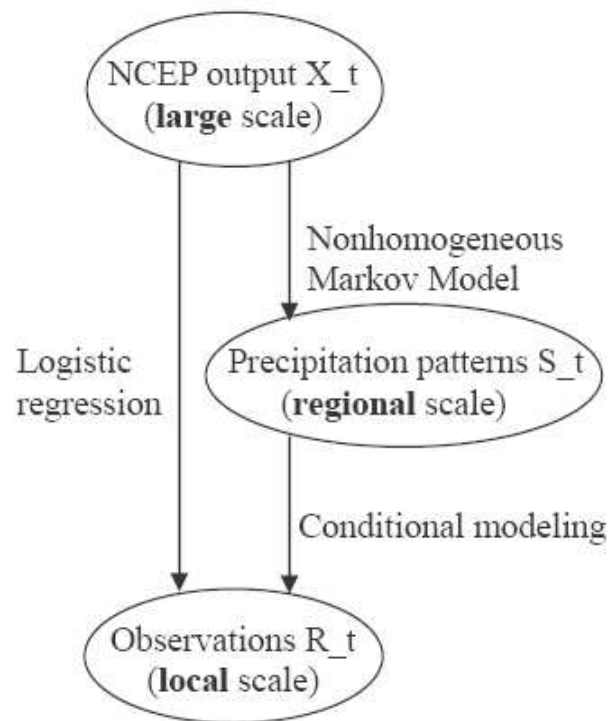


Figure 3: Schematic graph explaining the main components of our downscaling scheme.

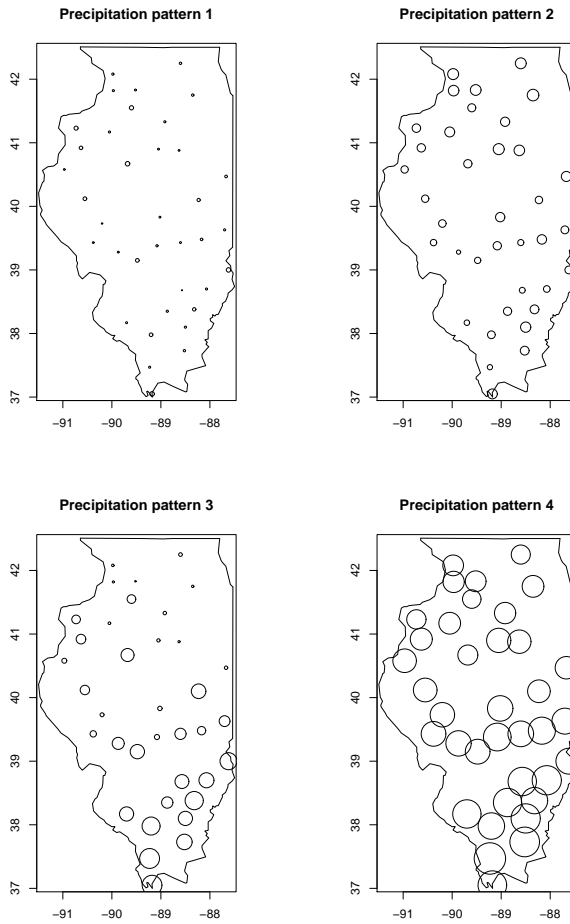


Figure 4: Four station-based precipitation patterns over Illinois derived by the Vrac *et al.* (2006) HAC method, with area proportional to mean rainfall for each cluster.

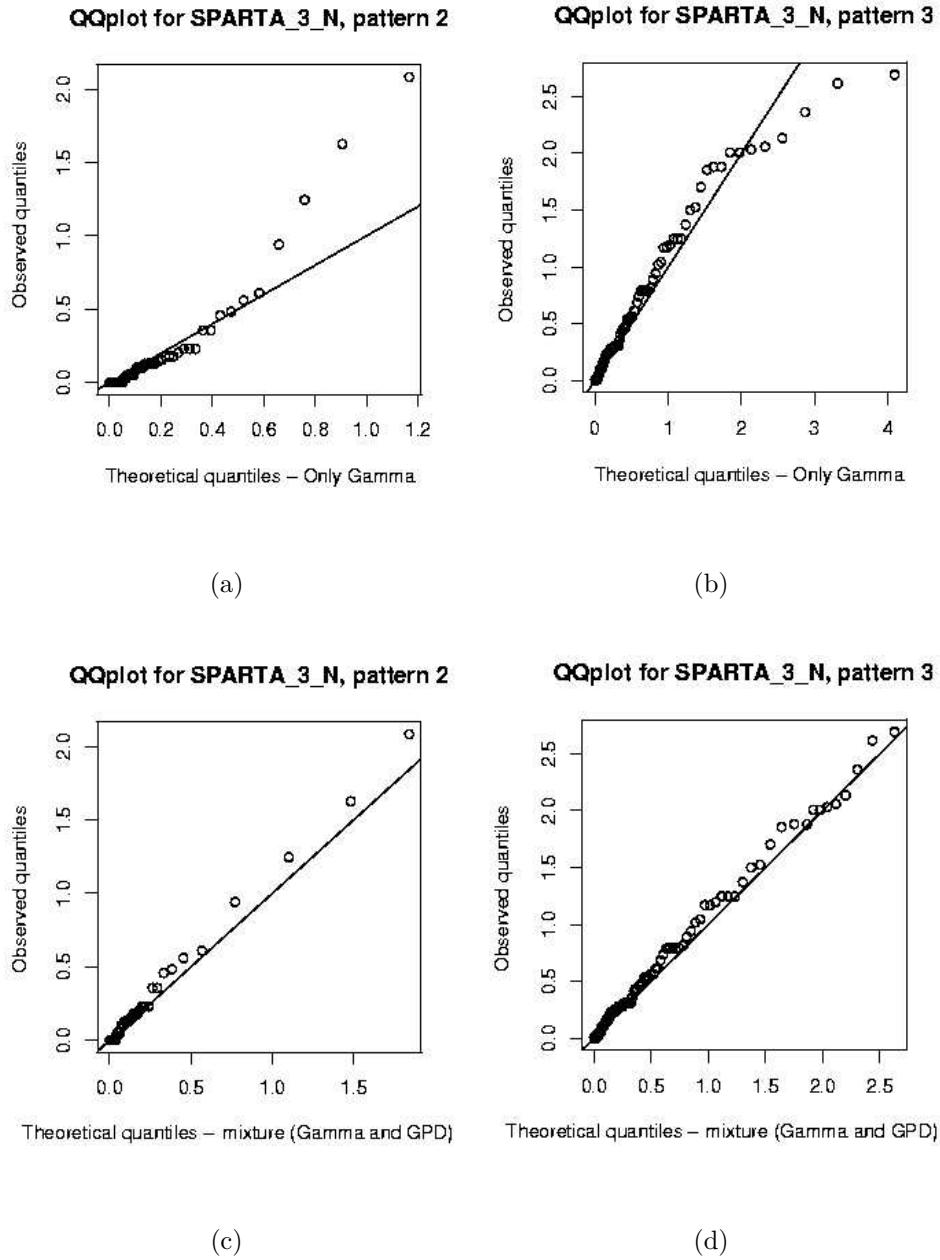


Figure 5: QQplots of precipitation patterns 2 and 3 for station “Sparta”, for function h_β in (11) as a Gamma distribution in (a) and (b) and h_β as a mixture (5) in (c) and (d). Units are cm.

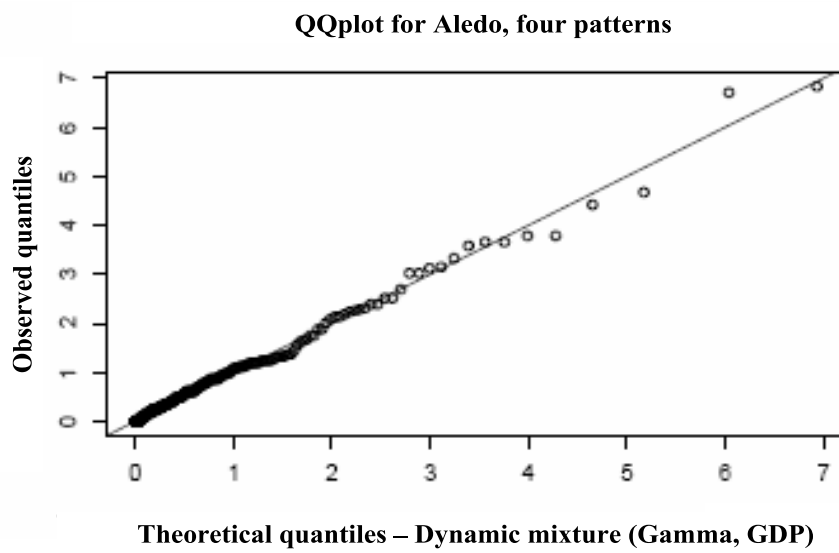
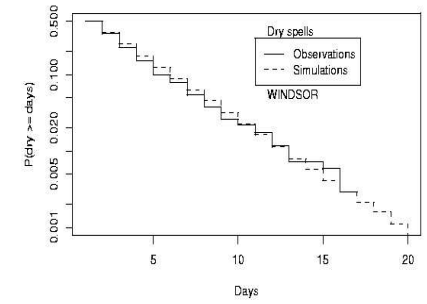
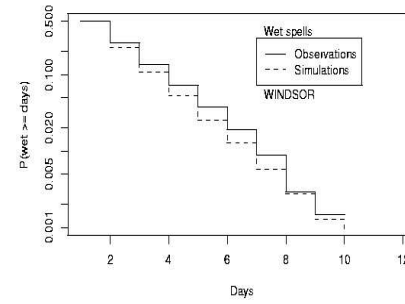
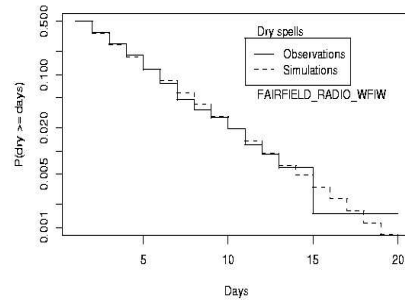
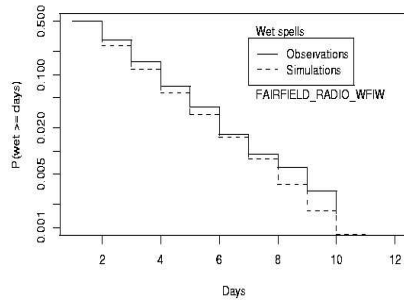
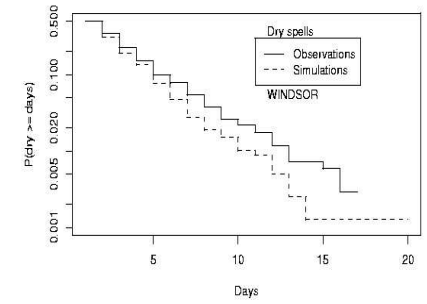
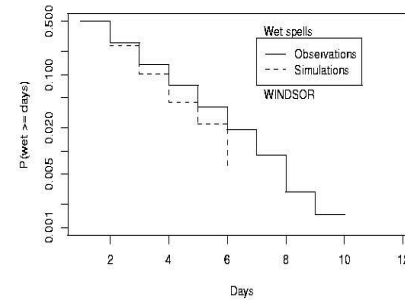
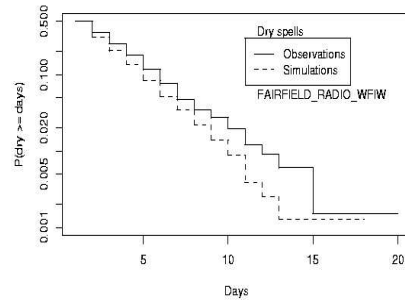
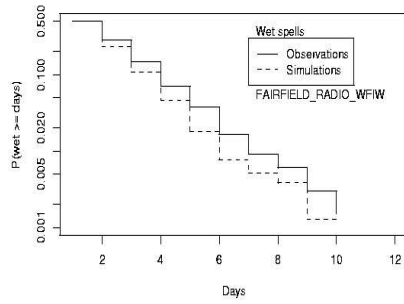


Figure 6: QQplot for Aledo with four patterns and model (iii)*, i.e., Gamma distributions for pattern 1 and mixtures for patterns 2-4 with one ξ per pattern and $\tau = 0$.



(a)

(b)



(c)

(d)

Figure 7: Wet and dry spells probabilities (in log-scale) obtained for Fairfield and Windsor. Upper panels (a) and (b): the “4 patterns” approach; lower panels (c) and (d): the “no pattern” approach.