# Regional frequency analysis conditioned on large-scale atmospheric or oceanic fields

Benjamin Renard, U. Lall

HAL Id: hal-01192599
https://hal.science/hal-01192599

Submitted on 3 Sep 2015

### RESEARCH ARTICLE

**Key Points:**
- Conditioning hydrologic variables on spatial fields rather than climate indices
- Multisite probabilistic prediction of the target hydrologic variable
- Predictive skill is illustrated with a case study on Mediterranean floods

**Correspondence to:**
B. Renard,
benjamin.renard@irstea.fr

# Regional frequency analysis conditioned on large-scale atmospheric or oceanic fields

**Benjamin Renard[1] and Upmanu Lall[2]**

[1]Irstea, UR HHLY Hydrology-Hydraulics, Lyon, France, [2]Department of Earth and Environmental Engineering, Columbia University, New York, USA

**Abstract** Many studies report that hydrologic regimes are modulated by large-scale modes of climate variability such as the El Niño Southern Oscillation (ENSO) or the North Atlantic Oscillation (NAO). Climate-informed frequency analysis models have therefore been proposed to condition the distribution of hydrologic variables on climate indices. However, standard climate indices may be poor predictors in some regions. This paper therefore describes a regional frequency analysis framework that conditions the distribution of hydrologic variables directly on atmospheric or oceanic fields, as opposed to predefined climate indices. This framework is based on a two-level probabilistic model describing both climate and hydrologic data. The climate data set (predictor) is typically a time series of atmospheric of oceanic fields defined on a grid over some area, while the hydrologic data set (predictand) is typically a regional data set of station data (e.g., annual average flow at several gauging stations). A Bayesian estimation framework is used, so that a natural quantification of uncertainties affecting hydrologic predictions is available. A case study aimed at predicting the number of autumn flood events in 16 catchments located in Mediterranean France using geopotential heights at 500 hPa over the North-Atlantic region is presented. The temporal variability of hydrologic data is shown to be associated with a particular spatial pattern in the geopotential heights. A cross-validation experiment indicates that the resulting probabilistic climate-informed predictions are skillful: their reliability is acceptable and they are much sharper than predictions based on standard climate indices and baseline predictions that ignore climate information.

## 1. Introduction

The variability of the climate system is a primary driver of the variability of hydrological regimes. Consequently, identifying meaningful climate variables that may be associated with the changing regional hydrologic state is an important building block for understanding and prediction of changing hydrologic risk. The goal of this paper is to propose a method that identifies the underlying ocean-atmosphere pattern and uses it for hydrologic prediction at selected time scales and lead times.

### 1.1. Large-Scale Climate Variability

Climate is expressed through specific modes of variability, such as the North Atlantic Oscillation (NAO) or the El Niño Southern Oscillation (ENSO), among many others. These modes of variability reflect large-scale climatic processes affecting large areas, and their temporal variability may include some low-frequency component ranging from a few years to a few decades. For instance, the NAO describes the strength of the atmospheric pressure difference between the Icelandic low and the Azores anticyclone [*van Loon and Rogers*, 1978; *Hurrell*, 1995; *Hurrell and VanLoon*, 1997]. The ENSO describes the interaction between the tropical oceans and atmosphere which manifests as band-limited interannual oscillations in sea level pressure (SLP) and sea surface temperature (SST) anomalies in the equatorial Eastern Pacific, with global teleconnections [e.g., *Kousky et al.*, 1984].

Such modes of variability are quantified by means of climate indices: for instance, the Niño3.4 index is defined as the SST anomaly in a particular region of the equatorial Pacific Ocean (see https://www.ncdc.noaa.gov/teleconnections/enso/indicators/sst.php). A NAO index is derived [*Barnston and Livezey*, 1987] as the first component of a rotated principal component analysis (PCA) of standardized 500 hPa geopotential height anomalies over the Northern Hemisphere (20N-90N; see http://www.cpc.ncep.noaa.gov/products/precip/CWlink/ENSO/verf/new.nao.shtml).

### 1.2. The Effect of Large-Scale Climate Variability on Hydrology

Modes of climate variability such as NAO or ENSO are associated with atmospheric/oceanic circulation patterns over a large spatial extent. Consequently, their effects on hydrologic variables have been reported in many regions of the globe. As an illustration, positive NAO phases correspond to warm and wet winters in Northern Europe and precipitation deficits in Southern Europe and North Africa [*Hurrell and VanLoon*, 1997]. The effect of NAO on streamflow variables has also been reported in many European regions [*Shorthouse and Arnell*, 1999; *Pociask-Karteczka*, 2006], in particular in Northern Europe [e.g. *Wilby et al.*, 1997; *Kiely*, 1999; *Stahl et al.*, 2001; *Kingston et al.*, 2006], and in the Iberian peninsula [e.g., *Trigo et al.*, 2004; *Vicente-Serrano and Cuadrat*, 2007]. The impact of ENSO has an even broader spatial extent, since its effects have been reported in many regions of the world (e.g., *Kripalani and Kulkarni* [1997] in South-East Asia; *Gershunov and Barnett* [1998] in the U.S.; *Grimm and Tedeschi* [2009] in South America; *Cai et al.* [2010] in Australia; *Shaman and Tziperman* [2010] in the Mediterranean; and *Philippon et al.* [2012] in South Africa).

Climate-informed frequency analysis frameworks have been proposed to account for the effect of climate variability. The general strategy is to abandon the assumption that the target hydrologic variable is identically distributed; instead, its distribution is conditioned on selected climate indices. For instance, ENSO-dependent frequency analysis models have been used by *El Adlouni et al.* [2007] and *Shang et al.* [2011] in California, *Aryal et al.* [2009] and *Sun et al.* [2014] in Australia, among many others.

### 1.3. Limitations of Standard Climate Indices

Standard climate indices such as NAO or ENSO result from a dimensionality reduction exercise: they attempt to summarize the information contained in a full spatial field into a single number. The associated information loss may result in a loss of predictive power for hydrology. In particular, some regions seem to be out of the area of influence of these climate indices. As an illustration, *Giuntoli et al.* [2013] reported that the influence of winter NAO is hardly discernable for French hydrologic regimes, despite the fact that this influence is evident in nearby regions of Northern and Southern Europe. Similarly, *Grantz et al.* [2005] discussed the limited predictive capability of standard climate indices such as the Southern Oscillation Index (SOI) or the Pacific Decadal Oscillation index (PDOI) for some catchments in the Western U.S. Further discussion of the limitations of standard climate indices can also be found in *Westra and Sharma* [2009].

*Lavers et al.* [2013] illustrated the reason why standard climate indices are sometimes poor predictors. They mapped the correlation between monthly precipitation at one given location in Europe and sea-level pressure over an extended North Atlantic region. While a clear NAO pattern appears in some seasons and some regions (e.g. Northern and Southern Europe in winter), distinct patterns appear for other locations or seasons, suggesting that atmospheric circulation patterns other than that described by NAO indices may be more useful for predicting precipitation in such locations/seasons.

### 1.4. Identifying Climate Patterns Relevant to Regional Hydrology

The topic of identifying statistical relationships between climate fields and hydrologic variables has been largely addressed in the literature, which is unsurprising given the limitations of standard climate indices described above. Several approaches have been explored to identify modes of climate variability that are relevant for hydrologic prediction in some specific target region. Arguably, the most natural approach is based on the use of correlation maps, as illustrated for instance by *Grantz et al.* [2005] and *Lavers et al.* [2013]. This is first useful as an exploratory tool, to identify in the oceanic or atmospheric field relevant regions of influence for the target hydrologic variable. Moreover, it may be used to define "customized" climate indices, e.g., by averaging the atmospheric or oceanic field over regions of high correlation as done in *Grantz et al.* [2005]. These customized climate indices can then be used as explanatory variables in a regression model to predict the hydrologic variable.

Alternative methods stemming from the data mining literature can also be used. Such methods include the Canonical Correlation Analysis (CCA), or the Singular Value Decomposition (SVD) of a cross-covariance matrix, as discussed by *Bretherton et al.* [1992]. These methods attempt to find the linear combination of explanatory variables (e.g. SST) that best explains a linear combination of the target variables (e.g., a regional hydrologic data set). More recently, *Westra et al.* [2008] and *Westra and Sharma* [2009] proposed using the Independent Component Analysis (ICA) to build multisite predictive models using oceanic field information (SST).

### 1.5. The Quest for Explicit Probabilistic Models

While the approaches described above may yield skillful predictions for specific hydrologic variables [e.g. *Grantz et al.*, 2005], some limitations restrict their use. In particular:

1. They are based on correlation/covariance computations, which are not necessarily optimal for heavily non-Gaussian data (e.g., extreme data, count data, or binary occurrence data).

2. While these methods are of interest to identify relevant climate patterns, they do not directly provide probabilistic predictions at several sites. This is an important distinction: while establishing and describing relationships between climate and hydrology is a valuable endeavor, it is not the same thing as building a probabilistic predictive model based on climate information.

3. A complete treatment of uncertainties, including the uncertainty in the identified explanatory climate pattern, is difficult.

These limitations necessitate the derivation of a fully probabilistic model to describe both hydrologic and climate data, with the aim of identifying the most relevant mode of climate variability for a region and simultaneously using it for making hydrologic prediction at several target sites. An explicit probabilistic model for hydrologic/climate data has the following advantages: (i) it is easily interpretable since it directly describes the original data set rather than some component(s) extracted from it; (ii) it allows making predictions at individual target sites; (iii) it allows using standard statistical tools for inference, prediction, and uncertainty quantification, e.g., Bayesian methods; (iv) it makes all assumptions explicit, which allows isolating them and subject each of them to empirical scrutiny; and (v) it offers the flexibility to modify model structure and assumptions if needed, without modifying the whole inference framework.

The main objective of this paper is to present a regional frequency analysis framework based on such an explicit probabilistic model.

### 1.6. Outline of the Paper

The paper is organized as follows. The theoretical basis of the proposed probabilistic model, the inference equations, and the use of the model for prediction are provided in section 2. A case study assessing the relationship between geopotential heights over the North-Atlantic region and the frequency of floods in 16 catchments in Mediterranean France is provided in section 3. Current limitations and directions for future work are then discussed in section 4, before summarizing the main outcomes of this work in the concluding section.

## 2. Theory

### 2.1. Notation

Let $\boldsymbol{y} = (y(x,t))_{x=1:N_x, t=1:N_t}$ denote observations of the target hydroclimatic variable to be predicted. Typically, $\boldsymbol{y}$ represents the hydrologic data, observed at $N_x$ locations over $N_t$ time steps. For instance $y(x,t)$ might be the annual mean flow recorded at site $x$ for year $t$.

Moreover, let $\boldsymbol{\phi} = (\phi(s,t))_{s=1:N_s, t=1:N_t}$ denote observations from an explanatory hydroclimatic variable, which will be used to predict the target variable. Typically, $\phi$ represents the climate data observed at $N_s$ locations during $N_t$ time steps. It is assumed that climate data have been preliminarily centered (by removing the empirical mean at each location $s$). Note that the space index is not the same for the two fields. Typically climate data may be taken from reanalyses [*Kalnay et al.*, 1996; *Uppala et al.*, 2005], and are hence available on a regular grid (with $N_s$ grid points), while hydrologic data are station data ($N_x$ rain gauges or gauging stations). Without loss of generality, one can use climate data that leads the hydrologic data by a certain time $T$ (e.g., if $y(x,t)$ denotes autumn precipitation totals, $\phi(s,t)$ may denote summer mean SST for the same year). As long as the sampling interval is the same and the time vector length is $N_t$ for both data sets, they do not even need to refer to the same averaging intervals (e.g., if $y(x,t)$ denotes the annual minimum streamflow, $\phi(s,t)$ may denote winter mean geopotential heights). The case where the climate and the hydrologic data are concurrent may be useful for downscaling information from General Circulation Model (GCM) climate simulations for retrospective or climate change scenarios. The case where the climate predictors lead the hydrologic variables is of interest for hydrologic forecasts at seasonal or other lead times. In

both cases, the probabilistic formulation allows one to assess the potential risk profile associated with a hydrologic state, accounting for the uncertainty of parameter estimation given the training data.

Last, the shorthand notation ":" is used to denote all the elements from one dimension of a matrix. For instance, $\phi(:, t)$ denotes the vector containing the spatial field of climate data observed at time $t$, while $\phi(s, :)$ denotes the vector containing the time series at gridpoint $s$. All vectors are row vectors unless noted otherwise.

### 2.2. Full Probabilistic Model
### 2.2.1. Probabilistic Assumptions

A two-level probabilistic model is proposed for describing the variability of predictand hydrologic data and predictor climate data. The first level of the model addresses the predictand data, and can be formalized as follows:

Level 1: Predictand hydrologic data

$$y(x, t) \sim D(\theta(x, t), \boldsymbol{\beta}(x)), \quad \text{where :} \tag{1a}$$

$$\theta(x, t) = \lambda(x)(1 + \tau(t)) \tag{1b}$$

$$\text{Conditional independence in time : } (y(x, t_1) \perp y(x, t_2)) | \tau(t_1), \tau(t_2) \quad \forall t_1 \neq t_2 \tag{1c}$$

$$\text{Conditional independence in space : } (y(x_1, t) \perp y(x_2, t)) \mid \tau(t) \; \forall x_1 \neq x_2 \tag{1d}$$

Equation (1a) states that hydrologic data are realizations from a distribution $D$ whose parameters vary in both space and time. More precisely, the $N_D$ parameters of $D$ are split in two groups:

1. $\theta(x, t)$ is the single parameter varying in both space and time.

2. $\boldsymbol{\beta}(x)$ encompasses the $N_D - 1$ remaining parameters, varying in space but not in time.

The assumption that a single parameter may vary in both space and time is made to avoid tedious notation here: in principle, describing spatiotemporal variations in several parameters of $D$ is possible. Typically, $\theta$ will correspond to a location parameter of $D$ (e.g., the mean for a Gaussian distribution, the location parameter for a Gumbel distribution). However, if one wishes to investigate how the *variability* of the hydrologic data evolves in time, one may also consider a scale parameter as part of $\theta$ (e.g., the standard deviation for a Gaussian distribution, the scale parameter for a Gumbel distribution).

In equation (1b), the spatiotemporal variation of $\theta$ is described using the simple formula $\theta(x, t) = \lambda(x)(1 + \tau(t))$, which can be interpreted as a space-time separability condition:

1. $\lambda(x)$ corresponds to a local effect, and denotes the site-specific marginal value of the parameter $\theta$ (i.e. the value of $\theta(x, t)$ obtained when $\tau(t) = 0$).
2. the term $(1 + \tau(t))$ corresponds to temporal modulations of the local effect $\lambda(x)$. The multiplicative nature of such modulations allows interpreting the values of the temporal pattern $\tau(t)$ as relative deviations (e.g., $\tau(t) = -0.2$ means that $\theta(x, t)$ is 20% below normal).

Note that since the temporal pattern $\tau$ does not vary in space, the model in equation (1b) effectively assumes that temporal modulations are identical for all sites. This is a rather strong hypothesis that will be further discussed in the discussion section 4.2. In addition, note that a link function may be used in equation (1b) to restrict the range of parameter $\theta$: for instance, applying equation (1b) to $\log(\theta)$ ensures that $\theta$ is always positive [see *McCullagh and Nelder*, 1989, for more details on link functions].

Equation (1c) assumes hydrologic data are independent in time given the values taken by the temporal pattern $\tau(t)$. In other words, it is assumed that any temporal dependence in the hydrologic data is induced by the temporal dependence of the temporal pattern.

Last, equation (1d) assumes hydrologic data are independent in space, conditional on the temporal pattern $\tau(t)$. Note that conditioning on the temporal pattern is central in this assumption. Indeed, observations from nearby sites are likely to be dependent, but it is assumed that this dependence is explained by the common temporal pattern modulating the value of the parameter $\theta$ in time.

The second level of the model is devoted to predictor climate data, and can be formalized as follows:

Level 2: Predictor climate data

$$\phi(:,t) \sim MG(\mu(:,t), \boldsymbol{V}) \tag{2a}$$

$$\mu(s,t) = \psi(s) * \tau(t) \tag{2b}$$

$$\boldsymbol{V} = \sigma^2 \Sigma_\gamma, \qquad \Sigma_\gamma(i,j) = \Gamma(||s_i - s_j||, \gamma) \tag{2c}$$

$$\text{Conditional independence in time}: (\phi(s,t_1) \perp \phi(s,t_2)) | \tau(t_1), \tau(t_2) \quad \forall t_1 \neq t_2 \tag{2d}$$

Equation (2a) states that at each time step, the climate field $\phi(:,t)$ is a realization from a Multivariate Gaussian (MG) distribution, with mean vector $\mu(:,t)$ and covariance matrix $\boldsymbol{V}$. The mean vector is decomposed as the product of a time-invariant spatial pattern $\psi(s)$ and a space-invariant temporal pattern $\tau(t)$ (Equation (2b)). This is the same pattern that was used in equation (1), and it provides the link between the predictor and the predictand. The covariance matrix $\boldsymbol{V}$ is written as $\sigma^2 \Sigma_\gamma$ (equation (2c)), where $\Sigma_\gamma$ is a correlation matrix and is parameterized by a spatially stationary distance-based dependence relationship with unknown parameters $\gamma$. Such a relationship can be chosen amongst the many admissible correlogram models existing in the geostatistical literature [e.g. *Chiles and Delfiner*, 1999]. Note that in order to keep notation simple, equation (2c) is effectively restricted to modeling stationary and isotropic fields. However, using more advanced geostatistical models is feasible (e.g., to introduce anisotropy to describe zonal or meridional structures, or nonstationarity of the correlogram in space). Last, equation (2d) assumes that the climate fields are independent in time, conditionally on the temporal pattern.

The factorization formalized in equations (2a)–(2c) can be thought of as a signal-noise decomposition of the climate predictor field: the time-varying mean is the signal, and is factorized using a large-scale spatial pattern $\psi(s)$ associated with the temporal pattern $\tau(t)$. The covariance matrix is the noise, and describes the remaining unexplained variability, which may have a local-scale spatial dependence structure.

### 2.2.2. Interpretation of the Probabilistic Model
The key element in the probabilistic model described in section 2.2.1 is the temporal pattern $\tau(t)$, which is present in both levels and therefore makes the link between them. In level 1, $\tau(t)$ represents a common temporal pattern that affects the distribution of the hydrologic variable in all sites within the target region. The existence of such a common temporal variability might reflect the fact that all sites are subject to the same climate forcing. The question is then to identify this forcing, i.e., to identify the spatial pattern $\psi(s)$ in the climate data that induces this temporal variability $\tau(t)$. This can be achieved by performing a space-time decomposition of the climate data with respect to the temporal pattern $\tau(t)$.

Interestingly, the second level of the model is similar to the probabilistic principal component analysis (pPCA) model proposed by *Tipping and Bishop* [1999]. More precisely, equations (2a) and (2b) correspond to a single-component pPCA, while equation (2c) is slightly more general than pPCA by allowing spatial dependence. In this respect, the proposed model could be interpreted as an extension of PCA, with the first level of the model making the link with the target hydrologic data. However, the aim is not to explain the *internal* climate variability (as PCA does), but rather to maximize the explanation of the *external* hydrologic variable by the identification of the leading mode in the climate data.

Last, note that the space-time decomposition formulas in equations (1b) and (2b) could be modified in some case studies. In particular, the formula $\theta(x,t) = \lambda(x)(1 + \tau(t))$ of equation (1b) expresses temporal variability as relative deviations. The rationale behind this choice is that the order of magnitude of hydrologic data may strongly vary from site to site (e.g. due to different catchment sizes for a streamflow variable). This formulation is a particular case of the more general model $\theta(x,t) = (\lambda_0 + \lambda_1(x))(\tau_0 + \tau_1(t))$. However, the latter model is not identifiable and requires adding two identifiability constraints. The constraints we used here ($\lambda_0 = 0$ and $\tau_0 = 1$) were chosen because they yield an intuitive interpretation for the parameters ($\lambda(x)$ are local "normal" values, $\tau(t)$ are relative deviations from the normal values). However, alternative choices would also be perfectly legitimate: for instance, the constraints $\sum_{x=1}^{N_x} \lambda_1(x) = 0$ and $\sum_{t=1}^{N_t} \tau_1(t) = 0$ would lead to interpreting $\lambda_0$ and $\tau_0$ as spatial and temporal averages.

### 2.2.3. Inference
The model described in section 2.2.1 requires estimating the following unknown quantities:

1. the local effects, λ (size $N_x$)

2. the temporal pattern, τ (size $N_t$)

3. the time-invariant parameters of $D$, $\boldsymbol{\beta}$ (size $(N_D-1)*N_x$)

4. the spatial pattern, ψ (size $N_s$)

5. the parameters describing the geostatistical properties of the climate fields, $\sigma$ and $\gamma$ (size depending of the geostatistical model).

The posterior distribution of these quantities, given observed hydrologic and climate data **y** and $\phi$, can be derived as follows:

$$p(\lambda, \tau, \boldsymbol{\beta}, \psi, \sigma, \gamma | \mathbf{y}, \phi) \propto p(\mathbf{y}|\lambda, \tau, \boldsymbol{\beta}, \psi, \sigma, \gamma, \phi)p(\lambda, \tau, \boldsymbol{\beta}, \psi, \sigma, \gamma | \phi)$$

$$\propto p(\mathbf{y}|\lambda, \tau, \boldsymbol{\beta}, \psi, \sigma, \gamma, \phi)p(\phi|\lambda, \tau, \boldsymbol{\beta}, \psi, \sigma, \gamma)p(\lambda, \tau, \boldsymbol{\beta}, \psi, \sigma, \gamma)$$

$$\propto \underbrace{p(\mathbf{y}|\lambda, \tau, \boldsymbol{\beta})}_{\text{hydrologic likelihood}} \underbrace{p(\phi|\tau, \psi, \sigma, \gamma)}_{\text{climate likelihood}} \underbrace{p(\lambda, \tau, \boldsymbol{\beta}, \psi, \sigma, \gamma)}_{\text{priors}} \tag{3}$$

The first two lines of equation (3) correspond to two successive applications of the Bayes theorem. The third line corresponds to simplifications in the computation of hydrologic and climate likelihoods. Indeed, given the first level of the model in equation (1), the hydrologic likelihood can be computed as:

$$p(\mathbf{y}|\lambda, \tau, \boldsymbol{\beta}) = \prod_{x=1}^{N_x} \prod_{t=1}^{N_t} p_D(y(x,t)|\lambda(x)(1+\tau(t)), \boldsymbol{\beta}(x)) \tag{4}$$

where $p_D(z|\boldsymbol{\eta})$ is the probability density function (pdf) of the distribution $D$ with parameters $\boldsymbol{\eta}$, evaluated at value $z$.

Similarly, given the second level of the model in equation (2), the climate likelihood can be computed as:

$$p(\phi|\tau, \psi, \sigma, \gamma) = \prod_{t=1}^{N_t} p_{MG}\left(\phi(:,t)|\psi(:) * \tau(t), \sigma^2 \boldsymbol{\Sigma}_\gamma\right) \tag{5}$$

where $p_{MG}(\mathbf{z}|\boldsymbol{\mu}, \mathbf{V})$ is the joint pdf of the $N_s$-dimensional multivariate Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\mathbf{V}$, evaluated at vector $\mathbf{z}$.

The posterior distribution in equation (3) is high dimensional and does not yield, in general, explicit estimators. It is therefore explored by means of a Markov Chain Monte Carlo (MCMC) sampler. The adaptive MCMC sampler used in this paper is described by *Renard et al.* [2006]. Alternative samplers, specifically geared toward high-dimensional inferences, may be more efficient [e.g., *Haario et al.*, 2005; *Laloy and Vrugt*, 2012; *Hoffman and Gelman*, 2014]. Moreover, the implementation of a dedicated MCMC sampler, taking advantage of the specific two-level structure of the model, may also be considered. This, however, lies beyond the scope of this paper and is left for future work.

### 2.2.4. Prediction
The ultimate aim of the proposed model is to make a probabilistic prediction for the hydrologic variable $y(x, t^*)$ at prediction time $t^*$, conditionally on the observed climate field $\phi(:, t^*)$. In this predictive context, we assume that inference has already been performed as described in preceding section 2.2.3, and that we observe the predictor climate field $\phi(:, t^*)$. However, the hydrologic data $y(:, t^*)$ are not observed—they are the quantity to be predicted.

This prediction is based on equation (1), and therefore requires estimating the value of the temporal pattern $\tau(t^*)$, which is unknown since the prediction time $t^*$ is not part of the calibration data set. The simplest way to achieve this is to remark that, conditionally on $\hat{\psi}$, $\hat{\sigma}$, and $\hat{\gamma}$, the maximum-likelihood estimator of $\tau(t^*)$ is explicit and can be computed as (see proof in Appendix section A1):

$$\hat{\tau}(t^*) = \frac{\phi(:,t^*)\left(\hat{\sigma}^2 \boldsymbol{\Sigma}_{\hat{\gamma}}\right)^{-1} \hat{\psi}^T}{\hat{\psi}\left(\hat{\sigma}^2 \boldsymbol{\Sigma}_{\hat{\gamma}}\right)^{-1} \hat{\psi}^T} \tag{6}$$

Equation (6) only provides a point estimate of the temporal pattern $\tau(t^*)$. However, a distributional estimate can easily be obtained by propagating the MCMC samples of parameters $\psi$, $\sigma$, and $\gamma$ in equation (6).

### 2.3. An Exploratory Tool

In this section, we describe an exploratory tool based on simplified probabilistic assumptions yielding explicit estimators of the spatial and temporal patterns. This enables simple and fast computations that can be useful for preliminary analyses and to facilitate the tuning of MCMC samplers.

#### 2.3.1. Simplified Probabilistic Assumptions

The first level assumes Gaussian hydrologic data, with a common mean $\tau(t)$ and a common variance equal to one:

Level 1: Hydrologic data

$$\tilde{y}(x,t) \sim N(\tau(t), 1) \tag{7a}$$

$$\text{Independence in space and time, conditionally on } \tau(t) \tag{7b}$$

In practice, these assumptions are too restrictive to be applied to the raw hydrologic data **y**. Consequently, the exploratory tool rather uses normal-score transformed data $\tilde{y}$:

$$\tilde{y}(x,t) = g^{-1}\left(\hat{F}_x(y(x,t))\right) \tag{8}$$

where $g$ is the cumulative distribution function (cdf) of the standard normal distribution and $\hat{F}_x$ is the empirical cdf estimated at site $x$.

The second level is similar to that of the full model of section 2.2, with the exception that a spatial independence assumption is now made:

Level 2: Climate data

$$\phi(s,t) = N(\psi(s) * \tau(t), \sigma^2) \tag{9a}$$

$$\text{Independence in space and time, conditionally on } \tau(t) \tag{9b}$$

#### 2.3.2. Inference

Under the assumptions described in equation (7), the maximum-likelihood estimator of the temporal pattern is simply equal to the spatial mean of transformed data $\tilde{y}$ (see proof in Appendix section A2.1):

$$\forall t = 1 : N_t \quad \tilde{\tau}(t) = \frac{1}{N_x} \sum_{x=1}^{N_x} \tilde{y}(x,t) = mean(\tilde{y}(:,t)) \tag{10}$$

Conditionally on this estimated temporal pattern $\tilde{\tau}$, the maximum-likelihood estimate of the spatial pattern can be derived from equation (9) (see proof in Appendix section A2.2):

$$\forall s = 1 : N_s \quad \hat{\psi}(s) = \frac{\sum_{t=1}^{N_t} \phi(s,t)\tilde{\tau}(t)}{\sum_{t=1}^{N_t} \tilde{\tau}^2(t)} = \frac{sd(\phi(s,:))}{sd(\tilde{\tau}(:))} cor(\phi(s,:), \tilde{\tau}(:)) \tag{11}$$

where $sd(.)$ denotes the empirical standard deviation of a vector and $cor(.,.)$ the empirical correlation between two vectors. Equation (11) therefore states that the estimated spatial pattern is nothing more than a weighted correlation map between the estimated temporal pattern (stemming from hydrologic data) and the spatial fields of climate data.

#### 2.3.3. Prediction

As in section 2.2.4, the temporal pattern at prediction time $t^*$ can easily be estimated as follows (see proof in Appendix section A2.3):

$$\hat{\tau}(t^*) = \frac{\sum_{s=1}^{N_s} \phi(s,t^*)\hat{\psi}(s)}{\sum_{s=1}^{N_s} \hat{\psi}^2(s)} = \frac{sd(\phi(:,t^*))}{sd(\hat{\psi}(:))} cor(\phi(:,t^*), \hat{\psi}(:)) \tag{12}$$

Note the distinction between $\tilde{\tau}$, the temporal pattern estimated from the hydrologic data (equation (10)), and $\hat{\tau}$, the temporal pattern reconstructed from the climate data (equation (12)).

## 3. Case Study

This case study illustrates the application of the proposed model. The general objective is to predict the number of flood events in Mediterranean catchments, given observed geopotential heights over the North-Atlantic region. Moreover, we compare these predictions with the ones obtained from standard climate indices.

### 3.1. Data
### 3.1.1. Climate Data

Geopotential heights at 500 hPa over the North-Atlantic region (100°W–40°E, 20°N–70°N) for the period 1948–2011 are used. Data are available at a daily time step on a 2.5° regular grid from the NCEP-NCAR reanalysis [*Kalnay et al.*, 1996]. Daily data are averaged over the autumn season (October–December) and centered to yield the climate data $\phi(s, t)$ ($s = 1{:}1197$ gridpoints, $t = 1{:}64$ autumns) that are used as explanatory variables for hydrology.

We choose geopotential heights at 500 hPa because many standard climate indices are derived from this particular variable. In particular, we also use in this case study three climate indices that have been shown to have a noticeable effect on hydrological variables in Europe: the NAO [e.g., *Shorthouse and Arnell*, 1999], the Eastern Atlantic Western Russia oscillation [EAWR, see e.g., *Ionita*, 2014], and the Scandinavian Pattern [SCAND, see e.g., *Wibig*, 1999]. The three time series for the period 1950–2013 have been downloaded from NOAA climate prediction center (http://www.cpc.ncep.noaa.gov/data/teledoc/telecontents.shtml), and averaged over the autumn season.

### 3.1.2. Hydrologic Data

Daily streamflow series for 16 catchments located in Mediterranean France are used (Figure 1), for the same period 1948–2011. The catchment sizes range from 10 to 621 km$^2$. These catchments are typical of Mediterranean hydrologic regimes, with very intense rainfall events producing flash floods during the autumn season (the most active period ranging from October to December). The variable of interest is the number of flood events during autumn (OND). This variable is extracted from daily series using a Peak-Over-Threshold (POT) approach [see *Lang et al.*, 1999, for details and general guidelines]. At each site, the threshold is taken as the 95th percentile of daily streamflows. Moreover, a minimum duration of 6 days between two successive flood events is imposed to avoid counting the same event several times. The number of flood events in autumn at all 16 sites hence constitutes the hydrologic data $y(x, t)$ ($x = 1{:}16$ sites, $t = 1{:}64$ autumns) used as the target variable to be predicted from climate. Note that unlike climate data, missing values exist in hydrologic data, since some stations do not cover the whole period 1948–2011. The treatment of missing data will be described in the following sections.
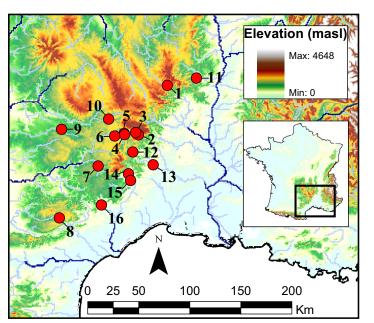
### 3.2. Exploratory Analyses

Before attempting to infer the full model of section 2.2, preliminary analyses are useful to confirm that there is some degree of predictability in the data set, and to suggest possible simplifications. This is achieved by means of the exploratory tool described in section 2.3. The treatment of missing hydrologic data accounts for the fact that the estimated temporal pattern in equation (10) is a spatial average: it is hence reasonable to constrain the number of available stations to a minimum value. We therefore restrict the
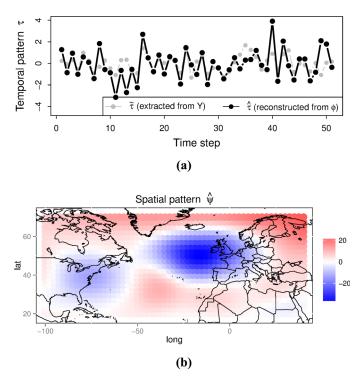


**Figure 1.** Location of the gauging stations used in the case study.

**(a)**



**(b)**

**Figure 2.** Exploratory analysis: estimated temporal (a) and spatial (b) patterns.

analysis to years with at least half of the 16 stations providing nonmissing values. This yields a subset of 51 years from the initial 64-year period 1948–2011. Also note that because the hydrological variable of interest only takes integer values, one needs to decide how to treat ties in the empirical cdf used for the normal-score transformation of equation (8). We decided to follow the "random rank" approach (e.g. for the data $y=(5,8,6,6)$, the ranks may be either $(1,4,2,3)$ or $(1,4,3,2)$, with the assignment of ranks 2 and 3 being chosen at random).

Figure 2a shows the two estimates of the temporal pattern: the pattern estimated from hydrologic data ($\tilde{\tau}$ in equation (10)) and the pattern reconstructed from climate data ($\hat{\tau}$ in equation (12)). Despite some differences, the two patterns are in acceptable agreement (the correlation coefficient is approximately 0.6), which indicates that the temporal variability observed in the hydrologic data can be reconstructed (to some reasonable extent) using climate information.

Figure 2b shows the spatial pattern, which is characterized by a strong negative anomaly over the midlatitudes of the Eastern Atlantic, and positive anomalies at higher latitudes (especially over Scandinavia). This pattern can be interpreted in relation to the target hydrologic data: strong low pressures off the coast lead to a counterclockwise circulation that would steer moisture toward Mediterranean France. In addition, high pressures over Scandinavia and North-Eastern Europe limit the possible trajectories of this moisture flux.

The exploratory tool is also used to assess the sensitivity of the results to the resolution of the climate grid. Indeed, the full-resolution grid comprises 1197 gridpoints, which will be computationally penalizing in the context of the full analysis using MCMC sampling. However, the pattern in Figure 2b has a large spatial scale and is quite smooth, suggesting that similar results might be obtained with a lower-resolution grid. Figure 3 therefore compares the temporal patterns $\hat{\tau}$ (reconstructed from climate data) obtained with various grid resolutions. It shows that even a low-resolution 10°-grid is sufficient to reconstruct this temporal pattern, with barely noticeable differences with the full-resolution 2.5°-grid. This interesting result suggests that the spatial resolution may not be of primary importance when one studies large-scale climate patterns. In fact, when the relevant spatial pattern is large and smooth, using more finely gridded products may just introduce redundant information that dramatically increases the dimension of the model, thereby complicating its inference. Since we may not know a priori what resolution is adequate, an exploratory analysis such as the one performed here is very useful to inform this decision.

### 3.3. Hydrologic Predictions Using the Full Model

This section implements the analysis using the full model, using all available hydrologic data over the 64-year period 1948–2011. Moreover, given the results of the exploratory analysis described above, we subsample the climate data on a low-resolution 10°-grid (90 gridpoints, compared to the 1197 gridpoints in the original 2.5° resolution).

#### 3.3.1. Model Specifications

We assume a Poisson distribution for the hydrologic data, with a rate parameter varying in both space and time as described in equation (1b). Note that because the Poisson distribution has a single parameter, there
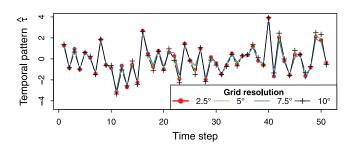
**Figure 3.** Exploratory analysis: sensitivity of the estimated temporal pattern to the spatial resolution of the climate grid.

is no "additional parameter" $\beta$ in this particular case study. The treatment of missing hydrologic data is straightforward in the full model: the terms corresponding to missing values are simply ignored in the double product of the likelihood equation (4).

Moreover, we use a simple exponential correlogram for describing spatial dependence. Following the notation of equation (2c), this can be formalized as:

$$\Sigma_\gamma(i,j) = \Gamma(||s_i - s_j||, \gamma) = \exp\left(-||s_i - s_j||/\gamma\right) \tag{13}$$

In order to compare the predictions from the full model with the predictions obtained using standard climate indices, we also use the following climate-informed model for the hydrologic data:

$$y(x,t) \sim Poisson(\lambda(x)(1 + \omega C(t))) \tag{14}$$

where $C(t)$ is one the climate indices (NAO, EAWR or SCAND) and parameter $\omega$ controls the effect of the climate index on the distribution of hydrologic data.

Last, in order to compare the climate-informed predictions with baseline predictions, we also use the following climate-independent model for hydrologic data:

$$y(x,t) \sim Poisson(\lambda(x)) \tag{15}$$

This model simply estimates a time-invariant Poisson distribution at each site. The resulting predictions can be viewed as "climatology" predictions, and can be considered as the baseline upon which any climate-informed prediction should improve.

### 3.3.2. Estimation of the Full Model

Samples (100,000) from the posterior distribution in equation (3) are generated using MCMC sampling (see section 2.2.3). The Gelman-Rubin statistic [*Gelman and Rubin*, 1992] is below 1.2 for all inferred quantities with this number of iterations. The dimensionality of the inference is equal to 172 in this case study (16 stations + 64 time steps + 90 gridpoints + 2 geostatistical parameters). Computing times are of the order of 1 h on a standard laptop (2.40 GHz CPU).

Figure 4 shows the estimation of local effects $\lambda(x)$ for all 16 stations, which can be interpreted as the "normal" number of flood events expected on each site. Moderate variations between sites can be observed, with the posterior medians ranging from approximately 1.5 to 2.5.

Figure 5 shows the estimated temporal and spatial patterns. Figure 5a focuses on the temporal pattern, with the line representing the modal estimate (maximizing the posterior pdf) and the boxplots representing the uncertainty. Figure 5a shows that the tempo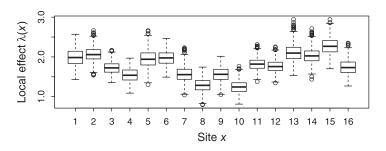ral pattern can be identified quite precisely, and the uncertainty in individual estimates at each time step is quite small compared to the temporal variability of the pattern. Figures 5b and 5c show the corresponding spatial pattern, both in terms of modal estimate (Figure 5b) and uncertainty (Figure 5c). The modal estimate is similar to the pattern that was obtained during the exploratory analysis (Figure
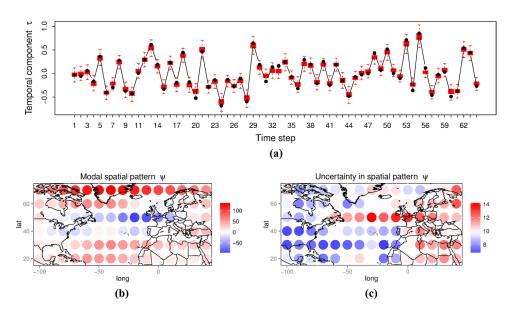


**Figure 4.** Estimation of local effects. Boxplots from the posterior distribution of $\lambda(x)$ at each site $x$ are shown.

**Figure 5.** Estimation of temporal and spatial patterns. (a) Temporal pattern (boxplots represent the posterior distribution, the line and points represent the modal estimate); (b) spatial pattern (modal estimate); and (c) uncertainty in the spatial pattern (measured by the posterior standard deviation).

2b), with negative anomalies over the midlatitudes of the Eastern Atlantic, and positive anomalies at higher latitudes. The uncertainty pattern (Figure 5c) also reveals an interesting structure, with areas of lower uncertainty being located in the western part of the spatial domain, while areas of larger uncertainty are found in the eastern part.

It is interesting to compare the temporal pattern inferred for the full model, $\hat{\tau}(t)$, with the ones induced by the climate indices, $\hat{\omega}C(t)$ ($\hat{\omega}$ is the modal estimate from 100,000 MCMC samples). Figure 6 shows the result of this comparison and yields several interesting insights. First, the temporal variability associated with EAWR and NAO is rather weak, suggesting that these indices have little effect on the distribution of hydrologic data. Index SCAND yields a larger temporal variability, which, however, remains smaller than the one resulting from the full model. The predictions from the full model will therefore be more dynamic than the predictions from any of the climate indices. Second, the information carried by the temporal pattern from the full model ($\hat{\tau}(t)$) is distinct from the information carried by climate indices: correlation coefficients with NAO, EAWR, and SCAND temporal patterns are equal to 0.60, 0.11, and 0.34, respectively. In other words, the mode of climate variability identified with the full model does not correspond to any of the predefined modes (NAO, EAWR, and SCAND).

### 3.3.3. Prediction
The full model is now applied to make a probabilistic prediction of the hydrologic variable (number of flood events in autumn), given the observed climate data (autumn average of geopotential heights). Figure 7 illustrates the principle of this climate-informed prediction. Given the observed geopotential field, the corresponding value for the temporal pattern can be predicted using equation (6). In turn, this value can be plugged in equation (1b) to estimate the distribution of the hydrologic variable at one particular site. Note that throughout this section, we are actually using predictive distributions: instead of simply using "values" corresponding to point estimates of inferred quantities, we systematically propagate forward the uncertainties embedded in MCMC samples in order to take them into account in the predictions.

Figure 7 shows the predictive distributions at site 5, conditional on three particular realizations of the geopotential field:

1. The field for 2000 is similar to the spatial pattern shown in Figure 5b. This results in a "high" value of the predicted temporal pattern (modal estimate $\hat{\tau}(t^*) \approx 0.7$), which in turn results in a "higher-than-normal" prediction for the number of flood events (green distribution in Figure 7).

2. The field for 1992 shows only small anomalies, resulting in a near-zero value of the predicted temporal pattern (modal estimate $\hat{\tau}(t^*) \approx -0.1$), which in turn results in a "normal" prediction for the number of flood events (red distribution in Figure 7).
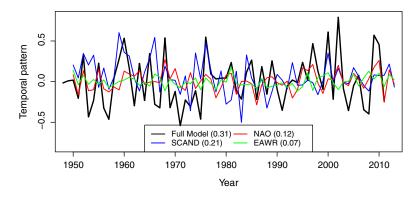
**Figure 6.** Comparison of the temporal patterns induced by climate indices (EAWR, NAO, SCAND) with the temporal pattern from the full model. Numbers in brackets are the standard deviations of each time series.

3. The field for 1975 shows strong anomalies that are very dissimilar to the spatial pattern shown in Figure 5b. The predicted temporal pattern is hence largely negative (modal estimate $\hat{\tau}(t^*) \approx -0.6$), which in turn results in a "lower-than-normal" prediction for the number of flood events (black distribution in Figure 7).

### 3.3.4. Cross Validation

In order to assess the quality of climate-informed probabilistic predictions, a leave-one-out cross-validation experiment is carried out. All results in the remainder of this section are hence validation results in a truly predictive context, with the value of the hydrological variable at prediction time being ignored for both estimation and prediction.

The quality of these probabilistic predictions is evaluated using two approaches. The first approach is based on the log pseudo marginal likelihood (*LPML*), which is a measure of predictive ability [*Gelfand et al.*, 1992]. Let $p_{t*}(z|\boldsymbol{y}(:,-t^*),\phi)$ denote the pdf of the predictive distribution at time $t^*$ (evaluated at $z$), derived as described in previous section 3.3.3. The notation $\boldsymbol{y}(:,-t^*)$ is used to recall that in the cross-validation context used here, the hydrologic data at prediction time $t^*$ are not used for inference and prediction. For one given site $x$, the *LPML* is then defined as:
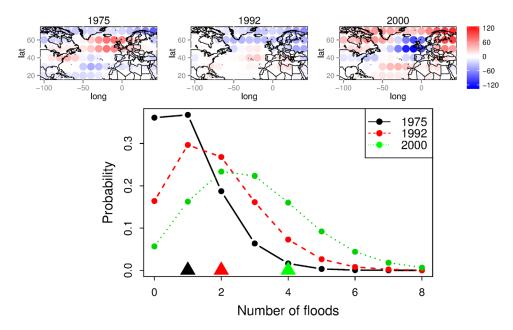


**Figure 7.** Climate-informed prediction: predictive distribution of the number of floods for site 5, conditional on the observed geopotential height fields for the years 1975, 1992, and 2000. Triangles represent the observed number of floods.
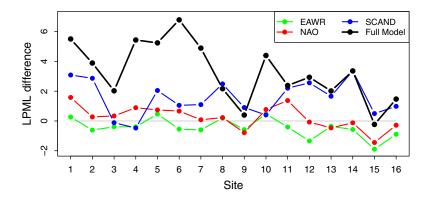
**Figure 8.** Cross validation: difference of log pseudo marginal likelihoods (LPML) between climate-informed models and the climate-independent baseline model.

$$LPML(x) = \sum_{t^*=1}^{Nt} \log\left(p_{t*}\left(y(x,t^*)|\boldsymbol{y}(:,-t^*),\boldsymbol{\phi}\right)\right) \quad (16)$$

The predictive ability of two competing models can then be compared by computing the difference $\Delta LPML$ between their $LPML$s. The exponential of this difference can be interpreted as a pseudo Bayes factor [*Gelfand et al.*, 1992].

Figure 8 shows the $LPML$ difference between the climate-informed models (the full model + the models using climate indices) and the climate-independent baseline model (15). According to Kass and Raftery's scale [*Kass and Raftery*, 1995], the evidence for the models using EAWR or NAO indices is weak ($\Delta LPML < 1$) for almost all sites, suggesting that they are not performing noticeably better than the baseline climate-independent model. Predictions using the SCAND index show a larger improvement over the baseline, with $LPML$ differences corresponding to a positive evidence ($1 < \Delta LPML < 3$) on many sites. Finally, the full model shows the largest improvement over the baseline: evidence is positive ($\Delta LPML > 1$) for most sites, and is often strong ($\Delta LPML > 3$) or very strong ($\Delta LPML > 5$).

The second approach for evaluating the quality of probabilistic predictions is based on reliability/sharpness diagrams, which are commonly used in the context of probabilistic forecasting [e.g. *Gneiting et al.*, 2007]. The reliability/sharpness diagrams in Figure 9 are based on the probability of occurrence of a given event (for instance, "observing two floods or more") and can be interpreted as follows:

1. The *x* axis corresponds to predicted probabilities, discretized into bins (here, 10 bins of width 0.1 are used). The blue line corresponds to the frequency of each bin, i.e. the frequency with which the event "$N>1$" was predicted with probability between 0 and 0.1, between 0.1 and 0.2, etc. The blue line hence illustrates the sharpness of the prediction: sharp ("courageous") predictions correspond to predicted probabilities spanning the whole 0–1 segment, while nonsharp predictions correspond to weakly varying predicted probabilities. Consequently, the peakier the blue line, the less sharp the prediction.

2. In addition to being sharp, the predictions should be reliable in the following sense: across all cases where the event "$N>1$" was predicted with probability 0.2, the event should indeed occur with frequency $\sim$20%. Reliability is hence illustrated by the red staircase curve, which counts the relative frequency of the event "$N>1$" within each probability bin. A reliable prediction yields a red curve close to the diagonal [see *Bröcker*, 2007, for additional discussions on the reliability diagram].

3. A possible visual bias in Figure 9 stems from the fact that the number of elements within each bin strongly varies. Consequently, an observed frequency of 10% does not carry the same information depending on whether it has been computed within a bin of 10 or 10,000 elements. In order to account for this issue, 95% predictive intervals around the diagonal are added to the plot. These intervals correspond to the expected frequency of occurrence among the *n* elements within a given bin, under the reliability assumption. They can easily be computed as binomial quantiles.

Figure 9 shows these diagrams for the event "$N>1$", for all sites pooled together, and for all models (the climate-informed models + the climate-independent model of equation (15)). Note that site-specific
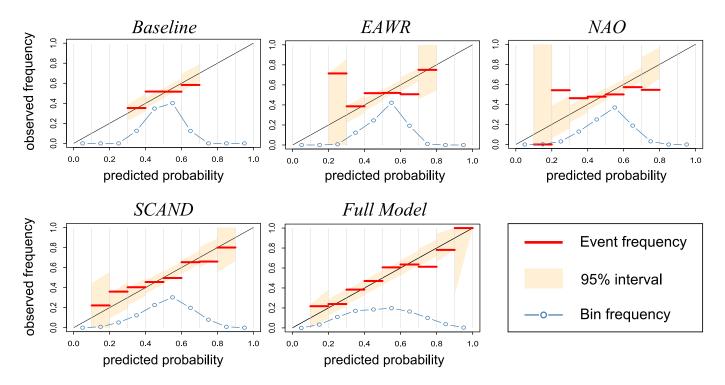
**Figure 9.** Cross validation: reliability/sharpness diagrams associated with the prediction of the event "$N>1$."

reliability diagrams could be derived, but in practice the number of data (64 at most) is too low to yield any meaningful conclusion.

Predictions from the climate-independent model have a low sharpness, as shown by the peaky blue curve. Note that for one given site, this probability does not vary in time with the climate-independent model: consequently, the variability observed in Figure 9 is purely spatial (the event does not have the same probability from site to site because of the local parameters $\lambda(x)$ in equation (15)). The reliability (illustrated by the red staircase curve) is acceptable.

Predictions based on EAWR or NAO indices only show a marginal improvement over the baseline in terms of sharpness. NAO-based predictions even show a worsening in terms of reliability. Predictions based on the SCAND index appear more skillful: they improve upon the baseline predictions in terms of sharpness (flatter blue curve), with no noticeable loss of reliability (red staircase curve in acceptable agreement with the 1:1 line). Finally, predictions from the full model show the largest improvement over the baseline. They are the sharpest of all predictions, and this improvement comes at no cost in terms of reliability.

## 4. Discussion

This section discusses current limitations of the proposed model and highlights avenues for further improvement and generalization.

### 4.1. Model Checking Tools

As in any model-based approach, the assumptions made to build the probabilistic model should be scrutinized and, if proven unrealistic, revised. In the case study of section 3.3, this was achieved in an indirect way by performing a cross-validation experiment: the overall reliability of climate-informed predictions suggests that the assumptions of the probabilistic model are at least reasonable in this particular study. However, more specialized diagnostics would be valuable to isolate individual assumptions of the model and scrutinize them more thoroughly. Such assumptions include:

1. The assumption of conditional independence in space for the hydrologic data: while the existence of a common temporal pattern is likely to explain a part of the spatial dependence, some amount of dependence might remain even after conditioning on the temporal pattern.

2. The assumption that the temporal pattern is common to all sites (see discussion in next section).

3. The functional form assumed in equation (1b): such a simple linear relationship might not be optimal. While using an alternative nonlinear formulation is perfectly feasible, identifying the most relevant relationship requires developing dedicated diagnostics. For instance, data-driven analyses could explore the temporal variation of the spatial mean, the spatial variability of the time mean, the magnitude of space-time interactions, etc. Flexible modeling approaches such as *k*-nearest neighbors and local likelihood methods are also of interest [e.g., *Sankarasubramanian and Lall*, 2003; *Steinschneider and Brown*, 2012].

4. The geostatistical assumptions in equation (2c): in particular, identifying the most relevant correlogram function and developing anisotropic models are topics of particular interest. Indeed, such assumptions may have a nonnegligible impact on predictions since they directly act on the predicted temporal pattern through equation (6). They may impact the estimated uncertainty of the prediction if they are misspecified.

Scrutinizing all these assumptions require developing diagnostic tools tailored to the proposed two-level model: this will be explored in future work.

### 4.2. On the Common Temporal Pattern Assumption

The assumption that the temporal pattern is common to all sites requires working in a hydrologic region where one expects a common effect of climate. In the case study of section 3, the region was chosen based on expertise, and the successful cross-validation exercise suggests that this choice was reasonable. An interesting development would be to develop methods for automatically delineating such a region, whose size should be the result of a tradeoff: indeed, the region should be small enough to concur with the common temporal pattern assumption. On the other hand, a too small region only comprising a couple of sites might yield a temporal pattern mainly reflecting local peculiarities, whose predictability from climate might be reduced.

Note that this topic is somehow related to the notion of homogeneity used in standard RFA [e.g. *Burn*, 1997; *Ouarda et al.*, 2001, among many others]: we aim at delineating a "homogeneous region" with respect to the temporal pattern. But this notion of homogeneity fundamentally differs from the one used in standard RFA, where the objective is to select sites having the same time-invariant distribution (up to some scaling factor). Whether or not these two types of homogenous regions would broadly coincide remains an open question.

### 4.3. Prediction at Ungauged Sites

The local effects $\lambda(x)$ in equation (1b) are site-specific parameters. While this provides the flexibility to describe local peculiarities, this does not allow using the model for prediction at an ungauged site. This could be improved by regionalizing these parameters, through a regression with explanatory variables (e.g. catchment size, elevation, etc.) or a hierarchical model (describing the stochastic variations of $\lambda(x)$ in space) or a combination of them [*Micevski et al.*, 2006; *Cooley et al.*, 2007; *Lima and Lall*, 2009; *Renard*, 2011]. Such development would allow making climate-informed predictions at ungauged sites and, more generally, mapping climate effects.

### 4.4. Multicomponent Models

The key assumption in the proposed model is that all hydrologic sites are impacted by a common temporal pattern (equation (1b)), which itself is driven by one specific mode of climate variability (equation (2)). However multiple influences of distinct modes of climate variability have been largely reported in the literature [e.g. *Verdon-Kidd and Kiem*, 2009, in Australia]. In the context of the proposed model, this would call for considering several temporal patterns, yielding for instance the relationship: $\theta(x,t) = \lambda(x)\big(1 + \tau_1(t) + \tau_2(t)\big)$. Each temporal pattern could correspond to distinct components extracted from a unique climate variable (analogously to the successive components defined in PCA). Alternatively, they could correspond to unique components extracted from two distinct climate variables (e.g. geopotential heights and SST). The development of a multicomponent version of our model will be undertaken in future work.

### 4.5. Applications

The main objective of this paper was to describe the probabilistic model and to illustrate its application based on a case study. However, a detailed investigation of the benefit of the proposed model for operational applications lies beyond the scope of this paper and was therefore not addressed in any depth in the case study of section 3. Future work will investigate potential applications, in particular in terms of seasonal forecasting and downscaling. Regarding the former application, introducing a lag between climate and hydrologic data (which was not done in the case study of section 3) would enable making climate-conditional forecasts. Many examples of such analyses can be found in the literature [*Wedgbrow et al.*, 2002; *Sankarasubramanian and Lall*, 2003; *Wilby et al.*, 2004; *Grantz et al.*, 2005; *Westra et al.*, 2008; *Lima and Lall*, 2010, among many others]. Regarding the downscaling application, the relationship identified between climate and hydrology could be used to make hydrologic predictions conditionally on some future projection of climate variables [*Tisseuil et al.*, 2010; *Tramblay et al.*, 2012b; *Tramblay et al.*, 2012a]. Note that unlike seasonal forecasting, this application does not require introducing a lag between climate and hydrologic data.

## 5. Conclusions

The objective of this paper was to describe a time-varying regional frequency analysis framework, based on conditioning the parameters of at-site distributions on the values taken by some large-scale atmospheric or oceanic field. This method is based on a two-level probabilistic model, which identifies: (i) a common temporal pattern affecting a regional data set of hydrologic variables (level 1) and (ii) the spatial pattern in the climate data set associated with this temporal variability (level 2). This allows making a probabilistic prediction for the hydrologic variable conditionally on an observed atmospheric or oceanic field. This is to be compared with an approach making prediction conditionally on some predefined climate indices (e.g. SOI, NAO, etc.). Such indices, while explaining an important part of the internal climate variability, do not necessarily provide the best explanation for a given region and a given hydrologic variable. The proposed model allows identifying the most relevant mode of climate variability for the region/variable of interest. It also improves on an approach such as the Canonical Correlation Analysis by providing a mechanism for a full uncertainty analysis and the ability to work with non-Gaussian and Gaussian distributions for the predictand.

Inference is performed within a Bayesian-MCMC framework. While this has important advantages (regarding uncertainty quantification in particular), this also constitutes a significant computational challenge. Consequently, a computationally efficient exploratory tool is also proposed in this paper to perform preliminary analyses.

A case study illustrates the application of the proposed model for predicting the number of flood events in 16 catchments located in Mediterranean France using geopotential heights at 500 hPa over the North-Atlantic region. The model identifies a particular spatial pattern in the geopotential heights data (negative anomalies over the midlatitudes of the Eastern Atlantic, positive anomalies in Scandinavia) that is associated with the temporal variability of hydrologic data. Results from a cross-validation experiment indicate that the resulting climate-informed predictions are skillful, with an acceptable reliability and an improved sharpness compared to baseline predictions that ignore climate information, or to predictions based on predefined climate indices (EAWR, NAO, and SCAND).

Since our approach is based on an explicit probabilistic model, it inherits both advantages and drawbacks from such model-based approaches. An important advantage is the flexibility of the framework: the probabilistic assumptions made to derive its constitutive components can be modified if need be. For instance, nonlinear relationships between the hydrologic and climate components could be investigated. Similarly, an extension of the currently single-component model to a multicomponent setup is also a promising perspective. Moreover, the use of parametric probabilistic assumptions allows embedding the model within a Bayesian framework, which in turns yields a natural and built-in quantification of uncertainties. However, as any model-based approach, our model does make several restrictive assumptions regarding the distribution of data, the form of the hydrology-climate relationship, etc. Such assumptions should be individually scrutinized, which requires developing diagnostic tools tailored to the particular structure of the proposed framework. This is a topic of primary importance that will be explored in future work.

## Appendix A

### A1. Prediction: Maximum Likelihood Estimator of the Temporal Pattern

Let $(\alpha_{i,j})_{i,j=1:N_s}$ denote the elements of the matrix $(\hat{\sigma}^2 \Sigma_{\hat{\gamma}})^{-1}$. At prediction time $t^*$, the log-likelihood of the observed field $\phi(:,t^*)$ is equal to:

$$
L(\phi(:,t^*)|\tau(t^*)) = -\frac{N_s}{2} \log(2\pi) - \frac{1}{2} \log\left(\det\left(\hat{\sigma}^2 \Sigma_{\hat{\gamma}}\right)\right)
$$
$$
-\frac{1}{2} \sum_{i,j=1}^{N_s} \alpha_{i,j} \left(\phi(j,t^*) - \tau(t^*) * \hat{\psi}(j)\right)\left(\phi(i,t^*) - \tau(t^*) * \hat{\psi}(i)\right)
$$

(A1)

Deriving this quantity with respect to $\tau(t^*)$ yields:

$$
\frac{\partial L}{\partial \tau(t^*)} = -\frac{1}{2} \sum_{i,j=1}^{N_s} \alpha_{i,j} \left(-\hat{\psi}(j)\left(\phi(i,t^*) - \tau(t^*) * \hat{\psi}(i)\right) - \hat{\psi}(i)\left(\phi(j,t^*) - \tau(t^*) * \hat{\psi}(j)\right)\right)
$$

(A2)

The log-likelihood therefore has an extremum for:

$$
\frac{\partial L}{\partial \tau(t^*)} = 0
$$

$$
\iff \sum_{i,j=1}^{N_s} \alpha_{i,j} \left(\hat{\psi}(j)\left(\phi(i,t^*) - \tau(t^*) * \hat{\psi}(i)\right) + \hat{\psi}(i)\left(\phi(j,t^*) - \tau(t^*) * \hat{\psi}(j)\right)\right) = 0
$$

$$
\iff 2\sum_{i,j=1}^{N_s} \alpha_{i,j} \hat{\psi}(j) \phi(i,t^*) = 2\tau(t^*) \sum_{i,j=1}^{N_s} \alpha_{i,j} \hat{\psi}(j) \hat{\psi}(i)
$$

(A3)

$$
\iff \tau(t^*) = \frac{\sum_{i,j=1}^{N_s} \alpha_{i,j} \hat{\psi}(j) \phi(i,t^*)}{\sum_{i,j=1}^{N_s} \alpha_{i,j} \hat{\psi}(j) \hat{\psi}(i)}
$$

$$
\iff \tau(t^*) = \frac{\phi(:,t^*)\left(\hat{\sigma}^2 \Sigma_{\hat{\gamma}}\right)^{-1} \hat{\psi}^T}{\hat{\psi}\left(\hat{\sigma}^2 \Sigma_{\hat{\gamma}}\right)^{-1} \hat{\psi}^T}
$$

Moreover, the second derivative of the log-likelihood is equal to:

$$
\frac{\partial^2 L}{\partial \tau(t^*)} = -\sum_{i,j=1}^{N_s} \alpha_{i,j} \hat{\psi}(j) \hat{\psi}(i) = -\hat{\psi}\left(\hat{\sigma}^2 \Sigma_{\hat{\gamma}}\right)^{-1} \hat{\psi}^T
$$

(A4)

The second derivative is therefore negative because the matrix $\left(\hat{\sigma}^2 \Sigma_{\hat{\gamma}}\right)^{-1}$ is positive-definite: this proves that the value in equation (A3) maximizes the log-likelihood.

### A2. Exploratory Tool: Maximum Likelihood Estimators

#### A2.1. Temporal Pattern Extracted From Hydrologic Data

The log-likelihood of normal-score transformed data $\tilde{y}$ is equal to:

$$
L(\tilde{y}|\tau) = \sum_{x=1}^{N_x} \sum_{t=1}^{N_t} -\frac{1}{2} \log(2\pi) - \frac{1}{2}(\tilde{y}(x,t) - \tau(t))^2
$$

(A5)

For any time step $k$ $(=1,..,N_t)$, the root of the partial derivative with respect to $\tau(k)$ is:

$$
\frac{\partial L}{\partial \tau(k)} = 0 \iff \sum_{x=1}^{N_x} (\tilde{y}(x,k) - \tau(k)) = 0
$$

$$
\iff \tau(k) = \frac{\sum_{x=1}^{N_x} \tilde{y}(x,k)}{N_x} = mean(\tilde{y}(:,k))
$$

(A6)

Moreover, the second derivative of the log-likelihood is equal to $-N_x$ and is hence negative, which proves that the value in equation (A6) maximizes the log-likelihood.

### A2.2. Spatial Pattern

Conditionally on the estimated temporal pattern $\tilde{\tau}$, the log-likelihood of climate data $\phi$ is equal to:

$$L(\phi|\psi) = \sum_{t=1}^{N_t} \sum_{s=1}^{N_s} -\frac{1}{2}\log(2\pi) - \log(\sigma) - \frac{(\phi(s,t) - \psi(s) * \tilde{\tau}(t))^2}{2\sigma^2} \tag{A7}$$

For any gridpoint $k$ $(=1, .., N_s)$, the root of the partial derivative with respect to $\psi(k)$ is:

$$\frac{\partial L}{\partial \psi(k)} = 0 \iff \frac{1}{\sigma^2} \sum_{t=1}^{N_t} \tilde{\tau}(t)(\phi(k,t) - \psi(k) * \tilde{\tau}(t)) = 0$$

$$\iff \psi(k) = \frac{\sum_{t=1}^{N_t} \tilde{\tau}(t)\phi(k,t)}{\sum_{t=1}^{N_t} [\tilde{\tau}(t)]^2} \tag{A8}$$

Moreover, the second derivative of the log-likelihood is equal to $-\frac{1}{\sigma^2}\sum_{t=1}^{N_t}[\tilde{\tau}(t)]^2$ and is hence negative, which proves that the value in equation (A8) maximizes the log-likelihood.

### A2.3. Temporal Pattern Reconstructed From Climate Data

The demonstration is identical to that in Appendix A1, simply replacing the matrix $\Sigma_{\hat{\gamma}}$ by the identity matrix.

## References

Aryal, S. K., B. C. Bates, E. P. Campbell, Y. Li, M. J. Palmer, and N. R. Viney (2009), Characterizing and modeling temporal and spatial trends in rainfall extremes, *J. Hydrometeorol.*, *10*(1), 241–253.

Barnston, A. G., and R. E. Livezey (1987), Classification, seasonality and persistence of low-frequency atmospheric circulation patterns, *Mon. Weather Rev.*, *115*(6), 1083–1126.

Bretherton, C. S., C. Smith, and J. M. Wallace (1992), An intercomparison of methods for finding coupled patterns in climate data, *J. Clim.*, *5*(6), 541–560.

Bröcker, J. (2007), Increasing the reliability of reliability diagrams, *Weather Forecast.*, *22*(3), 651–661.

Burn, D. H. (1997), Catchment similarity for regional flood frequency analysis using seasonality measures, *J. Hydrol.*, *202*(1–4), 212–230.

Cai, W., P. van Rensch, T. Cowan, and A. Sullivan (2010), Asymmetry in ENSO teleconnection with regional rainfall, its multidecadal variability, and impact, *J. Clim.*, *23*(18), 4944–4955.

Chiles, J.-P., and P. Delfiner (1999), *Geostatistics: Modeling Spatial Uncertainty*, 720 pp., John Wiley, Hoboken, N. J.

Cooley, D., D. Nychka, and P. Naveau (2007), Bayesian spatial modeling of extreme precipitation return levels, *J. Am. Stat. Assoc.*, *102*(479), 824–840.

El Adlouni, S., T. B. M. J. Ouarda, X. Zhang, R. Roy, and B. Bobée (2007), Generalized maximum likelihood estimators for the nonstationary generalized extreme value model, *Water Resour. Res.*, *43*, W03410, doi:10.1029/2005WR004545.

Gelfand, A. E., D. K. Dey, and H. Chang (1992), Model determination using predictive distributions with implementation via sampling-based-methods, paper presented at Bayesian Statistics 4, Oxford Univ. Press, Valencia, Spain, April 15–20, 1991.

Gelman, A., and D. B. Rubin (1992), Inference from iterative simulation using multiple sequences, *Stat. Sci.*, *7*(4), 457–472.

Gershunov, A., and T. P. Barnett (1998), ENSO influence on intraseasonal extreme rainfall and temperature frequencies in the contiguous United States: Observations and model results, *J. Clim.*, *11*(7), 1575–1586.

Giuntoli, I., B. Renard, J. P. Vidal, and A. Bard (2013), Low flows in France and their relationship to large scale climate indices, *J. Hydrol.*, *482*, 105–118.

Gneiting, T., F. Balabdaoui, and A. E. Raftery (2007), Probabilistic forecasts, calibration and sharpness, *J. R. Stat. Soc. Ser. B*, *69*, 243–268.

Grantz, K., B. Rajagopalan, M. Clark, and E. Zagona (2005), A technique for incorporating large-scale climate information in basin-scale ensemble streamflow forecasts, *Water Resour. Res.*, *41*, W10410, doi: 10.1029/2004WR003467.

Grimm, A. M., and R. G. Tedeschi (2009), ENSO and extreme rainfall events in South America, *J. Clim.*, *22*(7), 1589–1609.

Haario, H., E. Saksman, and J. Tamminen (2005), Componentwise adaptation for high dimensional MCMC, *Comput. Stat.*, *20*(2), 265–273.

Hoffman, M. D., and A. Gelman (2014), The No-U-Turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo, *J. Mach. Learn. Res.*, *15*, 1593–1623.

Hurrell, J. W. (1995), Decadal trends in the north-atlantic oscillation: Regional temperatures and precipitation, *Science*, *269*(5224), 676–679.

Hurrell, J. W., and H. VanLoon (1997), Decadal variations in climate associated with the north Atlantic oscillation, *Clim. Change*, *36*(3-4), 301–326.

Ionita, M. (2014), The impact of the East Atlantic/Western Russia pattern on the hydroclimatology of Europe from mid-einter to late spring, *Climate*, *2*(4), 296–309.

Kalnay, E., et al. (1996), The NCEP/NCAR 40-year reanalysis project, *Bull. Am. Meteorol. Soc.*, *77*(3), 437–471.

Kass, R. E., and A. E. Raftery (1995), Bayes Factors, *J. Am. Stat. Assoc.*, *90*(430), 773–795.

Kiely, G. (1999), Climate change in Ireland from precipitation and streamflow observations, *Adv. Water Resour.*, *23*(2), 141–151.

Kingston, D. G., D. M. Hannah, D. M. Lawler, and G. R. McGregor (2006), Interactions between large-scale climate and river flow across the northern North Atlantic margin, *IAHS AISH Publ.*, *308*, 350–355.

Kousky, V. E., M. T. Kagano, and I. F. A. Cavalcanti (1984), A review of the Southern Oscillation: Oceanic-atmospheric circulation changes and related rainfall anomalies, *Tellus Ser. A*, *36A*(5), 490–504.

Kripalani, R. H., and A. Kulkarni (1997), Rainfall variability over South-East Asia—connections with Indian monsoon and ENSO extremes: New perspectives, *Int. J. Climatol.*, *17*(11), 1155–1168.

Laloy, E., and J. A. Vrugt (2012), High-dimensional posterior exploration of hydrologic models using multiple-try DREAM(ZS) and high-performance computing, *Water Resour. Res.*, *48*, W01526, doi:10.1029/2011WR010608.

Lang, M., T. B. M. J. Ouarda, and B. Bobée (1999), Towards operational guidelines for over-threshold modeling, *J. Hydrol.*, *225*, 103–117.

Lavers, D., C. Prudhomme, and D. M. Hannah (2013), European precipitation connections with large-scale mean sea-level pressure (MSLP) fields, *Hydrol. Sci. J.*, *58*(2), 310–327.

Lima, C. H. R., and U. Lall (2009), Hierarchical Bayesian modeling of multisite daily rainfall occurrence: Rainy season onset, peak, and end, *Water Resour. Res.*, *45*, W07422, doi:10.1029/2008WR007485.

Lima, C. H. R., and U. Lall (2010), Climate informed long term seasonal forecasts of hydroenergy inflow for the Brazilian hydropower system, *J. Hydrol.*, *381*(1-2), 65–75.

McCullagh, P., and J. A. Nelder (1989), *Generalized Linear Models*, 2nd ed., 532 pp., Chapman and Hall, London, U. K.

Micevski, T., G. Kuczera, and S. W. Franks (2006), A Bayesian hierarchical regional flood model, in *Proceedings of the 30th Hydrology and Water Resources Symposium*, edited by E. Australia, Launceston, Engineers Australia, Tasmania, Australia, 4–7 December.

Ouarda, T. B. M. J., C. Girard, G. Cavadias, and B. Bobee (2001), Regional flood frequency estimation with canonical correlation analysis, *J. Hydrol.*, *254*, 157–173.

Philippon, N., M. Rouault, Y. Richard, and A. Favre (2012), The influence of ENSO on winter rainfall in South Africa, *Int. J. Climatol.*, *32*(15), 2333–2347.

Pociask-Karteczka, J. (2006), River hydrology and the North Atlantic oscillation: A general review, *Ambio*, *35*(6), 312–314.

Renard, B. (2011), A Bayesian hierarchical approach to regional frequency analysis, *Water Resour. Res.*, *47*, W11513, doi:10.1029/2010WR010089.

Renard, B., V. Garreta, and M. Lang (2006), An application of Bayesian analysis and MCMC methods to the estimation of a regional trend in annual maxima, *Water Resour. Res.*, *42*, W12422, doi:10.1029/2005WR004591.

Sankarasubramanian, A., and U. Lall (2003), Flood quantiles in a changing climate: Seasonal forecasts and causal relations, *Water Resour. Res.*, *39*(5), 1134, doi:10.1029/2002WR001593.

Shaman, J., and E. Tziperman (2010), An atmospheric teleconnection linking ENSO and southwestern European precipitation, *J. Clim.*, *24*(1), 124–139.

Shang, H. W., J. Yan, and X. B. Zhang (2011), El Nino-Southern Oscillation influence on winter maximum daily precipitation in California in a spatial model, *Water Resour. Res.*, *47*, W11507, doi:10.1029/2011WR010415.

Shorthouse, C., and N. W. Arnell (1999), The effects of climatic variability on spatial characteristics of European river flows., *Phys. Chem. Earth. Part B*, *24*, 7–13.

Stahl, K., S. Demuth, H. Hisdal, M. J. Santos, R. Veríssimo, and R. Rodrigues (2001), The North Atlantic Oscillation (NAO) and the drought, In *Assessment of the Regional Impact of Droughts in Europe*, final report, ARIDE, Inst. of Hydrol., Freiburg.

Steinschneider, S., and C. Brown (2012), Forecast-informed low-flow frequency analysis in a Bayesian framework for the northeastern United States, *Water Resour. Res.*, *48*, W10545, doi:10.1029/2012WR011860.

Sun, X., M. Thyer, B. Renard, and M. Lang (2014), A general regional frequency analysis framework for quantifying local-scale climate effects: A case study of ENSO effects on Southeast Queensland rainfall, *J. Hydrol.*, *512*, 53–68.

Tipping, M. E., and C. M. Bishop (1999), Probabilistic principal component analysis, *J. R. Stat. Soc. Ser. B*, *61*(3), 611–622.

Tisseuil, C., M. Vrac, S. Lek, and A. J. Wade (2010), Statistical downscaling of river flows, *J. Hydrol.*, *385*(1–4), 279–291.

Tramblay, Y., L. Neppel, J. Carreau, and E. Sanchez-Gomez (2012a), Extreme value modelling of daily areal rainfall over Mediterranean catchments in a changing climate, *Hydrol. Processes*, *26*(25), 3934–3944.

Tramblay, Y., W. Badi, F. Driouech, S. El Adlouni, L. Neppel, and E. Servat (2012b), Climate change impacts on extreme precipitation in Morocco, *Global Planet. Change*, *82–83*, 104–114.

Trigo, R. M., D. Pozo-Vazquez, T. J. Osborn, Y. Castro-Diez, S. Gamiz-Fortis, and M. J. Esteban-Parra (2004), North Atlantic oscillation influence on precipitation, river flow and water resources in the Iberian peninsula, *Int. J. Climatol.*, *24*(8), 925–944.

Uppala, S. M., et al. (2005), The ERA-40 re-analysis, *Q. J. R. Meteorol. Soc.*, *131*(612), 2961–3012.

van Loon, H., and J. C. Rogers (1978), The Seesaw in Winter Temperatures between Greenland and Northern Europe. Part I: General Description, *Mon. Weather Rev.*, *106*(3), 296–310.

Verdon-Kidd, D. C., and A. S. Kiem (2009), On the relationship between large-scale climate modes and regional synoptic patterns that drive Victorian rainfall, *Hydrol. Earth Syst. Sci.*, *13*(4), 467–479.

Vicente-Serrano, S., and J. Cuadrat (2007), North Atlantic oscillation control of droughts in north-east Spain: Evaluation since 1600 a.d, *Clim. Change*, *85*(3-4), 357–379.

Wedgbrow, C. S., R. L. Wilby, H. R. Fox, and G. O'Hare (2002), Prospects for seasonal forecasting of summer drought and low river flow anomalies in England and Wales, *Int. J. Climatol.*, *22*(2), 219–236.

Westra, S., and A. Sharma (2009), Probabilistic estimation of multivariate streamflow using independent component analysis and climate information, *J. Hydrometeorol.*, *10*(6), 1479–1492.

Westra, S., A. Sharma, C. Brown, and U. Lall (2008), Multivariate streamflow forecasting using independent component analysis, *Water Resour. Res.*, *44*, W02437, doi:10.1029/2007WR006104.

Wibig, J. (1999), Precipitation in Europe in relation to circulation patterns at the 500 hPa level, *Int. J. Climatol.*, *19*(3), 253–269.

Wilby, R. L., G. O'Hare, and N. Barnsley (1997), The North Atlantic Oscillation and British Isles climate variability, 1865–1996, *Weather*, *52*, 266–276.

Wilby, R. L., C. S. Wedgbrow, and H. R. Fox (2004), Seasonal predictability of the summer hydrometeorology of the River Thames, UK, *J. Hydrol.*, *295*(1-4), 1–16.