

Statistics of extreme ocean environments: Non-stationary inference for directionality and other covariate effects



Matthew Jones^a, David Randell^b, Kevin Ewans^c, Philip Jonathan^{b,*}

^a Department of Mathematical Sciences, Durham University, Durham DH1 3LE, United Kingdom

^b Shell Projects & Technology, Manchester M22 0RR, United Kingdom

^c Sarawak Shell Bhd., 50450 Kuala Lumpur, Malaysia

ARTICLE INFO

Article history:

Received 15 December 2015

Received in revised form

16 February 2016

Accepted 10 April 2016

Available online 29 April 2016

Keywords:

Extreme

Covariate

Non-stationary

Smoothing

Non-parametric

Spline

Gaussian process

mMALA

Kullback–Leibler

ABSTRACT

Numerous approaches are proposed in the literature for non-stationarity marginal extreme value inference, including different model parameterisations with respect to covariate, and different inference schemes. The objective of this paper is to compare some of these procedures critically. We generate sample realisations from generalised Pareto distributions, the parameters of which are smooth functions of a single smooth periodic covariate, specified to reflect the characteristics of actual samples from the tail of the distribution of significant wave height with direction, considered in the literature in the recent past. We estimate extreme values models (a) using Constant, Fourier, B-spline and Gaussian Process parameterisations for the functional forms of generalised Pareto shape and (adjusted) scale with respect to covariate and (b) maximum likelihood and Bayesian inference procedures. We evaluate the relative quality of inferences by estimating return value distributions for the response corresponding to a time period of $10 \times$ the (assumed) period of the original sample, and compare estimated return values distributions with the truth using Kullback–Leibler, Cramer–von Mises and Kolmogorov–Smirnov statistics. We find that Spline and Gaussian Process parameterisations, estimated by Markov chain Monte Carlo inference using the mMALA algorithm, perform equally well in terms of quality of inference and computational efficiency, and generally perform better than alternatives in those respects.

© 2016 Elsevier Ltd. All rights reserved.

1. Introduction

Accurate estimates of the likely extreme environmental loading on an offshore facility are vital to enable a design that ensures the facility is both structurally reliable and economic. This involves estimating the extreme value behaviour of meteorological and oceanographic (metocean) variables that quantify the various environmental loading quantities, primarily winds, wave, and currents. Examples of such parameters are significant wave height, mean wind speed and mean current speed. These characterise the environment for a given short period of time within which the environment is assumed to be stationary.

The long-term variability of these parameters is however non-stationary, in particular with respect to time, space and direction. From a temporal point of view metocean parameters generally have a strong seasonal variation, with an annual periodicity, and longer term variations due to decadal or semi-decadal climate variations. At any given location, the variability of a particular

parameter is also dependent on the direction; for example, wind forcing is typically stronger from some directions than others, and fetch and water depth effects can strongly influence the resulting magnitude of the waves. Clearly these effects will vary with location: a more exposed location will be associated with longer fetches, resulting in a more extreme wave climate.

When estimating the long-term variability of parameters, such as significant wave height, the non-stationary effects associated with e.g. direction and season can be incorporated by treating direction and season as covariates. The common practice is to perform extreme value analysis of hindcast data sets, which include many years of metocean parameters, along with their associated covariates. Such data sets have all the information needed for input to covariate analysis.

From a design perspective, the metocean engineer is often required to specify return values for directional sectors such as octants centred on the cardinal and semi-cardinal directions. These directional return value estimates must be consistent with the estimated omnidirectional return value. In a similar manner, return values may be required corresponding to particular seasons or months of the year, consistent with an all-year return value.

* Corresponding author.

E-mail address: philip.jonathan@shell.com (P. Jonathan).

Clearly, therefore, efficient and reliable inference for non-stationary extremes is of considerable practical interest, requiring estimation of (a) the rate and (b) the size of rare events. This work addresses the latter of these objectives.

A non-stationary extreme value model is generally superior to the alternative “partitioning” method sometimes used within the ocean engineering community. In the partitioning method, the sample is partitioned into subsets corresponding to approximately constant values of covariate(s); independent extreme value analysis is then performed on each subset. For example, in the current work we might choose to partition the sample into directional octants, and then estimate (8 independent stationary) extreme value models for each of the octants. There are two main reasons for favouring a non-stationarity model over the partitioning method. Firstly, the partitioning approach incurs a loss in statistical efficiency of estimation, since parameter estimates for subsets with similar covariate values are estimated independently of one another, even though physical insight would require parameter estimates to be similar. This problem worsens as the number of covariates and covariate subsets increases, and the sample size per subset decreases as a result. In the non-stationary model, we require that parameter estimates corresponding to similar values of covariates be similar, and optimise the degree of similarity during inference. For this reason, parameter uncertainty from the non-stationary model is generally smaller than from the partitioning approach. Secondly, the partitioning approach assumes that, within each subset, the sub-sample for extreme value modelling is homogeneous with respect to covariates. In general it is difficult to estimate what effect this assumption might have on parameter and return value estimates (especially when large intervals of values of covariates are combined into a subset). In the non-stationary model, we avoid the need to make this assumption.

Numerous articles have reported the essential features of extreme value analysis (e.g. [Davison and Smith, 1990](#)) and the importance of considering different aspects of covariate effects (e.g. [Northrop et al., 2016](#)). [Carter and Challenor \(1981\)](#) consider estimation of annual maxima from monthly data, when the distribution functions of monthly extremes are known. [Coles and Walshaw \(1994\)](#) describe directional modelling of extreme wind speeds using a Fourier parameterisation. [Scotto and Guedes-Soares \(2000\)](#) model the long-term time series of significant wave height with non-linear threshold models. [Anderson et al. \(2001\)](#) report that estimates for 100-year significant wave height from an extreme value model ignoring seasonality are considerably smaller than those obtained using a number of different seasonal extreme value models. [Chavez-Demoulin and Embrechts \(2006\)](#) describe smooth extreme value models in finance and insurance. [Chavez-Demoulin and Davison \(2005\)](#) provide a straight-forward description of a nonhomogeneous Poisson model in which occurrence rates and extreme value properties are modelled as functions of covariates. [Cooley et al. \(2006\)](#) use a Bayesian hierarchical model to characterise extremes of lichen growth. [Renard et al. \(2006\)](#) consider identification of changes in peaks over threshold using Bayesian inference. [Fawcett and Walshaw \(2006\)](#) use a hierarchical model to identify location and seasonal effects in marginal densities of hourly maxima for wind speed. [Mendez et al. \(2008\)](#) consider seasonal non-stationarity in extremes of NOAA buoy records. [Randell et al. \(2015a\)](#) discuss estimation for return values for significant wave height in the South China Sea using a directional-seasonal extreme value model. [Randell et al. \(2014\)](#) explore the directional characteristics of hindcast storm peak significant wave height with direction for locations in the Gulf of Mexico, North-West Shelf of Australia, Northern North Sea, Southern North Sea, South Atlantic Ocean, Alaska, South China Sea and West Africa. [Fig. 1](#) illustrates the essential features of samples such as these. The rate and magnitude of occurrences of storm

events vary considerably between locations, and with direction at each location. There are directional sectors with effectively no occurrences, there is evidence of rapid changes in characteristics with direction and of local stationarity with direction. Any realistic model for such samples needs to be non-stationary with respect to direction.

The objective of this paper is to evaluate critically different procedures for estimating non-stationary extreme value models. We quantify the extent to which extreme value analysis of samples of peaks over threshold exhibiting clear non-stationarity with respect to covariates, such as those in [Fig. 1](#) or simulation case studies in [Section 3](#) below, is influenced by a particular choice of model parameterisation or inference method. The 6 simulation case studies introduced in [Section 3](#) are constructed to reflect the general features of the samples in [Fig. 1](#), with the advantage that the statistical characteristics of the case studies are known exactly, allowing objective evaluation and comparison of competing methods of model parameterisation and inference. Our aim is that the results of this study are generally informative about any application of non-stationary extreme value analysis. We generate sample realisations from generalised Pareto distributions, the parameters of which are smooth functions of a single smooth periodic covariate. Then we estimate extreme value models (a) using Constant, Fourier, B-spline and Gaussian Process parameterisations for the functional forms of generalised Pareto parameters with respect to covariate and (b) maximum likelihood and Bayesian inference procedures. We evaluate the relative quality of inferences by estimating return value distributions for the response corresponding to a time period of $10 \times$ the (assumed) period of the original sample, and compare estimated return values distributions with the truth using Kullback–Leibler (e.g. [Perez-Cruz, 2008](#)), Cramer–von Mises (e.g. [Anderson, 1962](#)) and Kolmogorov–Smirnov statistics. We cannot hope to compare all possible parameterisations, but choose four parameterisations useful in our experience. Similarly, there are many competing approaches for maximum likelihood and Bayesian inference, and general interest in understanding their relative characteristics. For example, [Smith and Naylor \(1987\)](#) compare maximum likelihood and Bayesian inference for the three-parameter Weibull distribution. In this work, we choose to compare frequentist penalised likelihood maximisation (see [Section 2.3](#)) with two Markov chain Monte Carlo (MCMC) methods of different complexities. Non-stationary model estimation is a growing field. There is a huge literature on still further possibilities for parametric (e.g. Chebyshev, Legendre and other polynomial forms) and non-parametric (e.g. Gauss–Markov random fields and radial basis functions) model parameterisations with respect to covariates. Moreover, in extreme value analysis, pre-processing of a response to near stationarity (e.g. using a Box–Cox transformation) is preferred.

The outline of the paper is as follows. [Section 2](#) outlines the different model parameterisations and inference schemes under consideration. [Section 3](#) describes underlying model forms used to generate samples for inference, outlines the procedure for estimation of return value distributions and their comparison, and presents results of those comparisons. [Section 4](#) provides discussion and conclusions.

2. Estimating non-stationary extremes

Consider a random variable Y representing an environmental variable of interest such as significant wave height. The characteristics of Y are dependent on covariates such as (wave) direction, season, location and fetch. In this work we assume that a single periodic covariate θ (typically direction, or season) is sufficient to characterise the non-stationarity of Y . That is, we assume

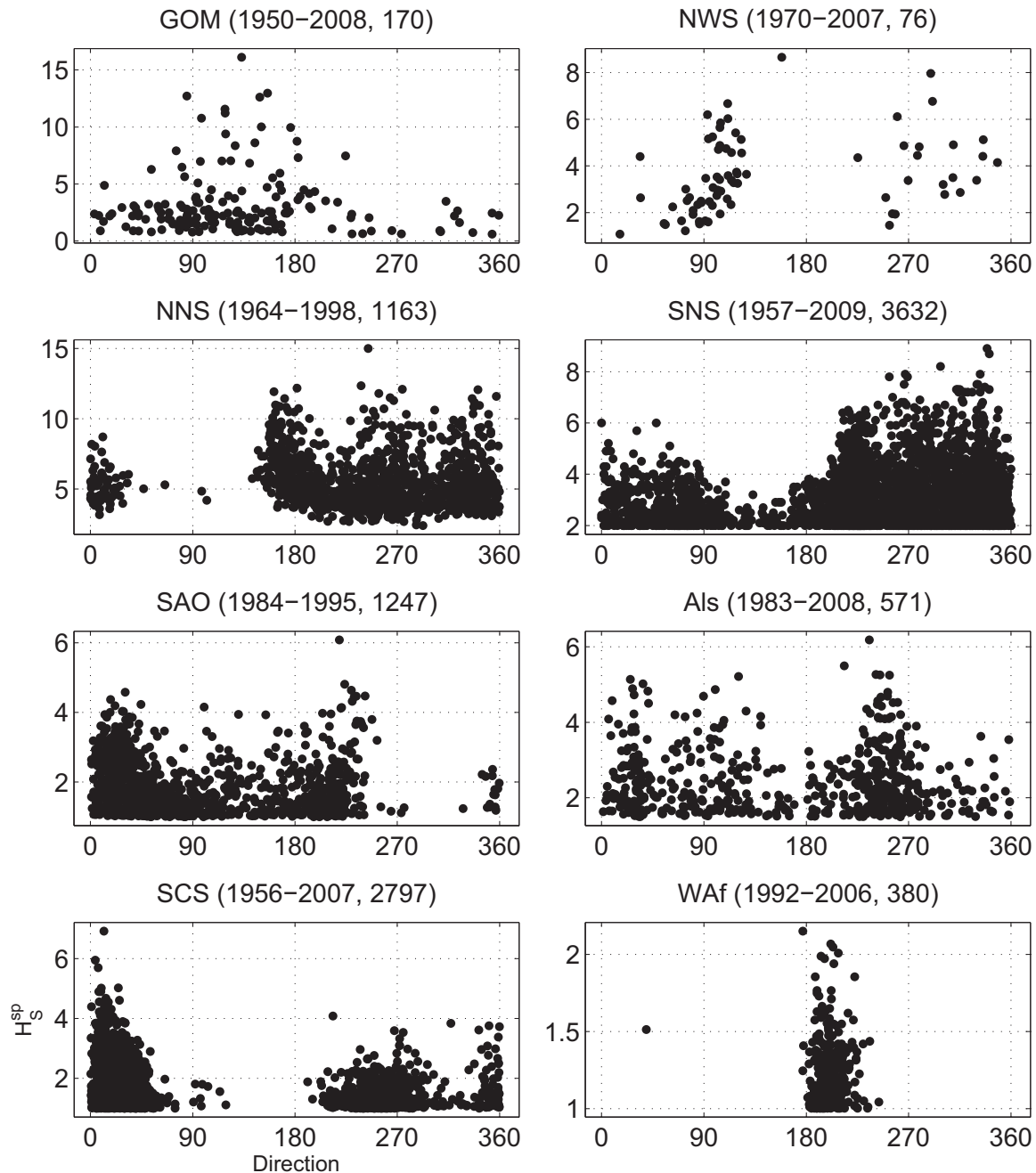


Fig. 1. Hindcast storm peak significant wave height on direction for 8 locations worldwide. From right to left, top to bottom: Gulf of Mexico (GOM), North-West Shelf of Australia (NWS), Northern North Sea (NNS), Southern North Sea (SNS), South Atlantic Ocean (SAO), Alaska (Als), South China Sea (SCS) and West Africa (Waf). Panel titles give the location, the sample period and storm peak sample size. Refer to [Randell et al. \(2014\)](#) for details of data sources.

that $Y|\theta$ has a stationary distribution. For exceedances $Y - \mu(\theta)$ of some high threshold $\mu(\theta)$, extreme value theory suggest that the conditional distribution of $Y - \mu(\theta)$ given that $Y > \mu(\theta)$ can be approximated by the generalised Pareto distribution

$$\Pr(Y > y | Y > \mu(\theta)) = \frac{1}{\sigma(\theta)} \left(1 + \frac{\xi(\theta)}{\sigma(\theta)} (y - \mu(\theta)) \right)^{-1/\xi(\theta)}.$$

For design purposes, return value distributions are typically estimated using peaks-over-threshold of significant wave height. The characteristics of these peaks, each corresponding to a different storm event, vary typically with respect to wave direction. Conditional on wave direction, the peaks are reasonably assumed

to be independent of one another. To evaluate the relative performance of different approaches to non-stationary extreme value analysis, it is therefore natural to use random simulation from generalised Pareto models with known directional characteristics.

2.1. Generalised pareto model

We assume we observe a sample $y = \{y_1, \dots, y_N\}$ of peaks over threshold drawn independently from a generalised Pareto distribution, the parameters of which are functions of corresponding observed covariate values $\theta = \{\theta_1, \dots, \theta_N\}$. The sample likelihood is a product of generalised Pareto (GP) likelihoods for each of the observations

$$f(y|\theta, \xi, \sigma, \mu) = \prod_{i=1}^N f(y_i|\xi(\theta_i), \sigma(\theta_i), \mu(\theta_i))$$

$$= \prod_{i=1}^N \frac{1}{\sigma(\theta_i)} \left(1 + \xi(\theta_i) \frac{(y_i - \mu(\theta_i))}{\sigma(\theta_i)} \right)^{-1/\xi(\theta_i)-1}$$

where $\xi(\theta)$ and $\sigma(\theta)$ are the shape and scale parameters as functions of covariate. We do not attempt to estimate the threshold function $\mu(\theta)$, assuming it is 0 for all covariate values. We also assume that the rate of occurrence $\rho(\theta)$ of exceedances of μ varies with covariate, but that $\rho(\theta)$ is known (see Section 3.1). It is computationally advantageous (e.g. Cox and Reid, 1987, Chavez-Demoulin and Davison, 2005) to transform variables from (ξ, σ) to the asymptotically independent pair (ξ, ν) , where $\nu(\theta) = \sigma(\theta)(1 + \xi(\theta))$. Inference therefore amounts to estimating the smooth functions $\xi(\theta)$ and $\nu(\theta)$, although we usually choose to illustrate the analysis in terms of $\xi(\theta)$ and $\sigma(\theta)$. In practical application, estimation of $\mu(\theta)$ is itself generally also problematic (e.g. Scarrott and MacDonald, 2012), particularly in the presence of non-stationarity (Northrop and Jonathan, 2011), but as necessary for inference as reliable estimation of $\xi(\theta)$ and $\sigma(\theta)$. We choose to focus on the latter in this work.

2.2. Covariate parameterisations

To accommodate non-stationarity, we parameterise ξ and ν as linear combinations of unknown parameters β_ξ and β_ν respectively, where

$$\nu(\theta) = B_\nu(\theta)\beta_\nu, \quad \text{and} \quad \xi(\theta) = B_\xi(\theta)\beta_\xi$$

and $B_\nu(\theta)$ and $B_\xi(\theta)$ are row vectors of basis functions evaluated at θ . We consider four different forms of basis function, corresponding to Constant (stationary), Fourier, Spline and Gaussian Process parameterisations for $\xi(\theta)$ and $\nu(\theta)$, as described below.

Physical considerations suggest that we should expect GP model parameters to vary smoothly with covariate. In general, the basis parameterisation introduced here permits estimation of functional forms for GP model parameters which are too variable (or too rough) with respect to covariate. We therefore need a mechanism to restrict the roughness of functional forms, such that their roughness is optimal given the evidence in the data. For each parameterisation, we therefore specify roughness matrices Q_η (for $\eta = \xi, \nu$) to regulate the roughness of $\eta(\theta)$ with respect to θ during inference. This ensures that the elements of β_η weight the individual basis functions in such a way that the resulting estimate is optimally smooth in some sense. The form of the roughness penalty term R_η is $\frac{1}{2}\lambda_\eta\beta_\eta'Q_\eta\beta_\eta$, for some roughness coefficient λ_η . The penalty is incorporated directly within a penalised likelihood for maximum likelihood inference, and within a prior distribution for β_η in Bayesian inference, as described in Section 2.3.

2.2.1. Constant (stationary) parameterisation

In the Constant parameterisation, the values of ξ and ν do not vary with respect to θ . We therefore adopt a scalar basis function which is constant across all values of covariate, so that $B_\nu(\theta) = B_\xi(\theta) = 1$, and corresponding roughness matrices $Q_\nu = Q_\xi = 1$. We do not expect return value distributions estimated under this parameterisation to fare well in general in our comparison, since samples are generated from non-stationary distributions. Quality of fit is expected to be poor, at least in some intervals of covariate. However, many practitioners continue to use stationary extreme value models, perhaps with high thresholds to mitigate non-stationarity, in applications; inclusion of a stationary parameterisation provides a useful point of reference for comparison, therefore.

2.2.2. Spline parameterisation

Under a Spline parameterisation, the vector of basis functions for each of ν and ξ is made up of p local polynomial B-spline functions with compact support, joined at a series of knots evenly spaced in the covariate domain (e.g. Eilers and Marx, 2010)

$$B_\nu(\theta) = B_\xi(\theta) = (b_1(\theta) \cdots b_p(\theta)).$$

We specify roughness matrices $Q_\nu = Q_\xi = D^T D$ which penalise squared differences between adjacent elements of the coefficient vectors, where

$$D = \begin{pmatrix} -1 & 1 & 0 & \cdots & 0 \\ 0 & -1 & 1 & & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

is a $(p-1) \times p$ difference matrix. In this work we set p to 50.

2.2.3. Fourier parameterisation

We use basis vectors composed of sine and cosine functions of n_p different periods

$$B_\nu(\theta) = B_\xi(\theta) = (1 \sin(\theta) \sin(2\theta) \cdots \sin(n_p\theta) \cos(\theta) \cdots \cos(n_p\theta)).$$

The roughness matrix is computed by imposing a condition on the squared second derivative of the resulting parameter function. If we write

$$\eta(\theta) = \sum_{k=1}^{n_p} (a_{\eta k} \cos(k\theta) + b_{\eta k} \sin(k\theta))$$

where $\eta = \xi$ or ν , and $a_{\eta k}$ and $b_{\eta k}$ are the parameters from β_η corresponding respectively to the sine and cosine functions of period k . The roughness criterion (from Jonathan et al., 2013) becomes

$$R_\eta = \int_0^{2\pi} (\eta''(\theta))^2 d\theta = \sum_{k=1}^{n_p} k^4 (a_{\eta k}^2 + b_{\eta k}^2)$$

such that the penalty matrix can be written in matrix form as

$$Q_\eta = \text{diag}(0, 1, 2^4, \dots, k^4, \dots, n_p^4, 1, 2^4, \dots, k^4, \dots, n_p^4)$$

for the $p = 2n_p + 1$ Fourier parameters $(a_{\eta 0}, a_{\eta 1}, \dots, a_{\eta n_p}, b_{\eta 1}, \dots, b_{\eta n_p})$. In this work, we set the value of n_p to 25 so that $p = 51$.

2.2.4. Gaussian process parameterisation

For a set of p nodes $\{\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_p\}$ on the covariate domain, we use a Gaussian Process parameterisation (Rasmussen and Williams, 2006), and relate each covariate input to a knot using the following basis vectors

$$B_\nu(\theta) = B_\xi(\theta) = (I_1(\theta), \dots, I_p(\theta))$$

where the indicator functions $I_j(\cdot)$ are defined as

$$I_j(\theta) = \begin{cases} 1 & \text{if } |\theta - \hat{\theta}_j| < |\theta - \hat{\theta}_k| \forall k \neq j \\ 0 & \text{otherwise.} \end{cases}$$

Roughness matrices Q_η (where $\eta = \nu, \xi$) are defined by the coefficient correlation matrix V_η via $Q_\eta = V_\eta^{-1}$, and the elements of V_η generated by a periodic squared exponential covariance function (MacKay, 1998)

$$V_{\eta jk} = \exp \left(-\frac{2}{r_\eta^2} \sin \left(\frac{\hat{\theta}_j - \hat{\theta}_k}{2} \right)^2 \right)$$

where r_η are correlation lengths for each of the parameters, fixed to likely values by comparison with the covariate functions used to

generate the data; a value of r_η for 0.6 was used throughout. V_η^{-1} penalises on the angular difference between the j th and k th nodes, reducing in value from $\exp(0)$ at angular difference zero to $\exp(-2/r_\eta^2)$ at angular difference 180° . Estimating the Gaussian Process parameterisation on the partitioned covariate domain, as opposed to fitting it to each of the data inputs, greatly reduces the number of parameters to estimate, and is physically reasonable provided that p is sufficiently large. In this work, we use $p=50$ equally spaced nodes. Estimating a parameter for each data point would have made the computational burden for the Gaussian Process parameterisation significantly greater than that for any of the other parameterisations. For example, computations with roughness matrix Q_η (of dimension $p \times p$ on the partitioned covariate domain) are more efficient than those using the $N \times N$ version of Q_η ($N \approx 1000$) defined per data input.

2.2.5. Model complexity

We assume it is known from physical considerations that the extremal characteristics of the environmental variable of interest (e.g. significant wave height) vary smoothly with directional covariate. Specifically, we expect the form for the tail of the distribution to be homogeneous within a narrow directional sector of width $\approx 10^\circ$. This in turn suggests a suitable minimum complexity for the non-stationary model parameterisations considered in this work. For the Spline parameterisation, we set $p=50$ corresponding to 50 spline basis functions equally spaced on $[0, 360)$, with a distance between peaks of adjacent spline basis functions of 7.2° . We estimate the Gaussian Process on a regular partition of the covariate domain into $p=50$ bins; the distance between the centres of adjacent bins is again 7.2° . For the Fourier parameterisation, we use a Fourier order of $n_p=25$ (and $p=51$), such that half a wavelength of the highest frequency Fourier component again corresponds to 7.2° . In this way, we expect the directional resolution of these three model parameterisations to be comparable. Similarly, the number p of basis coefficients to be estimated in each of the three parameterisations is comparable. Therefore, we hope to focus fairly in the analysis below on the different inferential challenges presented by Spline, Fourier and Gaussian Process parameterisations and different estimation schemes for problems of comparable complexity. The stationary Constant parameterisation is clearly less complex (with $p=1$), but a useful baseline for comparison: we expect inferences from the Constant parameterisation to be relatively poor, and therefore make obvious the need to consider non-stationarity.

2.3. Inference procedures

We consider two methods for estimating parameters and return value distributions for the models and parameterisations described above, namely (a) maximum penalised likelihood estimation with bootstrapping to quantify uncertainties, and (b) (two forms of) Bayesian inference using Markov Chain Monte Carlo (MCMC). These are discussed below.

2.3.1. Maximum likelihood estimation

We use an iterative back-fitting optimisation (see Appendix) to minimise the penalised negative log likelihood $-L^*(y|\beta_\xi, \beta_\nu; \lambda_\xi, \lambda_\nu)$ with respect to β_ξ and β_ν for given roughness coefficients λ_ξ and λ_ν , where

$$\begin{aligned} -L^*(y|\beta_\xi, \beta_\nu; \lambda_\xi, \lambda_\nu) &= -L(y|\beta_\xi, \beta_\nu) + R_\xi + R_\nu \\ &= -L(y|\beta_\xi, \beta_\nu) + \frac{1}{2}\lambda_\xi\beta_\xi'Q_\xi\beta_\xi + \frac{1}{2}\lambda_\nu\beta_\nu'Q_\nu\beta_\nu. \end{aligned}$$

Here, $-L(y|\beta_\xi, \beta_\nu)$ is the negative log sample GP likelihood from Section 2.1 expressed as a function of ξ and ν , and R_ξ and R_ν are

additive roughness penalties. The values of λ_ξ and λ_ν are selected using cross-validation to maximise the predictive performance of the estimated model, and bootstrap resampling is used to quantify the uncertainty of parameter estimates. The original sample is resampled with replacement a large number of times, and inference repeated for each resample. We use the empirical distributions of parameter estimates and return values over resamples as approximate uncertainty distributions.

We refer readers interested in further information on penalised likelihood methods to the work of Green and Silverman (1994), Davison (2003), Ruppert et al. (2003), Eilers and Marx (2010), outlined in applications to metocean by Jonathan and Ewans (2013).

2.3.2. Bayesian inference

From a Bayesian perspective, all of β_ξ , β_ν , λ_ξ and λ_ν are treated as parameters to be estimated. Their joint posterior distribution given sample responses y and covariates θ can be written

$$\begin{aligned} f(\beta_\xi, \beta_\nu, \lambda_\xi, \lambda_\nu | y, \theta) &\propto f(y|\theta, \xi, \sigma, \mu) f(\beta_\nu|\lambda_\nu) f(\beta_\xi|\lambda_\xi) f(\lambda_\nu|a_\nu, b_\nu) \\ &\quad \times f(\lambda_\xi|a_\xi, b_\xi) \end{aligned}$$

where $f(y|\theta, \xi, \sigma, \mu)$ is the sample GP likelihood from Section 2.1 and prior distributions $f(\beta_\nu|\lambda_\nu)$, $f(\beta_\xi|\lambda_\xi)$, $f(\lambda_\nu|a_\nu, b_\nu)$ and $f(\lambda_\xi|a_\xi, b_\xi)$ are specified as follows. Parameter smoothness of GP shape and (modified) scale functions is encoded by adopting Gaussian priors for their vectors β_η of basis coefficients (for $\eta = \xi, \nu$), expressed in terms of parameter roughness R_η

$$f(\beta_\eta|\lambda_\eta) \propto \lambda_\eta^{1/2} \exp\left(-\frac{\lambda_\eta}{2}\beta_\eta'Q_\eta\beta_\eta\right).$$

The roughness coefficient λ_η can be seen, from a Bayesian perspective, as a parameter precision for β_η . It is assigned a Gamma prior distribution, which is conjugate with the prior Gaussian distribution for β_η . The values of hyper-parameters are set such that Gamma priors are relatively uninformative; a_η and b_η takes values of 10^{-3} throughout this study for all parameterisations. The Bayesian inference can be illustrated by the directed acyclic graph shown in Fig. 2.

Estimates for β_ξ , β_ν , λ_ξ and λ_ν are obtained by sampling the posterior distribution above using MCMC. We choose to adopt a Metropolis-within-Gibbs framework (e.g. Gamerman and Lopes, 2006), where each of the four parameters is sampled in turn conditionally on the values of others. The full conditional distributions $\lambda_\xi|\beta_\xi$ and $\lambda_\nu|\beta_\nu$ of precision parameters are Gamma by conjugacy, and are sampled exactly in a Gibbs step. Full conditional distributions for coefficients β_ξ and β_ν are not available in closed form; a more general Metropolis-Hastings (MH) scheme must therefore be used.

There are a number of potential alternative strategies regarding the MH step for β_η ($\eta = \xi, \nu$). We choose to examine two

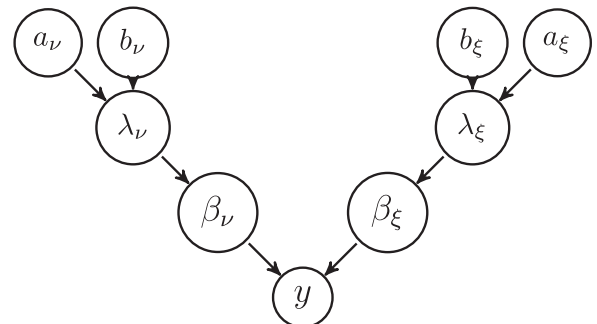


Fig. 2. Directed acyclic graphical representation of the Bayesian inference scheme.

possibilities: (a) a straightforward MH sampling of correlated Gaussian proposals for β_η , and (b) the mMALA algorithm of [Girolami and Calderhead \(2011\)](#), exploiting first- and second-derivative information from the log posterior to propose candidate values for the full vector of coefficients in high-probability regions. Implementations are described in the Appendix. Henceforth we refer to these two schemes as MH and mMALA respectively for brevity. The MH approach is simple to implement, but is likely to generate MCMC chains which mix relatively poorly. The mMALA scheme is expected to explore the posterior with considerably higher efficiency; however, its implementation requires knowledge of likelihood derivatives.

2.3.3. Comparing uncertainties

Parameter uncertainty is estimated by bootstrap resampling for ML inference and by sampling from the posterior distribution of parameters for Bayesian inference. These two approaches seek to estimate parameter uncertainty, but in different ways. Since Bayesian priors are chosen to be relatively uninformative, we expect – at least naively – that inferences concerning parameter estimates and return value distributions from the two approaches will be similar, but not the same. In general, the relationship between bootstrap uncertainty estimates and those from Bayesian inference is complex, and an open topic in the literature (e.g. [Fushiki et al., 2005](#)). A thorough theoretical analysis of this relationship in the current application is beyond the scope of the current work.

3. Evaluation of methods

This section describes evaluation of relative performance of different model parameterisations and inference schemes introduced in [Sections 2.2](#) and [2.3](#). We assess performance in terms of quality of estimation of distributions of return values corresponding to long return periods, estimated under models for large numbers of replicate samples of data from pre-specified underlying models.

We simulate 100 sample realisations, each of size 1000 from three different underlying models, described below and referred to henceforth as [Cases 1, 2 and 3](#), and further simulate 100 sample realisations of size 5000 from the same triplet of underlying models, referring to these as [Cases 4, 5 and 6](#) respectively. We next estimate extreme value models for all sample realisations, model parameterisations and inference schemes. We assume that any sample realisation (for any case) corresponds to a period \mathcal{T} years of observation. We then simulate 1000 replicates of return period realisations, each replicate consisting of observations of directional extreme values corresponding to a return period of $10 \times \mathcal{T}$, and estimate the distribution of the maximum observed (the $10\mathcal{T}$ -year maximum return value) for all model parameterisations and inference methods, by accumulation from the 1000 replicates. Return value simulations under models estimated using Bayesian inference proceed by sampling a different vector of model parameter estimates at random from the estimated joint posterior distribution of parameters and simulating the appropriate return period of events from the corresponding distribution, for each of the 1000 replicates. For ML inference, for each replicate, we sample the vector of model parameter estimates at random from a set of 100 parameter estimate vectors of generated by the bootstrap analysis. Return value distributions are estimated omnidirectionally (that is, including all directions) and for 8 directional octants centred on the cardinal and semi-cardinal directions (by considering only those observations from the return period realisation with the appropriate directional characteristics). For each

sample realisation from [Cases 4, 5 and 6](#), we estimate return value distributions for all parameterisations but for only mMALA inference, since as will be discussed in [Section 3.3](#) below, the computational effort associated with any of mMALA, MH and MLE (with bootstrap resampling) for these cases is large.

We quantify the quality of return value inference by comparing the empirical cumulative distribution function generated under the fitted model for each sample realisation with that from simulation under the known underlying case. We quantify the discrepancy between empirical distribution functions by estimating Kullback–Leibler, Cramer–von Mises and Kolmogorov–Smirnov statistics. We visualise relative performance by plotting the empirical cumulative distribution function of the test statistic over the 100 sample realisations, for each combination of case, model parameterisation and inference method. We also compare performance in terms of prediction of the 37.5th percentile of the $10\mathcal{T}$ -year return value distribution, since this is often used in metocean and coastal design applications; it corresponds approximately to the location of the mode of a Weibull distribution with shape parameter ≈ 2 . However, we are not only interested in quality of inference, but also in computational efficiency. This is evaluated and illustrated in [Section 3.3](#).

To complete the full analysis described here, comprised of 100 random samples of each of 6 cases, required running 3 dedicated workstations (exploiting each of 48 cores and 196 GB RAM per workstation) for approximately 10 weeks. All assessments of return value distributions, in terms of e.g. Kullback–Leibler divergence or a central percentile, are therefore based on 100 independent estimates per case. We note that the precision with which the median and quartiles of the distribution of Kullback–Leibler divergence, or the distribution of a central percentile (or is bias), is therefore considerably higher than that with which extreme quantiles are estimated.

In practical application to metocean design using a sample of measured or hindcast significant wave height data, for example, it is critical to demonstrate that an iterative simulation algorithm such as MH or mMALA has produced a chain which has itself converged to the stationary distribution, and whether that distribution has been adequately explored. This can be achieved, for example, by comparing inferences from multiple independent chains (e.g. [Brooks and Gelman, 1998](#)). This diagnosis is critical to such applications, and is the main evidence that MCMC inference is valid. In the current work, visual inspection of trace plots was used to confirm that an adequate period of burn-in had been specified for all combinations of parameterisation and inference scheme. However since the underlying true models are known, we choose to use comparison with the truth as the basis to assess the relative performance of different model parameterisations and inference schemes. The finding of this work is that mMALA as implemented here behaves more reasonably as a default approach to inference (requiring less user intervention and fine tuning) than MH. There is no doubt that chain convergence diagnostics would have indicated this also, and might also have prompted refinement of the MH scheme in particular to improve its performance.

3.1. Case studies considered

First, we describe model [Cases 1–6](#) used to generate sample realisations for extreme value modelling. For each case, potentially all of Poisson rate ρ of threshold exceedance, GP shape ξ and scale σ of exceedance size vary as a function of covariate θ . The extreme value threshold μ is fixed at zero throughout.

Case 1. For extreme value threshold $\mu(\theta) = 0$, we simulate 1000 observations with a uniform Poisson rate $\rho(\theta) = 1000/360$ per degree covariate, and a low order Fourier parameterisation of GP

shape $\xi(\theta) = \sin(\theta) + \cos(2\theta) + 2$ and scale $\sigma(\theta) = -0.2 + (\sin(\theta - 30))/10$.

Case 2. For extreme value threshold $\mu(\theta) = 0$ and the same Fourier parameterisation of GP shape and scale as in **Case 1**, a non-uniform Poisson rate $\rho(\theta) = \max(\sin(\theta) + 1.1, 0) \times 1000/c_p$, where $c_p = \int_0^{360} \max(\sin(\theta) + 1.1, 0) d\theta$ is used to simulate 1000 observations.

Case 3. For extreme value threshold $\mu(\theta) = 0$, the forms of each of $\rho(\theta)$, $\xi(\theta)$ and $\sigma(\theta)$ are defined by mixtures of between one and five Gaussian densities, as illustrated in **Fig. 3**. Sample size is 1000.

Cases 4, 5 and 6. These cases are identical to **Cases 1, 2 and 3** respectively, except that Poisson rate ρ is increased by a factor of five. Sample size is therefore 5000.

Fig. 3 illustrates typical sample realisations of **Cases 1, 2 and 3**. Parameter variation of GP shape ξ and scale σ with direction θ are identical in **Cases 1 and 2** Poisson rate ρ is constant in **Case 1** only. In **Cases 2 and 3**, ρ is very small at $\theta \approx 270^\circ$ leading to a sparsity of corresponding observations. ξ is largest (but negative) at $\theta \approx 120^\circ$ for **Cases 1 and 2**, leading to larger observations here. For **Case 3**, ξ is largest (and positive) at $\theta \approx 30^\circ$ leading to the heaviest tail in any of the cases considered. **Fig. 4** shows parameter estimates for ξ and σ , corresponding to the sample realisation of **Case 2** shown in **Fig. 3**, for different model parameterisations using mMALA inference. Visual inspection suggests that estimates of similar quality are obtained using all of Spline, Fourier and Gaussian Process parameterisations, but that the Constant parameterisation is poor. It is also apparent that identification of ξ is more difficult than σ . Corresponding plots (not shown) for maximum likelihood and Metropolis–Hastings inference show broadly similar

characteristics, as do plots for other realisations of the same case, and realisations of different cases. (Posterior) cumulative distribution functions of return values based on models for the sample realisations of **2** illustrated in **Figs. 3 and 4**, corresponding to a return period of 10 times the period of the original sample, are shown in **Fig. 5** for different model parameterisations and mMALA inference.

It can be seen omnidirectionally that the Constant parameterisation provides best agreement with the known return value distribution, despite the fact that parameter estimates in **Fig. 4** do not reflect the directional non-stationarity present. In **Jonathan et al. (2008)**, it is demonstrated that a stationary extreme value model may produce good estimates for omni-directional return value distributions in at least two situations. Firstly, when the extreme value threshold is set sufficiently high that only observations from the most extreme interval of the covariate domain are modelled, therefore threshold exceedances and estimated models will be effectively stationary. Secondly, it may be that different bias effects introduced by the stationarity assumption compensate for each other in estimation of return values at a certain return levels. For example, with reference to the sample illustrated in **Fig. 4**, the constant model underestimates the value of GP shape in general, but does a relatively good job of identifying the maximum values of shape and scale in directions around 120° . The constant model does a very poor job for directions around 270° . The uncertainty in parameter estimates from the constant model is also lower since the number of parameters to estimate is lower for that model. Reference to **Fig. 3** for **Case 2** shows also that the rate of occurrence of threshold exceedances is larger for directions where the constant model performs well. These characteristics are reflected in the corresponding return value estimates in **Fig. 5**. For directional octants, the constant model does

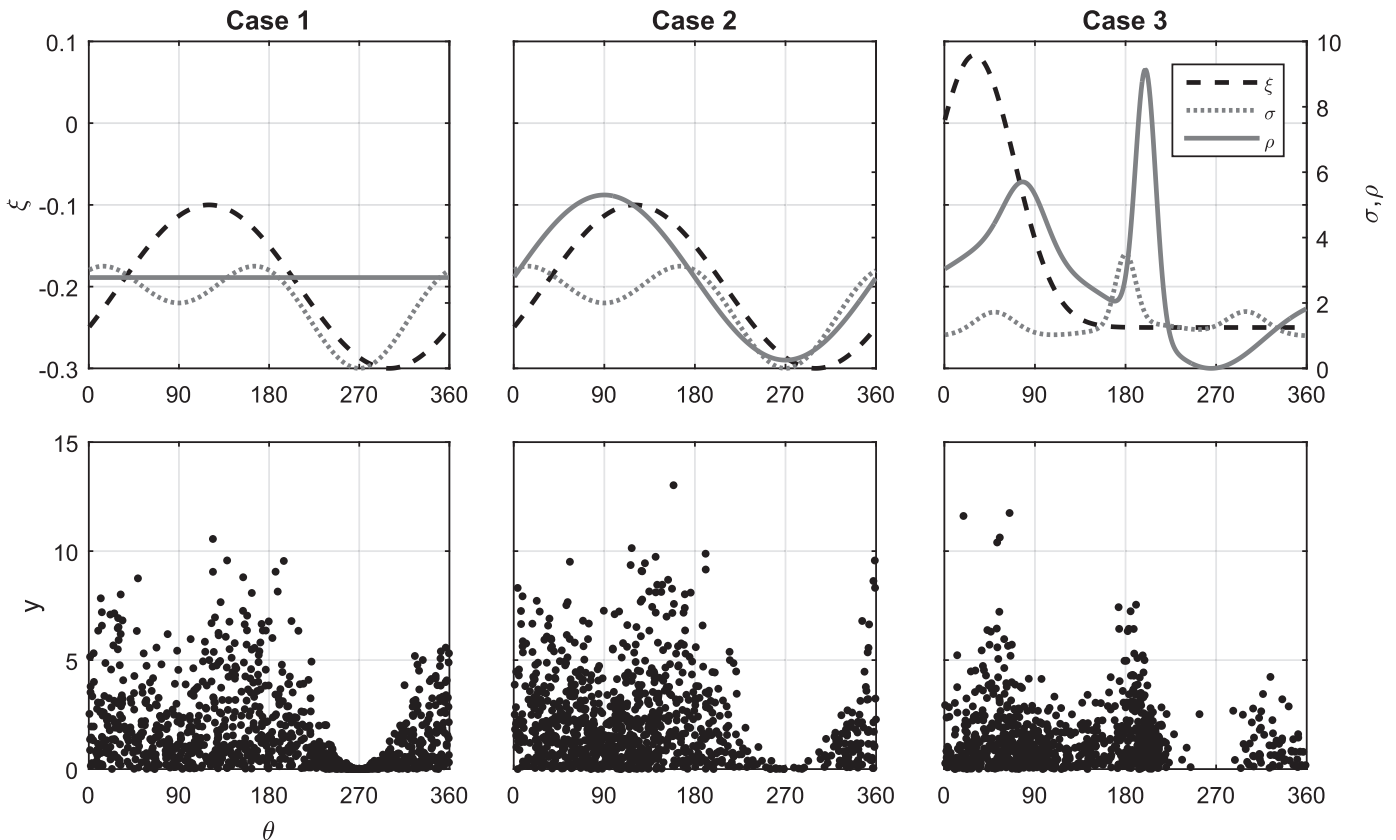


Fig. 3. Illustrations of sample realisations from each of **Cases 1** (left), **2** (centre) and **3** (right). Upper panels show parameter variation of GP shape ξ , scale σ and Poisson rate ρ with direction θ for each case. Lower panels show the 10th realisation of the corresponding simulated samples. ξ and σ for **Cases 4, 5 and 6** are identical to those of **Cases 1, 2 and 3** respectively. The value of ρ for **Cases 4, 5 and 6** is five times that of **Cases 1, 2 and 3** respectively.

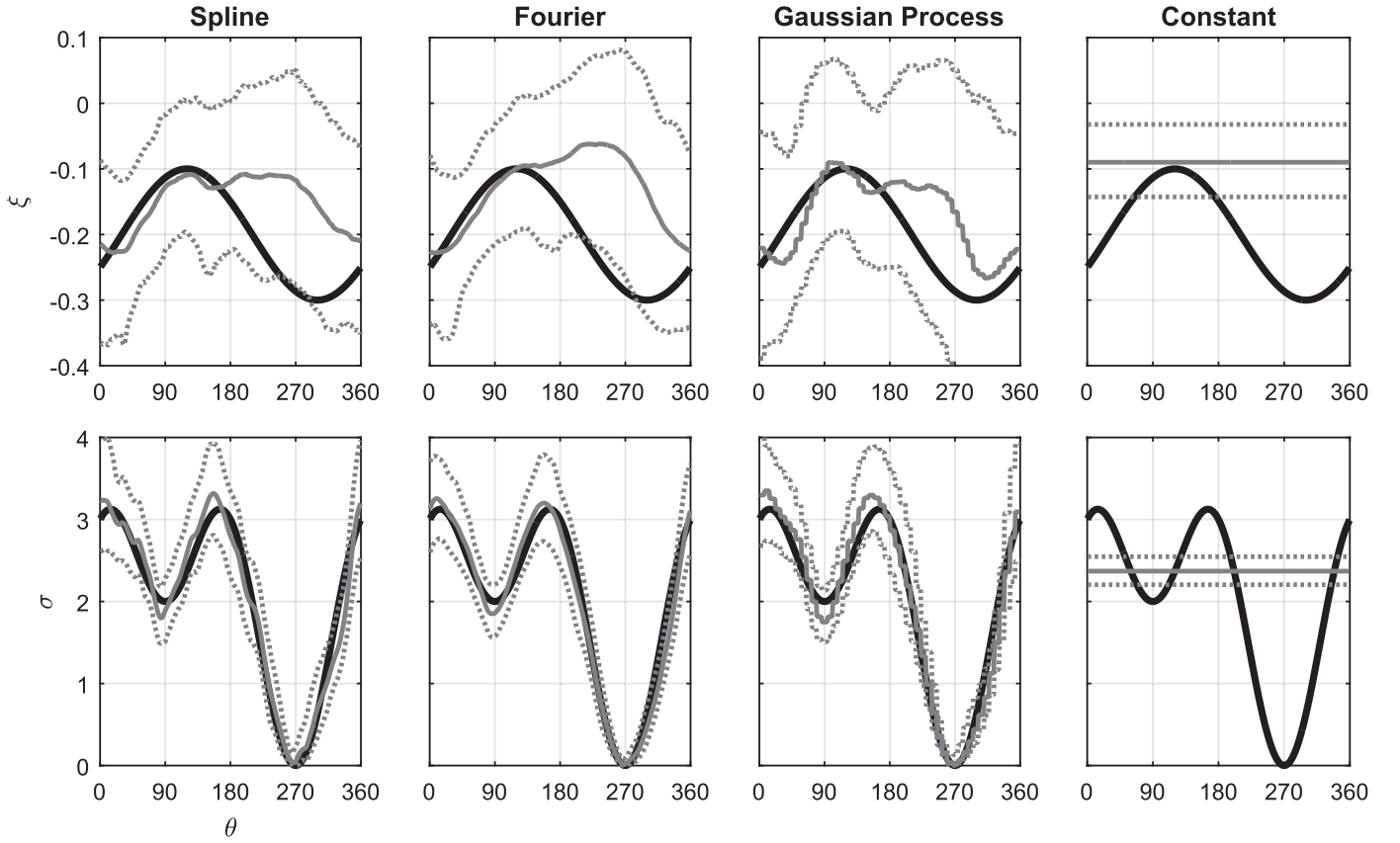


Fig. 4. Parameter estimates for GP shape ξ (upper) and scale σ (lower) for the sample realisation of Case 2 shown in Fig. 1, for different model parameterisations (left to right: Spline, Fourier, Gaussian Process, Constant) using mMALA inference. Each panel illustrates the true parameter (solid black), posterior median estimate (solid grey) with 95% credible interval (dashed grey).

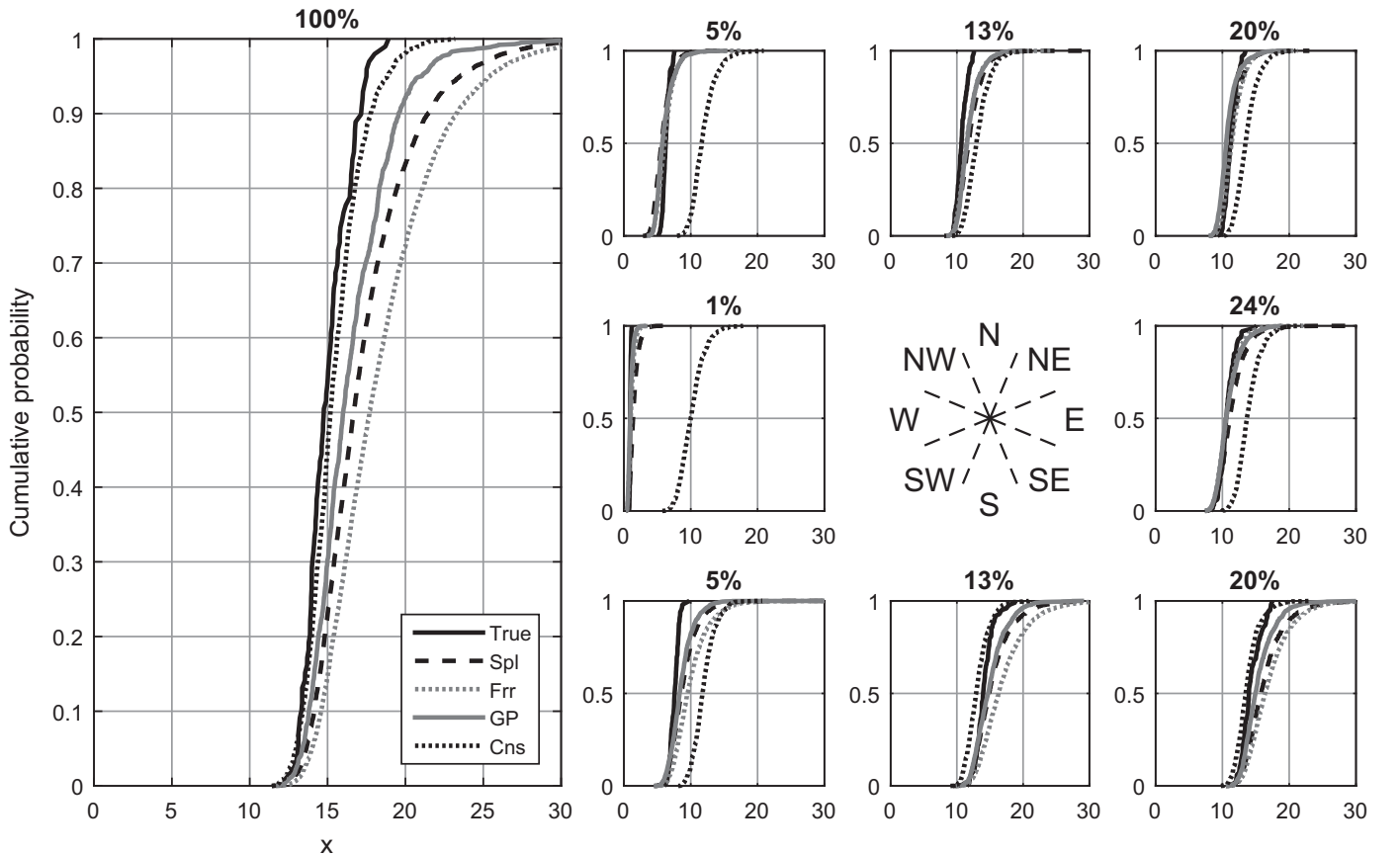


Fig. 5. Posterior cumulative distribution functions of return value for the sample realisation of Case 2 shown in Fig. 1, corresponding to a return period of 10 times the period of the original sample. The left hand panel shows the omnidirectional return value distribution, and right hand panels the corresponding directional estimates. The title for each panel gives the expected percentage of individuals in that directional sector. In each panel, estimates are given for different model parameterisations using mMALA inference. The true return value distribution is given in solid black.

relatively well for directions around 120° , but very poorly around 270° . However, the omnidirectional estimate is dominated by directions around 120° . The relatively good performance of the directional model around 120° , together with the large rate of occurrence of events there and the relatively small parameter uncertainties for the constant model, result in its providing good estimation (in this case) of the omnidirectional return value. This occurs despite the fact that the directional bias of the constant model is greater than that of any of the non-stationary models. In general however, it is not possible to know a priori how a constant model will perform in estimating the omnidirectional value. It is clear however that its directional bias will be larger than that of an appropriate directional model. Reference to Fig. 9 below, assuming that mMALA provides more satisfactory inference as implemented here, suggests that values of KL divergence are somewhat larger and more variable for the constant model compared to the Spline or Gaussian process parameterisations. Reference to Fig. 10 shows that the Constant model does very poorly for certain directional sectors. Omnidirectionally, and for 8 directional octant sectors, non-stationary model parameterisations perform similarly. However, the Constant parameterisation does particularly poorly for the western and north-western sectors, for which the rate of occurrence of events is relatively low, and both ξ and σ are near their minimum values. Fig. 6 illustrates uncertainty (over all 100 sample realisations) in the cumulative distribution function of return value for Case 2, corresponding to a return period of 10 times the period of the original sample, using the Spline parameterisation and mMALA inference. The median estimate for the return value distribution (over all 100 sample realisations of 2) is shown in

solid grey, with corresponding point-wise 95% uncertainty band in dashed grey. The true return value distribution is given in solid black. There is good agreement in all sectors. We explore differences in inferences for return value distributions more fully in Section 3.2.

3.2. Assessing quality of inference

The criteria used to compare distributions of return values are now described. Since, for comparison only, we only have access to samples from distributions, where necessary we project empirical distributions onto a linear grid using linear interpolation, and evaluate grid-based approximations to facilitate comparison. Then we compare empirical return value distributions using each of the following three statistics. The Kolmogorov–Smirnov criterion compares two distributions in terms of the maximum vertical distance between cumulative distribution functions, as $D_{ks}(F_0, F_1) = \sup_x |F_1(x) - F_0(x)|$. The Cramer–von Mises criterion evaluates the average squared difference of one distribution from a second, reference distribution, using $\int_{-\infty}^{\infty} (F_1(x) - F_0(x))^2 dx$. The Kullback–Leibler divergence compares distributions using the average ratio of logarithms of density functions $D_{kl}(F_0, F_1) = \int_{-\infty}^{\infty} \log\left(\frac{f_0(x)}{f_1(x)}\right) f_0(x) dx$; in this work, we use the approximation of Perez-Cruz (2008). The general characteristics of differences in return value inference due to model parameterisation and inference method were found to be similar for each of the three statistics. Only comparisons using Kullback–Leibler (KL) divergence are therefore reported here. We note that perfect agreement between $f_1(x)$ and $f_0(x)$ yields a

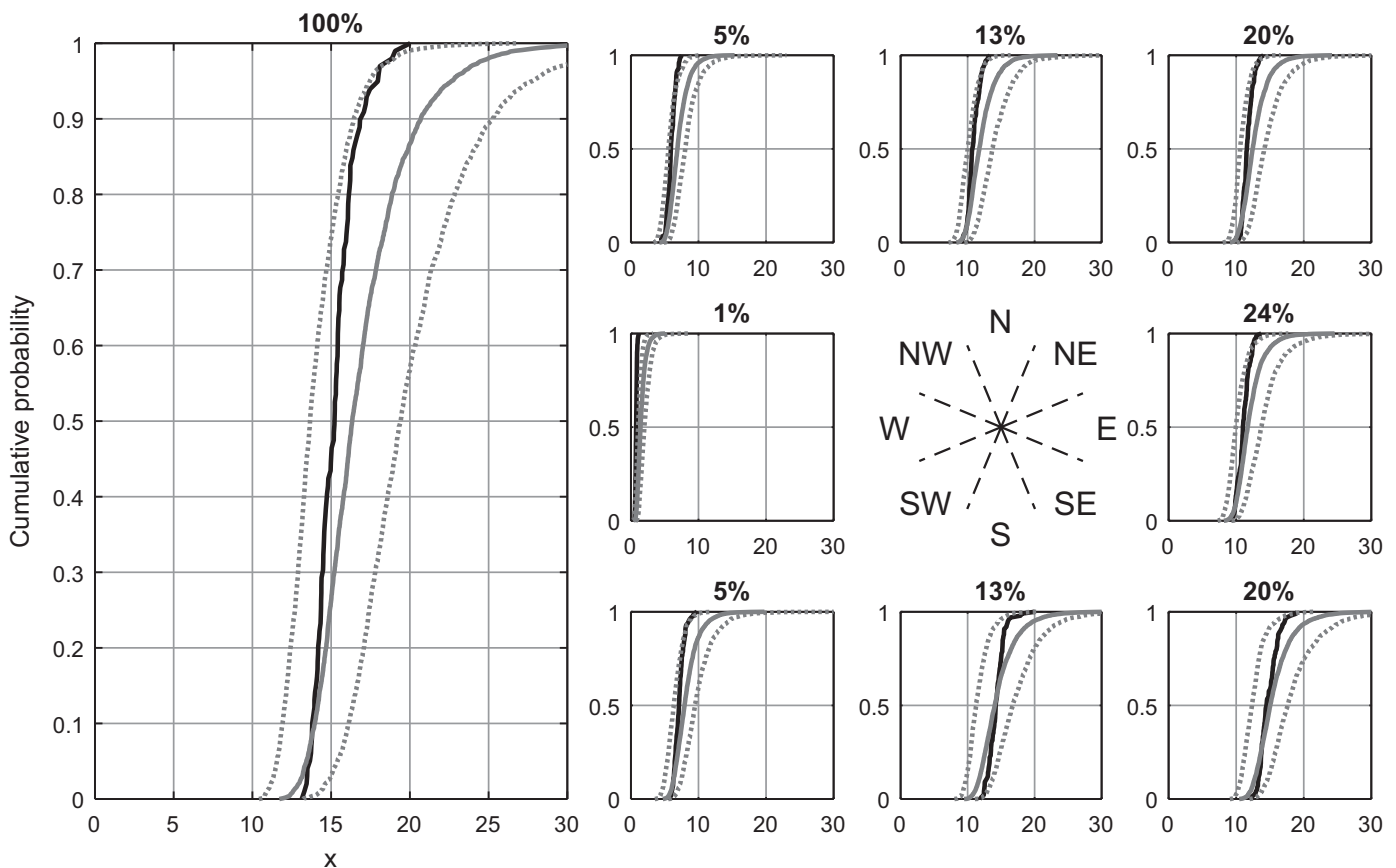


Fig. 6. Uncertainty (over all 100 sample realisations) in the cumulative distribution function of return value for Case 2, corresponding to a return period of 10 times the period of the original sample. The left hand panel shows the omnidirectional return value distribution, and right hand panels the corresponding directional estimates. The title for each panel gives the expected percentage of individuals in that directional sector. In each panel, the true return value distribution is given in solid black. The median estimate (over realisations) for return value distribution of the Spline model parameterisation using mMALA inference is shown in solid grey, with corresponding point-wise 95% uncertainty band in dashed grey.

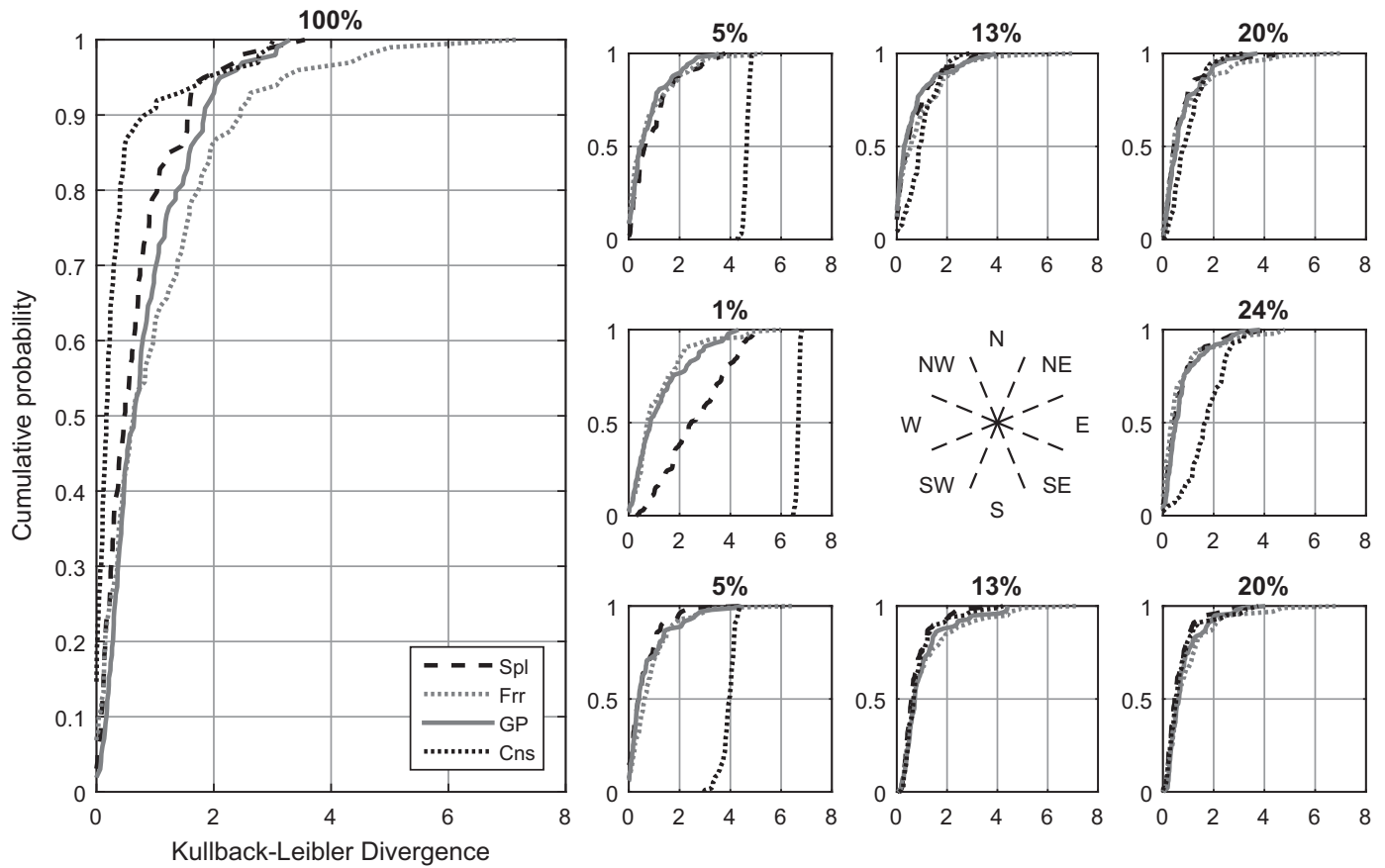


Fig. 7. Empirical cumulative distribution functions of the Kullback–Leibler divergence between return value distributions (corresponding to a return period of 10 times that the original sample) estimated under the true model and those estimated under models of sample realisations with different parameterisations and mMALA inference for Case 2. The title for each panel gives the expected percentage of individuals in that directional sector.

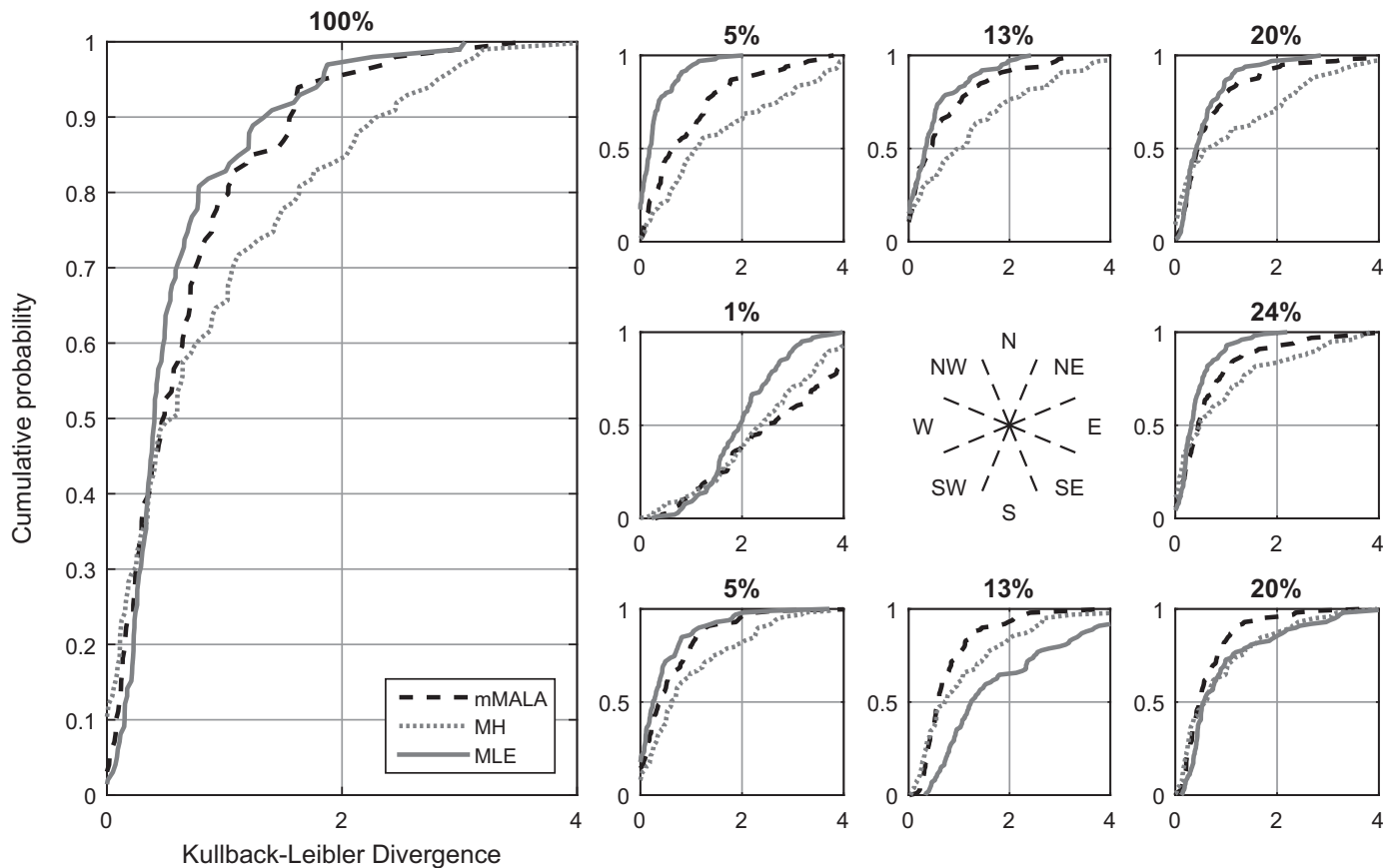


Fig. 8. Empirical cumulative distribution functions of the Kullback–Leibler divergence between return value distributions (corresponding to a return period of 10 times that the original sample) estimated under the samples from the true return value distribution and those estimated under models of sample realisations with Spline parameterisations and different inference procedures for Case 2. The title for each panel gives the expected percentage of individuals in that directional sector.

minimum KL divergence of zero.

For illustration, Fig. 7 shows empirical cumulative distribution functions for the KL divergence between return value distributions (corresponding to a return period of 10 times that the original sample) estimated under the true model and those estimated under models of sample realisations with different parameterisations and mMALA inference for 2. The distributions of KL divergence from all non-stationary model parameterisations appear to be very similar, as might be expected from consideration of figures similar to Fig. 5. However, the Constant model yields the best performance omnidirectionally in this case (since the corresponding distribution of KL divergence is shifted towards zero). In stark contrast, the Constant model does particularly badly in the eastern, south-western, western and north-western sectors. Fig. 8 gives empirical cumulative distribution functions of the KL divergence between return value distributions (corresponding to a return period of 10 times that of the original sample) estimated under the true return value distribution and those estimated under models of sample realisations with Spline parameterisations and different inference procedures for the same case. There appears to be little to choose between mMALA and MLE inference methods for this case, with MH somewhat poorer.

Fig. 9 summarises the characteristics of distributions for KL divergence corresponding to the omnidirectional return value distribution for all cases, model parameterisations and inference methods considered in this work. In general, we note that all non-

stationary parameterisations perform well with mMALA inference. With MH inference, performance is generally poorer, especially for Fourier parameterisation. MLE does better than MH. We note that the Constant parameterisation generally performs well for the omnidirectional return value, but there is some erratic behaviour, notably for Case 4. Fig. 10 is the corresponding plot for the (generally sparsely populated) western directional sector. The Spline and Fourier parameterisations with mMALA inference perform best. We note that the Fourier parameterisation does less well using MH and MLE, and that the Constant parameterisation behaves very erratically. We also note that, somewhat surprisingly, the Gaussian Process model performs considerably less well than the Spline and Fourier parameterisations. We surmise that this is due to mean-reversion in the absence of observations, compared with the other parameterisations which prefer interpolation to reduce parameter roughness. We note that the Constant parameterisation also performs badly for Case 4.

From a practitioner's perspective, it is also interesting to quantify the performance of different model parameterisations and inference methods in estimating some central value (e.g. the mean, median, or 37.5th percentile) of the distribution of the return value. We choose the 37.5th percentile, since this percentile is near the mode of the distribution, and commonly used in the met-ocean community.

Fig. 11 illustrates this comparison in terms of a box-whisker plot. In each panel, the true value, estimated by simulation under

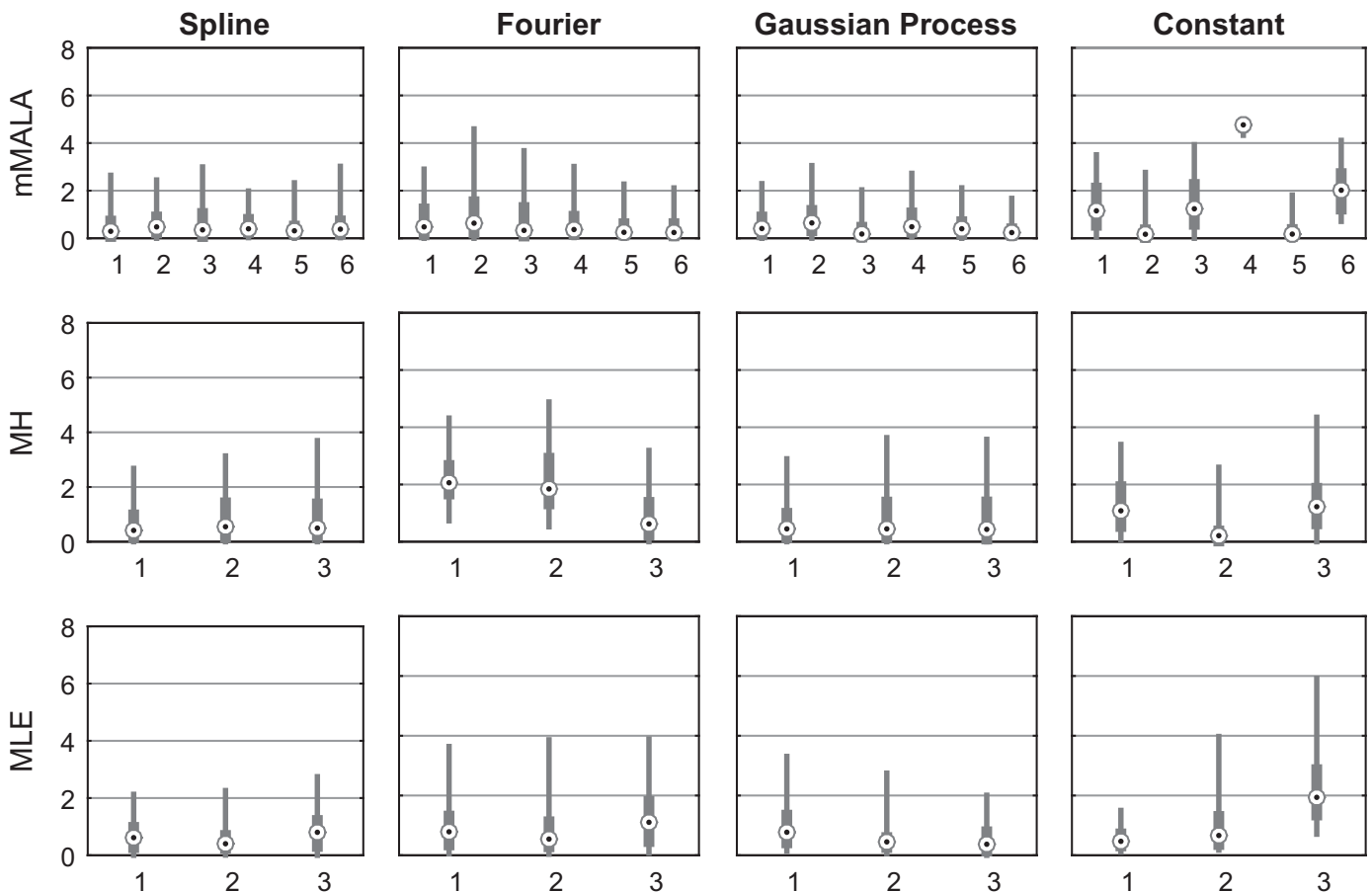


Fig. 9. Box-whisker comparison of samples of Kullback–Leibler (KL) divergence between omnidirectional return value distributions (corresponding to a return period of 10 times that the original sample) estimated under samples from the true return value distribution and those estimated under models of each of 100 sample realisations. mMALA inference is reported for all six cases (abscissa labels) and model parameterisations (columns of panels). Metropolis–Hastings (MH) inference and maximum likelihood estimation (MLE) are reported only for the smaller sample sizes. The sample of KL divergence is summarised by the median (white disc with black central dot), the interquartile range (grey rectangular box) and the (2.5%, 97.5%) interval (grey line).

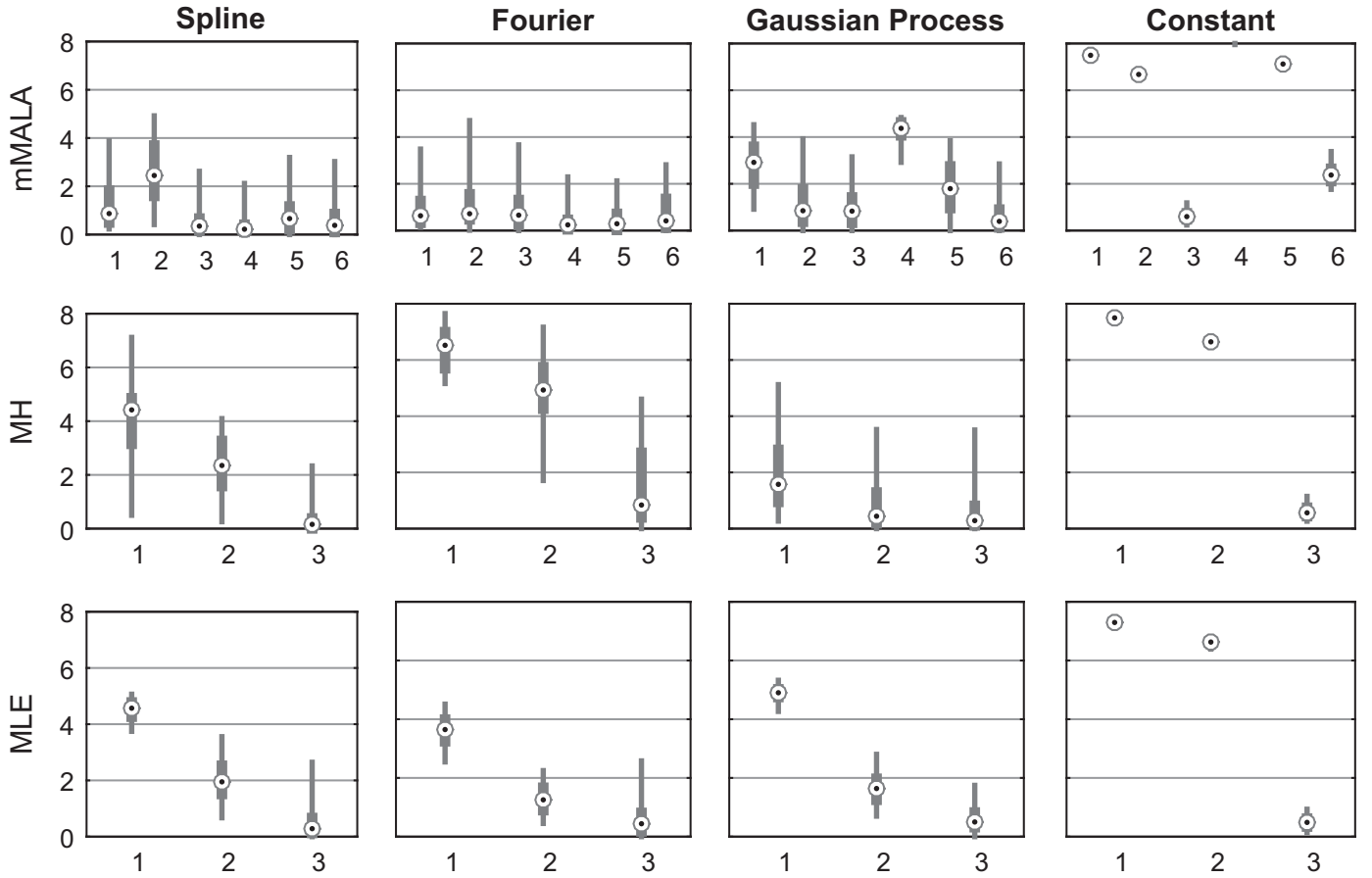


Fig. 10. Box-whisker comparison of samples of Kullback–Leibler (KL) divergence between western sector return value distributions (corresponding to a return period of 10 times that the original sample) estimated under samples from the true return value distribution and those estimated under models of each of 100 sample realisations. mMALA inference is reported for all six cases (abscissa labels) and model parameterisations (columns of panels). Metropolis–Hastings (MH) inference and maximum likelihood estimation (MLE) are reported only for the smaller sample sizes. The sample of KL divergence is summarised by the median (white disc with black central dot), the interquartile range (grey rectangular box) and the (2.5%, 97.5%) interval (grey line). The ordinate scale is the same as that of Fig. 9 to facilitate comparison. The Constant parameterisation for Case 4 with mMALA inference yields values of KL divergence larger than 8.

the true model, is shown as a black disc. The distribution of estimates from 100 different sample realisations of each case is summarised by the median (white disc with black central dot), the interquartile range (grey rectangular box) and the (2.5%, 97.5%) interval (grey line).

The performance of non-stationary parameterisations with mMALA is good, with the possible exception of Case 3 (and Case 6); yet MH inference tends to produce greater bias and variability in estimates of the 37.5th percentile. We may surmise that this difference may be due to the fact that mMALA exploits knowledge of likelihood gradient and curvature. The Constant model again performs more erratically. Corresponding box-whisker plots for the eastern directional octant (not shown) show similar characteristics: for a given inference scheme, all non-stationary models yield similar performance, but the Constant model overestimates throughout. It is interesting that MLE shows some bias in return value estimation for the omnidirectional 37.5th percentile, but this is not the case in general. True values of the 37.5th percentile for the western sector (see Fig. 12) are considerably lower than for the eastern sector, and lower again than the omnidirectional values. In this sector, the rate of occurrences of events is generally lower in all cases. Nevertheless, Fig. 12 has many similar features to Fig. 11. However, we note that MH struggles in combination with the Fourier parameterisation, probably since the latter has the whole of the covariate domain as its support; intelligent proposals (like

those used here in MLE and mMALA) are necessary. Estimates using the Constant model are erratic. Overall, we note that MLE and mMALA inference for all of Spline, Fourier and Gaussian Process parameterisations perform relatively well, and equally well.

3.3. Assessing efficiency of inference

The effective sample size (m^* , e.g. Geyer, 1992) gives an estimate of the equivalent number of independent iterations that a Markov chain Monte Carlo represents, and is defined by $m^* = m / (1 + 2 \sum_{k=1}^{\infty} c_k)$, where c_k is the autocorrelation of the MCMC chain at lag k , and m is the actual chain length. The effective sample size per hour is defined by m^*/T , where T is the elapsed computational time (in hours) for m steps of the chain. For maximum likelihood inference with bootstrap uncertainty estimation, since bootstrap resamples are independent of one another, we estimate effective sample size per hour as m_{BS}/T where m_{BS} is the number of bootstrap resamples used and T is now the total elapsed computational time (in hours) to execute analysis of the m_{BS} bootstrap resamples. Comparison of effective sample sizes per hour for different cases, parameterisations and inference methods gives some indication of relative computational efficiency, although objective comparison is difficult. In particular we note that software implementations in MATLAB exploiting common computational structures between different approaches have been used; these are almost certainly to the detriment of computational

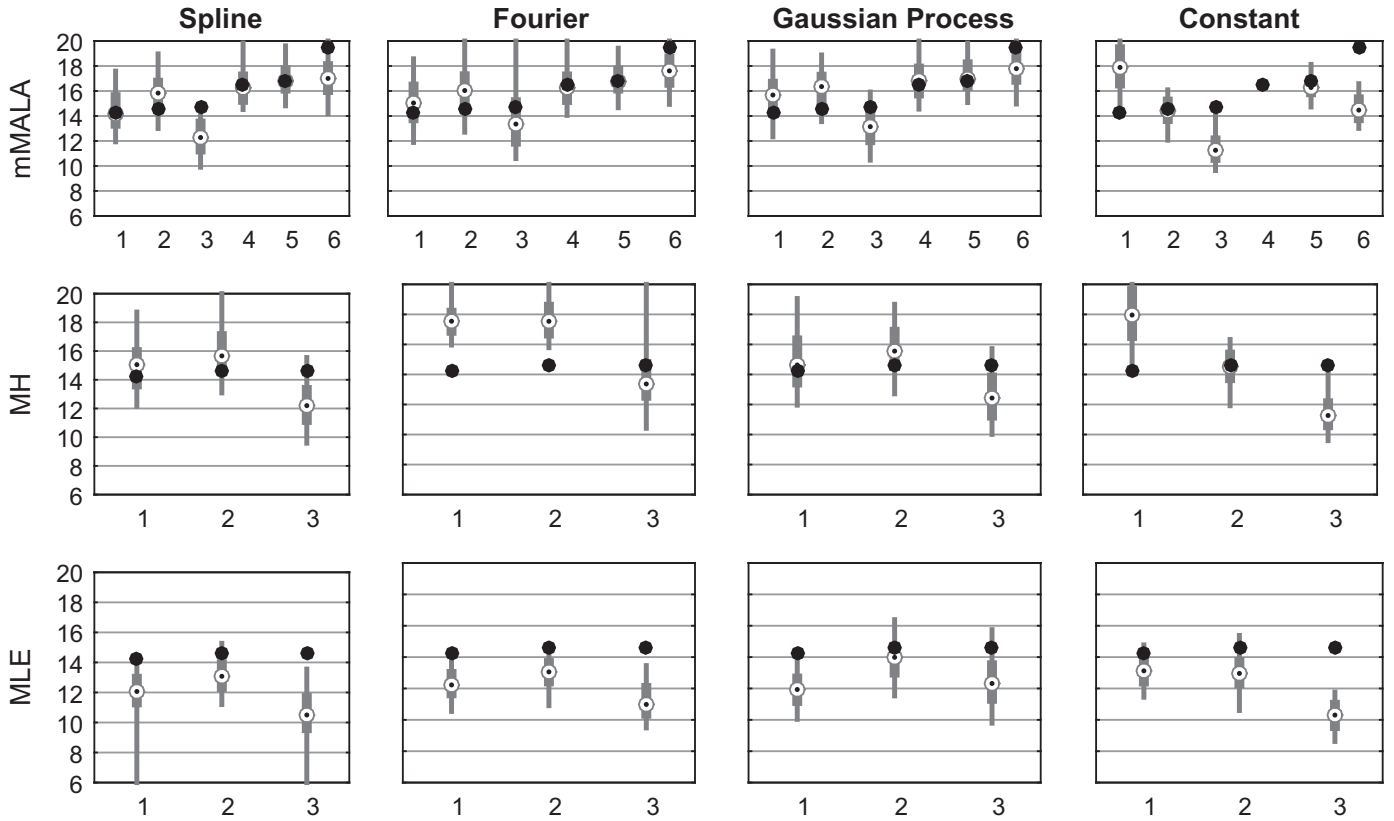


Fig. 11. Box-whisker comparison of estimates for the 37.5th percentile of the omnidirectional return value distribution (in metres) for different cases, model parameterisations and inference procedures. mMALA inference is reported for all six cases (abscissa labels) and model parameterisations (columns of panels). Metropolis–Hastings (MH) inference and maximum likelihood estimation (MLE) are reported only for the smaller sample sizes. In each panel, the estimate from simulation under the true model is shown as a black disc. The distribution of estimates from 100 different sample realisations is summarised by the median value (white disc with black central dot), the interquartile range (grey rectangular box) and the (2.5%, 97.5%) interval (grey line). The Constant parameterisation for Case 4 with mMALA inference yields values larger than 20 m.

efficiency for some of the approaches, particularly the Constant parameterisation. For this reason, we do not report effective sample size per hour for the Constant parameterisation. Computational run-times are also of course critically dependent on software and hardware resources used. We note that the focus of this work is primarily quality of inference, rather than its computational efficiency. Specific modelling choices, such as the set-up of the cross-validation strategy adopted for MLE and choice of burn-in length and proposal step-size for MH and mMALA within reasonable bounds may not influence inferences greatly, but will obviously however affect run times. Similarly the Spline, Fourier and Gaussian Process parameterisations used were chosen to be of similar complexity, but small differences may again influence relative computational efficiency of inference. With these caveats in mind, Fig. 13 illustrates the distribution of estimated effective sample size (ESS), and effective sample size per hour (ESS/hr) for different cases, model parameterisations and inference procedures.

The left hand side of Fig. 12 illustrates $\log_{10}(\text{ESS})$ for all combinations of parameterisations and inference schemes. For MLE inference, we choose to report the number of bootstrap resamples used, as described in Section 3. The effective sample sizes for mMALA inference are considerably larger than for MH for B-spline and Fourier parameterisations. For the Gaussian process parameterisation, mMALA still provides a larger ESS, but the difference between mMALA and MH is smaller. For mMALA inference, ESS is largest for the B-spline parameterisation; the Gaussian process parameterisation provides the smallest ESS on average. For MH

inference, ESS for B-splines and Fourier parameterisations is near 10, suggesting that the posterior density has not been sufficiently explored due to poor MCMC mixing using the MH algorithm as implemented. The value of ESS is approximately constant across the different cases examined for all combinations of model parameterisation and inference method. The right hand side of Fig. 12 indicates that for B-spline parameterisation, ESS/hr is also higher for mMALA than for MH. For Fourier and Gaussian process parameterisations, ESS/hr is comparable for mMALA and MH. The right hand side of Fig. 13 shows that, for mMALA inference, the ESS/hr is considerably lower for Cases 4, 5 and 6, indicating that inference using large sample sizes is slower. For this reason, in this work, we do not provide results for Cases 4, 5 and 6 using MLE and MH. Overall, comparing non-stationary parameterisations, ESS/hr is larger for Splines and Gaussian Processes than for Fourier. There is little difference in ESS/hr for different model parameterisations using MLE.

4. Discussion

Adequate allowance for non-stationarity is essential for realistic environmental extreme value inference. Ignoring the effects of covariates leads to unrealistic inference in general. The applied statistics and environmental literature provides various competing approaches to modelling non-stationarity. Adoption of non- or semi-parametric functional forms for generalised Pareto shape and scale in peaks over threshold modelling is far preferable in general

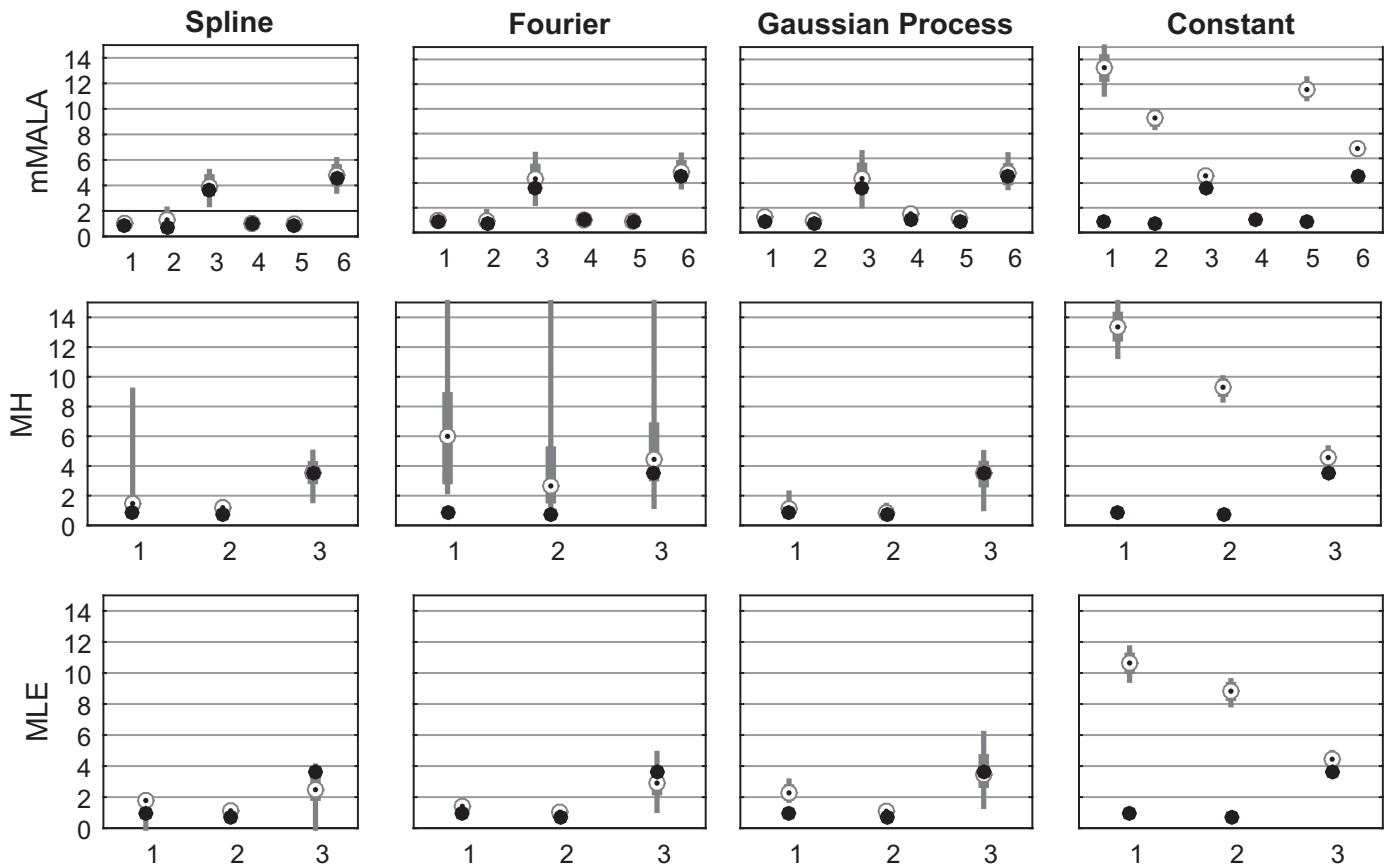


Fig. 12. Box-whisker comparison of estimates for the 37.5th percentile of the return value distribution (in metres) for the western directional sector (least populous for Cases 2, 3, 5 and 6), for different cases, model parameterisations and inference procedures. mMALA inference is reported for all six cases (abscissa labels) and model parameterisations (columns of panels). Metropolis–Hastings (MH) inference and maximum likelihood estimation (MLE) are reported only for the smaller sample sizes. In each panel, the estimate from simulation under the true model is shown as a black disc. The distribution of estimates from 100 different sample realisations is summarised by the median value (white disc with black central dot), the interquartile range (grey rectangular box) and the (2.5%, 97.5%) interval (grey line). The Constant parameterisation for Case 4 with mMALA inference yields values larger than 15 m.

in real world applications, than the assumption of a less flexible parameterisation. We find that B-spline and Gaussian Process parameterisations estimated by Bayesian inference (using mMALA) perform well in terms of quality and computational efficiency, and generally outperform alternatives in this work.

The Gaussian Process parameterisation is computationally unwieldy for larger problems, unless gridding on the covariate domain is performed. The Fourier parameterisation, utilising bases with global support (compared to Spline and Gaussian Process basis functions whose support is local on the covariate domain), is generally somewhat more difficult to estimate well in practice, showing greater instability to choices such as starting solution for maximum likelihood estimation. The Constant parameterisation performs surprisingly well in estimating the omnidirectional return values distribution in some cases, but is generally very poor in estimating directional variation.

Various choices of methods of inference are also available. Competing approaches include maximum (penalised) likelihood optimisation and Bayesian inference using Markov chain Monte Carlo sampling. It appears however that the major difference, in terms of practical value of inference, is not between frequentist and Bayesian paradigms but rather the advantage gained by exploiting knowledge of likelihood gradient and curvature. In addition, it appears that inference schemes which sample from a negative log likelihood surface randomly, rather than seeking its minimum deterministically, are more stable, and therefore more routinely implementable and useable. Moreover, Bayesian

inference gives a more intuitive framework for statistical learning and communication of uncertainty, particularly to a non-specialist audience.

Here, we have focussed on the estimation of non-stationary shape and scale parameters for the conditional distribution of independent peaks over threshold. In practical application, careful estimation of non-stationary extreme value threshold is at least as important for reliable inference. We emphasise that for practical application, any non-stationarity of extreme value threshold should be examined and identified either before or alongside non-stationarity of GP parameters. Anderson et al. (2001) note that the combination of non-stationary threshold and stationary GP shape and scale is sufficient for modelling a sample of significant wave heights in the North Sea. Physical and statistical intuition suggest, when considering an extreme value model for a quantity such H_s , that non-stationary estimates should be sought in order for each of (a) extreme value threshold, (b) then GP scale, and (c) finally GP shape, and adopted only if justified statistically. Note however that physically plausible oceanographic examples corresponding, for instance, to stationary extreme value threshold but non-stationary GP parameters are also conceivable. In the current work, we assume effectively that any non-stationarity in extreme value threshold as already been identified perfectly and its effect removed from the sample cases considered.

It appears that specification of a piecewise model for the whole sample (e.g. Behrens et al., 2004; MacDonald et al., 2011,

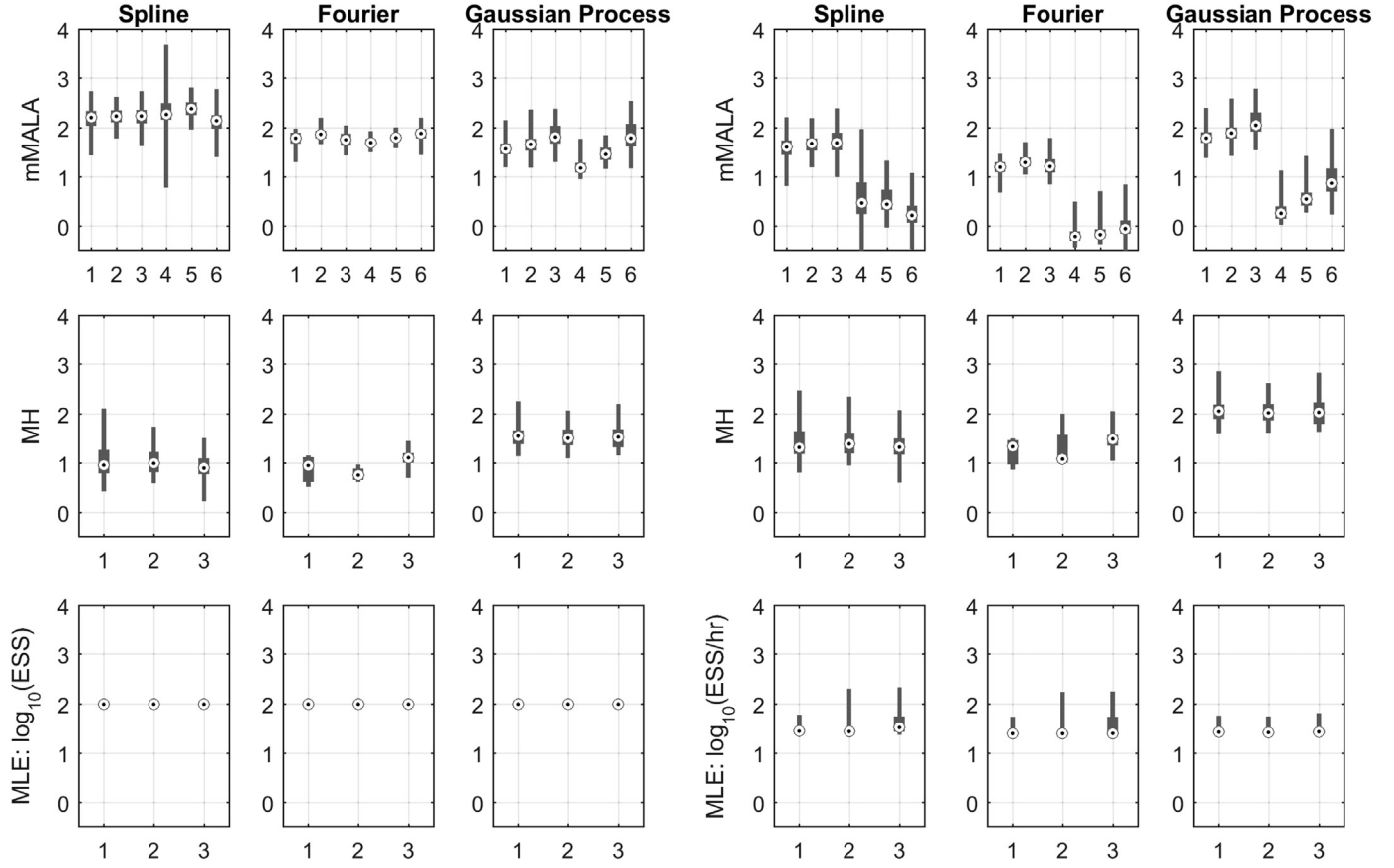


Fig. 13. Estimates for effective sample size (ESS, left hand side) and effective sample size per hour (ESS/hr, right hand side) on logarithm base 10 scale for different cases, non-stationary model parameterisations and inference procedures. mMALA inference is reported for all six cases (abscissa labels) and model parameterisations (columns of panels). Metropolis–Hastings (MH) inference and maximum likelihood estimation (MLE) are reported only for the smaller sample sizes. The distribution of estimates from 100 different realisations of the original sample is summarised by the median value (white disc with black central dot), the interquartile range (grey rectangular box) and the (2.5%, 97.5%) interval (grey line).

Randell et al., 2015b) incorporating the appropriate tail form, rather than a model for threshold exceedances in isolation is a promising route, since then threshold and tail parameters can be estimated together. However, even the simplest form of a whole sample model requires estimation of additional parameters for the model below the threshold; these parameters are also typically non-stationary with respect to covariates. Efficient and reliable non-stationary estimation using parameterisations and inference schemes similar to those presented here is key.

Acknowledgement

We are grateful to Kathryn Turnbull of Lancaster University, UK for useful comments.

Appendix A

A.1. Maximum likelihood estimation

For maximum likelihood estimation (MLE), we use a back-fitting (or iteratively re-weighted least-squares, IRLS) algorithm to estimate vectors of basis coefficients β_η ($\eta = \xi, \nu$) derived in Jonathan et al. (2014). For fixed value of smoothness parameter λ_η , we initialise coefficients to starting value $\beta_\eta^{(0)}$, and then iterate the

following step until convergence

$$\beta_\eta^{(i+1)} = \left(B_\eta^T W(\beta_\eta^{(i)}) B_\eta + \lambda_\eta Q_\eta \right)^{-1} \left(B_\eta^T V(\beta_\eta^{(i)}) + B_\eta^T W(\beta_\eta^{(i)}) B_\eta \beta_\eta^{(i)} \right)$$

where

$$V(\beta_\eta) = \nabla_\eta L(y|\Omega) \quad \text{and} \quad W(\beta_\eta) = -\nabla_\eta \nabla_\eta^T L(y|\Omega)$$

are derived at the end of this Appendix. This algorithm is similar to the mMALA algorithm (see below) used to generate proposals for the corresponding Metropolis–Hastings step in MCMC, in that both exploit first- and second-derivative information to move towards regions of high probability. The back-fitting iteration is of course deterministic, whereas the mMALA step is stochastic. In practice we use the expected values of $W(\beta_\eta)$ for ease of computation.

A.2. MCMC sampling algorithms

Denoting the set of parameters to be estimated by $\Omega = \{\beta_\xi, \beta_\nu, \lambda_\xi, \lambda_\nu\}$, inference proceeds by sampling from the full conditional distributions $f(\Omega_k | y, \theta, \Omega_{-k}, \Gamma)$ for each parameter in Ω in turn, where $\Gamma = \{a_\xi, b_\xi, a_\nu, b_\nu\}$ is the set of fixed hyperparameters for prior distributions. The form of the full conditional distribution varies depending on the type of parameter being estimated, as explained below.

Full conditional distributions for basis coefficients: Vector β_η ($\eta = \xi, \nu$) has the following conditional distribution

$$f(\beta_\eta | y, \theta, \Omega_{-\beta_\eta}, \Gamma) \propto f(y | \theta, \beta_\eta) f(\beta_\eta | \lambda_\eta)$$

which is not available in closed form, and therefore cannot be sampled directly in a Gibbs step. Instead, we generate samples using the Metropolis–Hastings algorithm: given current state $\beta_\eta^{(i)}$, we propose new parameter β_η^* from proposal distribution $f(\beta_\eta^* | \beta_\eta^{(i)})$ and evaluate the acceptance ratio

$$A(\beta_\eta^*, \beta_\eta^{(i)}) = \frac{f(\beta_\eta^* | y, \theta, \Omega_{-\beta_\eta}) f(\beta_\eta^{(i)} | \beta_\eta^*)}{f(\beta_\eta^{(i)} | y, \theta, \Omega_{-\beta_\eta}) f(\beta_\eta^* | \beta_\eta^{(i)})}$$

accepting the proposal with probability $q = \min(1, A)$, setting $\beta_\eta^{(i+1)} = \beta_\eta^*$. Otherwise we reject the proposal and set $\beta_\eta^{(i+1)} = \beta_\eta^{(i)}$. As outlined in Section 2.3 and detailed below, we consider two different methods for generating multivariate proposals for β_η . In the first approach (referred to as MH), we make an entirely stochastic Gaussian random walk proposal using a fixed covariance matrix; in the second (referred to as mMALA), we make a proposal which is partly deterministic and partly stochastic, accounting for local curvature of the likelihood surface.

For Metropolis–Hastings (MH) inference, we generate Gaussian random walk proposals of the form

$$\beta_\eta^* = \beta_\eta^{(i)} + (B_\eta^T B_\eta + \kappa_\eta Q_\eta)^{-1} \nu_\eta \epsilon$$

where ϵ is a vector of independent standard Normal random variables, and the values of step size ν_η and scale factor κ_η are adjusted to achieve reasonable acceptance rates of approximately 0.25. For inference using the Riemann manifold Metropolis-adjusted Langevin algorithm (mMALA, as implemented by Girolami and Calderhead, 2011) we propose using derivatives of the target distribution at the current sample. This promotes proposals in regions of higher probability, at the additional computational cost of computing necessary derivatives and matrix inverses. At iteration i of the sampling algorithm, where the current sample of the coefficients is $\beta_\eta^{(i)}$, proposals are made as

$$\beta_\eta^* = \beta_\eta^{(i)} + \frac{\nu_\eta^2}{2} G^{-1}(\beta_\eta^{(i)}) D(\beta_\eta^{(i)}) + \nu_\eta \sqrt{G^{-1}(\beta_\eta^{(i)})} \epsilon$$

where ϵ is a vector of independent standard Normal random variables, ν_η is (adjustable) step size, and

$$D(\beta_\eta^{(i)}) = \nabla_{\beta_\eta} L(\beta_\eta) \Big|_{\beta_\eta^{(i)}} \quad \text{and} \quad G(\beta_\eta^{(i)}) = -\nabla_{\beta_\eta} \nabla_{\beta_\eta}^T L(\beta_\eta) \Big|_{\beta_\eta^{(i)}}$$

are the negative gradient and negative Hessian of the log density, with

$$L(\beta_\eta) = \log f(\beta_\eta | y, \theta, \Omega_{-\beta_\eta}, \Gamma) \quad \text{and} \quad \nabla_{\beta_\eta} = (\partial/\partial\beta_{\eta 1}, \dots, \partial/\partial\beta_{\eta p})^T.$$

Computation of likelihood derivatives is described at the end of this Appendix. In practice we use the expected values of $G(\beta_\eta^{(i)})$ for ease of computation.

Full conditional distributions for prior precisions: Prior precision parameter λ_η ($\eta = \xi, \nu$) has the following conditional distribution

$$f(\lambda_\eta | y, \theta, \Omega_{-\lambda_\eta}, \Gamma) \propto f(\beta_\eta | \lambda_\eta) f(\lambda_\eta | a_\eta, b_\eta).$$

By construction, since the Gamma distribution is a conjugate prior for the precision of a Gaussian distribution, we know that the full conditional distribution is also Gamma, with updated parameters

$$\hat{a}_\eta = a_\eta + \frac{p_\eta}{2} \quad \text{and} \quad \hat{b}_\eta = b_\eta + \frac{1}{2} \beta_\eta^T Q_\eta \beta_\eta.$$

A.3. Derivatives of the posterior distribution

Here we find the derivatives of the log posterior distribution, required for maximum likelihood and mMALA inference. The log likelihood of the observed data under the generalised Pareto distribution is

$$\begin{aligned} L(y | \Omega) &= \sum_{i=1}^N [-\log\left(\frac{\nu_i}{1+\xi_i}\right)] & \text{for } \xi_i \neq 0 \\ &= \left\{ -\left(\frac{1}{\xi_i} + 1\right) \log\left(1 + \frac{\xi_i}{\nu_i}(1 + \xi_i)y_i\right) \right\} \\ &\quad \sum_{i=1}^N [-\log\left(\frac{\nu_i}{1+\xi_i}\right) - \frac{(1+\xi_i)y_i}{\nu_i}] & \text{for } \xi_i = 0. \end{aligned}$$

The log conditional distribution for the vector of basis coefficients β_η ($\eta = \xi, \nu$) is then the sum of this likelihood plus a contribution from the prior distribution

$$L(\beta_\eta) = \log f(\beta_\eta | y, \theta, \Omega_{-\beta_\eta}, \Gamma) = L(y | \Omega) - \frac{\lambda_\eta}{2} \beta_\eta^T Q_\eta \beta_\eta.$$

We note the equivalence between this expression and the penalised (negative log) likelihood used for maximum likelihood inference. The gradient of the log conditional distribution for vector of coefficients β_η ($\eta = \xi, \nu$) is

$$\nabla_{\beta_\eta} L(\beta_\eta) = \nabla_{\beta_\eta} L(y | \Omega) - \lambda_\eta Q_\eta \beta_\eta.$$

Using the chain rule, the likelihood gradient can be computed as

$$\nabla_{\beta_\eta} L(y | \Omega) = (\nabla_{\beta_\eta} (B_\eta \beta_\eta))^T (\nabla_{\beta_\eta} L(y | \Omega)) = B_\eta^T (\nabla_{\beta_\eta} L(y | \Omega)).$$

The components of $\nabla_{\xi} L(y | \Omega)$ are computed as

$$\frac{\partial}{\partial \xi_i} L(y) = \begin{cases} -\frac{1}{\xi_i^2 G_i} (1 - 2\xi_i)(G_i - 1) + \frac{1}{1+\xi_i} + \frac{1}{\xi_i} \log(G_i) & \text{for } \xi_i \neq 0 \\ -\frac{y_i}{\nu_i} + \frac{1}{1+\xi_i} & \text{for } \xi_i = 0 \end{cases}$$

where $G_i = 1 + \frac{\xi_i}{\nu_i}(1 + \xi_i)y_i$, and the components of $\nabla_{\nu} L(y | \Omega)$ are

$$\frac{\partial}{\partial \nu_i} L(y) = \begin{cases} \frac{1}{\nu_i} \left(1 - \left(\frac{1}{\xi_i} + 1\right) \frac{G_i - 1}{G_i}\right) & \text{for } \xi_i \neq 0 \\ \frac{1}{\nu_i} \left(1 - \frac{G_i - 1}{\xi_i}\right) & \text{for } \xi_i = 0. \end{cases}$$

Differentiating $\nabla_{\beta_\eta} L(\beta_\eta)$ ($\eta = \xi, \nu$) again gives the Hessian matrix

$$\nabla_{\beta_\eta} \nabla_{\beta_\eta}^T L(\beta_\eta) = \nabla_{\beta_\eta} \nabla_{\beta_\eta}^T L(y | \Omega) - \lambda_\eta Q_\eta.$$

Applying the chain rule

$$\nabla_{\beta_\eta} \nabla_{\beta_\eta}^T L(y | \Omega) = B_\eta^T (\nabla_{\beta_\eta} \nabla_{\beta_\eta}^T L(y | \Omega)) B_\eta.$$

Note that the components of $\nabla_{\eta} L(y | \Omega)$ and $(\nabla_{\eta} \nabla_{\eta}^T L(y | \Omega))$ are computed separately for $\eta = \xi$ and $\eta = \nu$. Further, the expected values of likelihood second derivatives with respect to ξ and ν are

$$-\mathbb{E}_Y \left[\frac{\partial^2}{\partial \xi_i \partial \xi_j} L(y | \Omega) \right] = \begin{cases} \frac{1}{(1+\xi_i)^2} & \text{for } i = j \\ 0 & \text{for } i \neq j \end{cases}$$

and

$$-\mathbb{E}_Y \left[\frac{\partial^2}{\partial \nu_i \partial \nu_j} L(y | \Omega) \right] = \begin{cases} \frac{1}{\nu_i^2 (1 + 2\xi_i)} & \text{for } i = j \\ 0 & \text{for } i \neq j \end{cases}$$

such that Hessian matrices are diagonal. Moreover, the expectations of all of the cross derivatives $\frac{\partial^2}{\partial \xi_i \partial \nu_j} L(y|\Omega)$ are zero, since estimates of ξ and ν are asymptotically independent by construction (e.g. Chavez-Demoulin and Davison, 2005).

References

- Anderson, C., Carter, D., Cotton, P., 2001. Wave Climate Variability and Impact on Offshore Design Extremes. Report commissioned from the University of Sheffield and Satellite Observing Systems for Shell International.
- Anderson, T.W., 1962. On the distribution of the two-sample Cramer–von Mises criterion. *Ann. Math. Stat.* 33, 1148–1159.
- Behrens, C.N., Lopes, H.F., Gamerman, D., 2004. Bayesian analysis of extreme events with threshold estimation. *Stat. Model.* 4, 227–244.
- Brooks, S.P., Gelman, A., 1998. General methods for monitoring convergence of iterative simulations. *J. Comput. Graph. Stat.* 7, 434–455.
- Carter, D.J.T., Challenor, P.G., 1981. Estimating return values of environmental parameters. *Q. J. R. Meteorol. Soc.* 107, 259.
- Chavez-Demoulin, V., Davison, A., 2005. Generalized additive modelling of sample extremes. *J. R. Statist. Soc. C* 54, 207–222.
- Chavez-Demoulin, V., Embrechts, P., 2006. Smooth extremal models in finance and insurance. *J. Risk Insur.* 71, 183–199.
- Coles, S., Walshaw, D., 1994. Directional modelling of extreme wind speeds. *Appl. Stat.* 43, 139–157.
- Cooley, D., Naveau, P., Jomelli, V., Rabatel, A., Grancher, D., 2006. A Bayesian hierarchical extreme value model for lichenometry. *Environmetrics* 17, 555–574.
- Cox, D.R., Reid, N., 1987. Parameter orthogonality and approximate conditional inference. *J. R. Stat. Soc. B* 49, 1–39.
- Davison, A., Smith, R.L., 1990. Models for exceedances over high thresholds. *J. R. Stat. Soc. B* 52, 393.
- Davison, A.C., 2003. *Statistical Models*. Cambridge University Press, Cambridge.
- Eilers, P.H.C., Marx, B.D., 2010. Splines, knots and penalties. *Wiley Intersci. Rev.: Comput. Stat.* 2, 637–653.
- Fawcett, L., Walshaw, D., 2006. A hierarchical model for extreme wind speeds. *J. R. Stat. Soc. C* 55, 631–646.
- Fushiki, T., Komaki, F., Aihara, K., 2005. Nonparametric bootstrap prediction. *Bernoulli* 11, 293–307.
- Gamerman, D., Lopes, H.F., 2006. *Markov Chain Monte Carlo: Stochastic Simulation for Bayesian Inference*. Taylor and Francis, London.
- Geyer, C.J., 1992. Practical Markov chain Monte Carlo. *Stat. Sci.* 7, 473–483.
- Girolami, M., Calderhead, B., 2011. Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *J. R. Stat. Soc. B* 73, 123–214.
- Green, P.J., Silverman, B., 1994. *Nonparametric Regression and Generalised Linear Models: A Roughness Penalty Approach*. Chapman and Hall, London, UK.
- Jonathan, P., Ewans, K.C., 2013. Statistical modelling of extreme ocean environments with implications for marine design: a review. *Ocean Eng.* 62, 91–109.
- Jonathan, P., Ewans, K.C., Forristall, G.Z., 2008. Statistical estimation of extreme ocean environments: the requirement for modelling directionality and other covariate effects. *Ocean Eng.* 35, 1211–1225.
- Jonathan, P., Randell, D., Ewans, K., 2013. Joint modelling of extreme ocean environments incorporating covariate effects. *Coast. Eng.* 79, 22–31.
- Jonathan, P., Randell, D., Wu, Y., Ewans, K., 2014. Return level estimation from non-stationary spatial data exhibiting multidimensional covariate effects. *Ocean Eng.* 88, 520–532.
- MacDonald, A., Scarrott, C.J., Lee, D., Darlow, B., Reale, M., Russell, G., 2011. A flexible extreme value mixture model. *Comput. Stat. Data Anal.* 55, 2137–2157.
- MacKay, D., 1998. Introduction to Gaussian processes. In: Bishop, C. (Ed.), *Neural Networks and Machine Learning*. Springer-Verlag, Berlin, pp. 84–92.
- Mendez, F.J., Menendez, M., Luceno, A., Medina, R., Graham, N.E., 2008. Seasonality and duration in extreme value distributions of significant wave height. *Ocean Eng.* 35, 131–138.
- Northrop, P., Jonathan, P., 2011. Threshold modelling of spatially-dependent non-stationary extremes with application to hurricane-induced wave heights. *Environmetrics* 22, 799–809.
- Northrop, P., Jonathan, P., Randell, D., 2016. Threshold modeling of nonstationary extremes. In: Dey, D., Yan, J. (Eds.), *Extreme Value Modeling and Risk Analysis: Methods and Applications*. Chapman and Hall/CRC, pp. 87–108.
- Perez-Cruz, F., July 2008. Kullback–Leibler divergence estimation of continuous distributions. In: *IEEE International Symposium on Information Theory, 2008 (ISIT 2008)*, pp. 1666–1670.
- Randell, D., Feld, G., Ewans, K., Jonathan, P., 2015a. Distributions of return values for ocean wave characteristics using directional-seasonal extreme value analysis. *Environmetrics* 26, 442–450.
- Randell, D., Turnbull, K., Ewans, K., Jonathan, P., 2015b. Bayesian inference for non-stationary marginal extremes, draft at (www.lancs.ac.uk/~jonathan) (submitted to *Environmetrics* August 2015).
- Randell, D., Zanini, E., Vogel, M., Ewans, K., Jonathan, P., 2014. Omnidirectional return values for storm severity from directional extreme value models: the effect of physical environment and sample size. In: *Proceedings of 33rd International Conference on Ocean, Offshore and Arctic Engineering*, San Francisco OMAE2014-23156.
- Rasmussen, C., Williams, C., 2006. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA.
- Renard, B., Lang, M., Bois, P., 2006. Statistical Analysis of Extreme Events in a Nonstationary Context via a Bayesian Framework. Case Study with Peak-Over-Threshold Data. *Stochastic Environmental Research and Risk Assessment*, vol. 21, pp. 97–112.
- Ruppert, D., Wand, M.P., Carroll, R.J., 2003. *Semiparametric Regression*. Cambridge University Press, Cambridge.
- Scarrott, C., MacDonald, A., 2012. A review of extreme value threshold estimation and uncertainty quantification. *Revstat* 10, 33–60.
- Scotto, M., Guedes-Soares, C., 2000. Modelling the long-term time series of significant wave height with non-linear threshold models. *Coast. Eng.* 40, 313–327.
- Smith, R.L., Naylor, J.C., 1987. A comparison of maximum likelihood and Bayesian estimators for the three-parameter Weibull distribution. *J. R. Stat. Soc. C* 36, 358–369.