



CrossMark
click for updates

Research

Cite this article: Roberts S, Osborne M, Ebden M, Reece S, Gibson N, Aigrain S. 2013 Gaussian processes for time-series modelling. *Phil Trans R Soc A* 371: 20110550. <http://dx.doi.org/10.1098/rsta.2011.0550>

One contribution of 17 to a Discussion Meeting Issue ‘Signal processing and inference for the physical sciences’.

Subject Areas:

statistics, mathematical modelling, applied mathematics, pattern recognition

Keywords:

Gaussian processes, time-series analysis, Bayesian modelling

Author for correspondence:

S. Roberts

e-mail: sjrob@robots.ox.ac.uk

Gaussian processes for time-series modelling

S. Roberts¹, M. Osborne¹, M. Ebden¹, S. Reece¹,
N. Gibson² and S. Aigrain²

¹Department of Engineering Science, and

²Department of Astrophysics, University of Oxford,
Oxford OX1 3PU, UK

In this paper, we offer a gentle introduction to Gaussian processes for time-series data analysis. The conceptual framework of Bayesian modelling for time-series data is discussed and the foundations of Bayesian non-parametric modelling presented for *Gaussian processes*. We discuss how domain knowledge influences design of the Gaussian process models and provide case examples to highlight the approaches.

1. Introduction

If we are to take full advantage of the richness of scientific data available to us, we must consider a principled framework under which we may reason and infer. To fail to do this is to ignore uncertainty and risk false analysis, decision-making and forecasting. What we regard as a prerequisite for intelligent data analysis is ultimately concerned with the problem of computing in the presence of uncertainty. Considering data analysis under the mathematics of modern probability theory allows us to exploit a profound framework under which information, uncertainty and risk for actions, events and outcomes may be readily defined. Much recent research hence focuses on the principled handling of uncertainty for modelling in environments that are dynamic, noisy, observation costly and time sensitive. The machinery of probabilistic inference brings to the field of time-series analysis and monitoring robust, stable, computationally practical and principled approaches that naturally accommodate these real-world challenges. As a framework for reasoning in the presence of uncertain, incomplete and delayed information, we appeal to Bayesian inference. This allows us to perform robust modelling even in highly uncertain situations, and has a long pedigree in inference. Being able to

include measures of uncertainty allows, for example, us to actively select where and when we would like to observe samples and offers approaches by which we may readily combine information from multiple noisy sources.

This paper favours the conceptual over the mathematical (of course, the mathematical details are important and elegant but would obscure the aims of this paper; the interested reader is encouraged to read the cited material and a canonical text such as Rasmussen & Williams [1]). We start §2 with a short overview of why *Bayesian* modelling is important in time-series analysis, culminating in arguments that provoke us to use non-parametric models. Section 3 presents a conceptual overview of a particular flavour of non-parametric model, the Gaussian process (GP), which is well suited to time-series modelling [1]. We discuss in more detail the role of *covariance functions*, the influence they have on our models and explore, by example, how the (apparently subjective) function choices we make are in fact motivated by domain knowledge. Section 5 presents real-world time-series examples, from sensor networks, changepoint data and astronomy, to highlight the practical application of GP models. The more mathematical framework of inference is detailed in §4.

2. Bayesian time-series analysis

We start by casting time-series analysis into the format of a *regression* problem, of the form $y(x) = f(x) + \eta$, in which $f()$ is a (typically) unknown function and η is a (typically white) additive noise process. The goal of inference in such problems is twofold: firstly to evaluate the putative form of $f()$ and secondly to evaluate the probability distribution of y_* for some x_* , i.e. $p(y_* | x_*)$. To enable us to perform this inference, we assume the existence of a dataset of *observations*, typically obtained as input–output pairs, $\mathbb{D} = (x_i, y_i)$ for example. For the purposes of this study, we make the tacit assumption that the inputs x_i (representing, e.g. time locations of samples) are known precisely, i.e. there is no *input noise*, but that observation noise is present on the y_i . When we come to analyse time-series data, there are two approaches we might consider. The first *function mapping* and the second *curve fitting*.

The mapping approach considers inference of a function f that maps some observed x to an outcome variable y *without* explicit reference to the (time) ordering of the data. For example, if we choose x to be a datum in a time series, and y to be the next datum, then inferring $f(x)$ models the relationship between one datum and its successor. Problematically, the mapping is (typically) static, so poorly models non-stationary time series, and there is difficulty in incorporating time-series domain knowledge, such as beliefs about smoothness and continuity. Furthermore, if the periods between samples are uneven, this approach fails to accommodate this knowledge with ease.

Curve fitting, on the other hand, makes the tacit assumption that y is ordered by x , the latter normally taken to be the time variable, with inference proceeding by fitting a curve to the set of x, y points. Prediction, for example, is thence achieved by extrapolating the curve that models the observed past data. The relationship between x and y is hence not fixed, but conditioned on observed data that (typically) lies close, in time, to the point we are investigating. In this study, we make the decision to concentrate on this approach, as we believe it offers a more profound model for much of the time-series data we are concerned with.

As a simple example to introduce the canonical concepts of Bayesian modelling, we consider a small set of data samples, located at $x = 0, 1, 2$, and associated observed target values. Least-squares regression on this data using a simple model (based on polynomial splines) gives rise to the curve shown as the line in figure 1a. We see that, naturally, this curve fits our observed data very well. What about the credibility of the model in regions where we see no data, importantly $x > 2$? If we look at a larger set of example curves from the same model, we obtain a family of curves that explains the observed data *identically* yet differ very significantly in regions where we have no observations, both interpolating between sample points, and in extrapolation. This simple example leads naturally to us considering a *distribution of curves*. Working with some

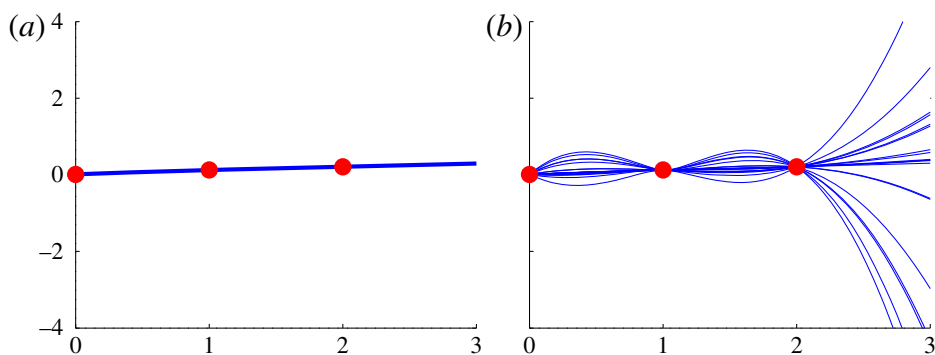


Figure 1. A simple example of curve fitting. (a) The least-squares fit of a simple spline to the observed data (circles). (b) Example curves with identical fit to the data as the least-squares spline. These curves have high similarity close to the data yet high variability in regions of no observations, both interpolating and, importantly for time series, as we extrapolate beyond $x = 2$. (Online version in colour.)

distribution over the curves, each of which offers an explanation for the observed data, is central to Bayesian modelling. We note that curves which lie towards the edges of this distribution have higher average curvature than those which lie close to the middle. In the simple example under consideration, there is an intimate relationship between curvature, complexity and Bayesian inference, leading naturally to posterior beliefs over models being a combination of how well observed data are explained and how complex the explanatory functions are. This elegant formalism encodes in a single mathematical framework such ideas as *Occam's razor*, such that simple explanations of observed data are favoured.

(a) Parametric and non-parametric models

The simple example above showed that there are many functions which can equally well explain data that we have observed. How should we choose from the bewildering array of mathematical functions that give rise to such explanatory curves? If we have strong prior knowledge regarding a system, then this (infinite-dimensional) function space may be reduced to a single family; perhaps the family of quartic polynomials may be the right choice. Such models are considered to be *parametric*, in the sense that a finite number of unknown parameters (in our polynomial example, these are the coefficients of the model) need to be inferred as part of the data modelling process. Although there is a very large literature (rightly so) on such parametric modelling methods, there are many scenarios in which we have little, or no, prior knowledge regarding appropriate models to use. We may, however, have seemingly less specific domain knowledge; for example, we may know that our observations are visible examples from an underlying process that is smooth, continuous and variations in the function take place over characteristic time scales (not too slowly yet not so fast) and have typical amplitude. Surprisingly, we may work mathematically with the infinite space of all functions that have these characteristics. Furthermore, we may even contemplate probability distributions over this function space, such that the work of modelling, explaining and forecasting data is performed by refining these distributions, so focusing on regions of the function space that are excellent contenders to model our data. As these functions are not characterized with explicit sets of parameters to be inferred (unlike our simple polynomial example, in which sets of coefficients need to be evaluated), this approach is referred to as a branch of *non-parametric* modelling.¹

¹This always feels rather disingenuous though, as these models do have *hyperparameters*, which we discuss later in this paper. These still need to be inferred! They are referred to as *hyperparameters*, as they govern such things as the scale of a distribution rather than acting explicitly on the functional form of the curves.

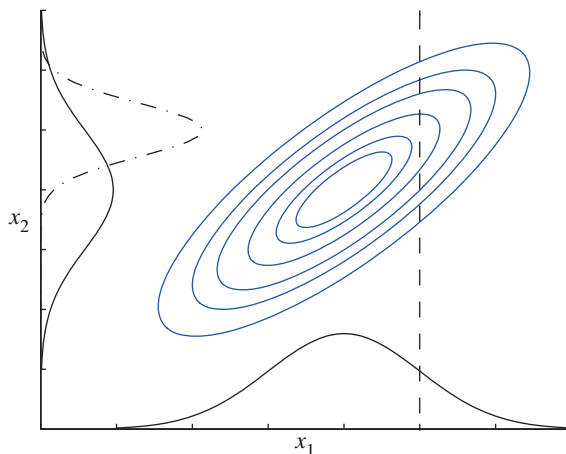


Figure 2. The conceptual basis of GPs starts with an appeal to simple multi-variate Gaussian distributions. A joint distribution (covariance ellipse) forms marginal distributions $p(x_1)$, $p(x_2)$ that are vague (black solid line). Observing x_1 at a value indicated by the vertical dashed line changes our beliefs about x_2 , giving rise to a conditional distribution (black dashed-dot line). Knowledge of the covariance lets us shrink uncertainty in one variable, based on observation of the other. (Online version in colour.)

As the dominant machinery for working with these models is that of probability theory, they are often referred to as *Bayesian non-parametric models*. We now focus on a particular member, namely the GP.

3. Gaussian processes

We start this introduction to GPs by considering a simple two-variable Gaussian distribution, which is defined for variables x_1, x_2 say, by a mean and a 2×2 covariance matrix, which we may visualize as a covariance ellipse corresponding to equal probability contours of the joint distribution $p(x_1, x_2)$. Figure 2 shows an example of a two-dimensional distribution as a series of elliptical contours. The corresponding *marginal* distributions $p(x_1)$ and $p(x_2)$ are shown as ‘projections’ of this along the x_1 and x_2 axes (solid black lines). We now consider the effect of observing one of the variables such that, for example, we observe x_1 at the location of the dashed vertical line in the figure. The resultant *conditional distribution* $p(x_2 | x_1 = \text{known})$, indicated by the dash-dotted curve, now deviates significantly from the marginal $p(x_2)$. Because of the relationship between the variables implied by the covariance, knowledge of one shrinks our uncertainty in the other.

To see the intimate link between this simple example and time-series analysis, we represent the same effect in a different format. Figure 3 shows the mean (black line) and $\pm\sigma$ (grey-shaded region) for $p(x_1)$ and $p(x_2)$. Figure 3a depicts our initial state of ignorance and figure 3b after we observe x_1 . Note how the observation changes the location and uncertainty of the distribution over x_2 . Why stop at only two variables? We can extend this example to arbitrarily large numbers of variables, the relationships between which are defined by an ever larger covariance. Figure 4 shows the posterior distribution for a 10 day example in which observations are made at locations 2, 6 and 8. Figure 4a shows the posterior mean and $\pm\sigma$ as in our previous examples. Figure 4b extends the posterior distribution evaluation densely in the same interval (here, we evaluate the distribution over several hundred points). We note that the ‘discrete’ distribution is now rather continuous. In principle, we can extend this procedure to the limit in which the locations of the x_i are infinitely dense (here, on the real line) and so the infinite joint distribution over them all is equivalent to a distribution over a function space. In practice, we will not need to work with

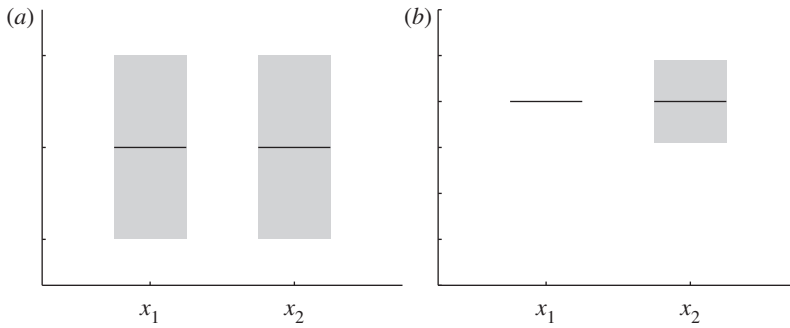


Figure 3. The change in distributions on x_1 and x_2 is here presented in a form more familiar to time-series analysis. (a) The initial, vague, distributions (the black line showing the mean and the grey shading $\pm\sigma$) and (b) subsequent to observing x_1 . The distribution over x_2 has become less uncertain and the most-likely ‘forecast’ of x_2 has also shifted.

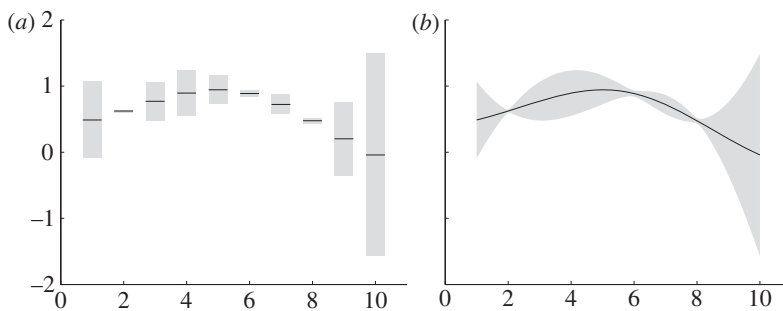


Figure 4. (a) The posterior distribution (the black line showing the mean and the grey shading $\pm\sigma$) for a 10 day example, with observations made at locations 2, 6 and 8. (b) Evaluates the posterior densely in the interval [1,10] showing how arbitrarily dense evaluation gives rise to a ‘continuous’ posterior distribution with time.

such infinite spaces, it is sufficient that we can choose to evaluate the probability distribution over the function at *any location on the real line* and that we incorporate any observations we may have at any other points. We note, crucially, that the locations of observations and points we wish to investigate the function are *not constrained* to lie on any predefined sample points; hence, we are working in continuous time with a GP.

(a) Covariance functions

As we have seen, the covariance forms the beating heart of GP inference. How do we formulate a covariance over arbitrarily large sets? The answer lies in defining a *covariance kernel function*, $k(x_i, x_j)$, which provides the covariance element between any two (arbitrary) sample locations, x_i and x_j say. For a set of locations $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$, we hence may define the *covariance matrix* as

$$\mathbf{K}(\mathbf{x}, \mathbf{x}) = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \cdots & k(x_n, x_n) \end{pmatrix}. \quad (3.1)$$

This means that the entire function evaluation, associated with the points in \mathbf{x} , is a draw from a multi-variate Gaussian distribution,

$$p(\mathbf{y}(\mathbf{x})) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \mathbf{K}(\mathbf{x}, \mathbf{x})), \quad (3.2)$$

where $\mathbf{y} = \{y_1, y_2, \dots, y_n\}$ are the dependent function values, evaluated at locations x_1, \dots, x_n , and μ is a *mean function*, again evaluated at the locations of the x variables (which we will briefly revisit later). If we believe that there is noise associated with the observed function values, y_i , then we may fold this noise term into the covariance. As we expect noise to be uncorrelated from sample to sample in our data, so the noise term adds only to the diagonal of \mathbf{K} , giving a modified covariance for noisy observations of the form

$$\mathbf{V}(\mathbf{x}, \mathbf{x}) = \mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I}, \quad (3.3)$$

where \mathbf{I} is the identity matrix and σ^2 is a *hyperparameter* representing the noise variance.

How do we evaluate the GP posterior distribution at some test datum, x_* say? We start with considering the joint distribution of the observed data \mathbb{D} (consisting of \mathbf{x} and associated values \mathbf{y}) augmented by x_* and y_* ,

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mu(\mathbf{x}) \\ \mu(x_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\mathbf{x}, \mathbf{x}) & \mathbf{K}(\mathbf{x}, x_*) \\ \mathbf{K}(x_*, \mathbf{x}) & k(x_*, x_*) \end{bmatrix}\right), \quad (3.4)$$

where $\mathbf{K}(\mathbf{x}, x_*)$ is the column vector formed from $k(x_1, x_*), \dots, k(x_n, x_*)$ and $\mathbf{K}(x_*, \mathbf{x})$ is its transpose. We find, after some manipulation, that the posterior distribution over y_* is Gaussian with mean and variance given by

$$m_* = \mu(x_*) + \mathbf{K}(x_*, \mathbf{x})\mathbf{K}(\mathbf{x}, \mathbf{x})^{-1}(\mathbf{y} - \mu(\mathbf{x})) \quad (3.5)$$

and

$$\sigma_*^2 = K(x_*, x_*) - \mathbf{K}(x_*, \mathbf{x})\mathbf{K}(\mathbf{x}, \mathbf{x})^{-1}\mathbf{K}(\mathbf{x}, x_*). \quad (3.6)$$

We may readily extend this to infer the GP at a set of locations outside our observations, at \mathbf{x}_* say, to evaluate the posterior distribution of $\mathbf{y}(\mathbf{x}_*)$. The latter is readily obtained once more by extending the above equations and using standard results for multi-variate Gaussians. We obtain a posterior mean and variance given by

$$p(\mathbf{y}_*) = \mathcal{N}(\mathbf{m}_*, \mathbf{C}_*), \quad (3.7)$$

where

$$\mathbf{m}_* = \mu(\mathbf{x}_*) + \mathbf{K}(\mathbf{x}_*, \mathbf{x})\mathbf{K}(\mathbf{x}, \mathbf{x})^{-1}(\mathbf{y} - \mu(\mathbf{x})) \quad (3.8)$$

and

$$\mathbf{C}_* = \mathbf{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{K}(\mathbf{x}_*, \mathbf{x})\mathbf{K}(\mathbf{x}, \mathbf{x})^{-1}\mathbf{K}(\mathbf{x}, \mathbf{x}_*)^T, \quad (3.9)$$

in which we use the shorthand notation for the covariance, $\mathbf{K}(\mathbf{a}, \mathbf{b})$, defined as

$$\mathbf{K}(\mathbf{a}, \mathbf{b}) = \begin{pmatrix} k(a_1, b_1) & k(a_1, b_2) & \cdots & k(a_1, b_n) \\ k(a_2, b_1) & k(a_2, b_2) & \cdots & k(a_2, b_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(a_n, b_1) & k(a_n, b_2) & \cdots & k(a_n, b_n) \end{pmatrix}. \quad (3.10)$$

If we believe (and in most situations we do) that the observed data are corrupted by a noise process, we would replace the $\mathbf{K}(\mathbf{x}, \mathbf{x})$ term above with, for example, $\mathbf{V}(\mathbf{x}, \mathbf{x})$ from equation (3.3) above.

What should the functional form of the kernel function $k(x_i, x_j)$ be? To answer this, we will start by considering what the covariance elements indicate. In our simple two-dimensional example, the off-diagonal elements define the correlation between the two variables. By considering time series in which we believe the informativeness of past observations, in explaining current data, is a function of how long ago we observed them, we then obtain *stationary* covariance functions that are dependent on $|x_i - x_j|$. Such covariance functions can be represented as the Fourier transform

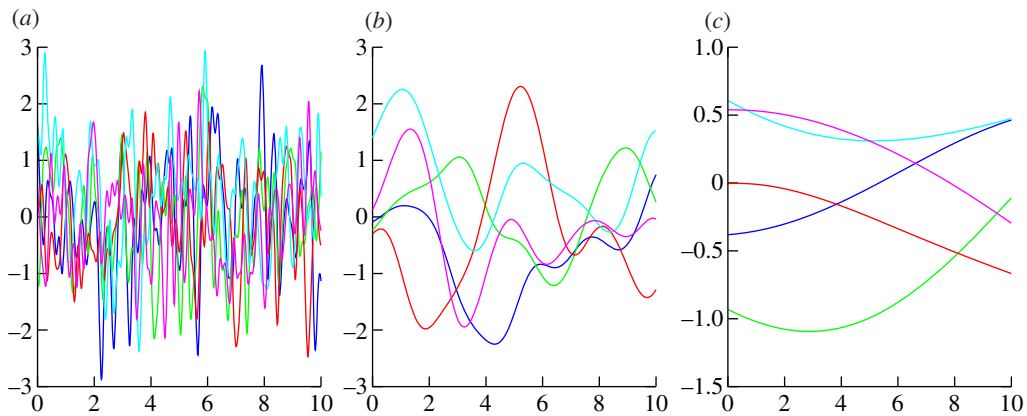


Figure 5. (a–c) Functions drawn from a GP with a squared exponential covariance function with output scale $h = 1$ and length scales $\lambda = 0.1$ (a), 1 (b), 10 (c). (Online version in colour.)

of a normalized probability density function (via Bochner’s theorem [1]); this density can be interpreted as the spectral density of the process. The most widely used covariance function of this class is arguably the squared exponential (SE) function, given by

$$k(x_i, x_j) = h^2 \exp \left[- \left(\frac{x_i - x_j}{\lambda} \right)^2 \right]. \quad (3.11)$$

In equation (3.11), we see two more hyperparameters, namely h, λ , which respectively govern the output scale of our function and the input, or time, scale. The role of inference in GP models is to refine vague distributions over many, very different curves, to more precise distributions that are focused on curves that explain our observed data. As the form of these curves is uniquely controlled by the hyperparameters, so, in practice, inference proceeds by refining distributions over them. As h controls the gain, or magnitude, of the curves, we set this to $h = 1$ to generate figure 5, which shows curves drawn from a GP (with an SE covariance function) with varying $\lambda = 0.1, 1, 10$ (figure 5a–c). The important question of *how* we infer the hyperparameters is left until later in this paper, in §4. We note that to be a valid covariance function, $k()$, implies only that the resultant covariance matrix, generated using the function, is guaranteed to be positive (semi-)definite. As a simple example, figure 6a shows a small sample of six observed data points, shown as dots, along with error bars associated with each. The seventh datum, with ‘?’ beneath it, is unobserved. We fit a GP with the SE covariance kernel (equation (3.11)). Figure 6b shows the GP posterior mean (black curve) along with $\pm 2\sigma$ (the posterior standard deviation). Although only a few samples are observed, corresponding to the set of \mathbf{x}, \mathbf{y} of equations (3.8) and (3.9), we here evaluate the function on a fine set of points, evaluating the corresponding y_* posterior mean and variance using these equations and hence providing interpolation between the noisy observations (this explains the past) and extrapolation for $x_* > 0$ that predicts the future. In this simple example, we have used a ‘simple’ covariance function. As the sum of valid covariance functions is itself a valid covariance function (more on this in §3c(i) later) so we may entertain more complex covariance structures, corresponding to our prior belief regarding the data. Figure 7 shows GP modelling of observed (noisy) data for which we use slightly more complex covariances. Figure 7a shows data modelled using a sum of SE covariances, one with a bias towards shorter characteristic time scales than the other. We see how this combination elegantly allows us to model a system with both long- and short-term dynamics. Figure 7b uses an SE kernel, with bias towards longer time-scale dynamics, along with a periodic component kernel (which we will discuss in more detail in §3c(ii)). Note here how extrapolation outside the data indicates a strong posterior belief regarding the continuance of periodicity.

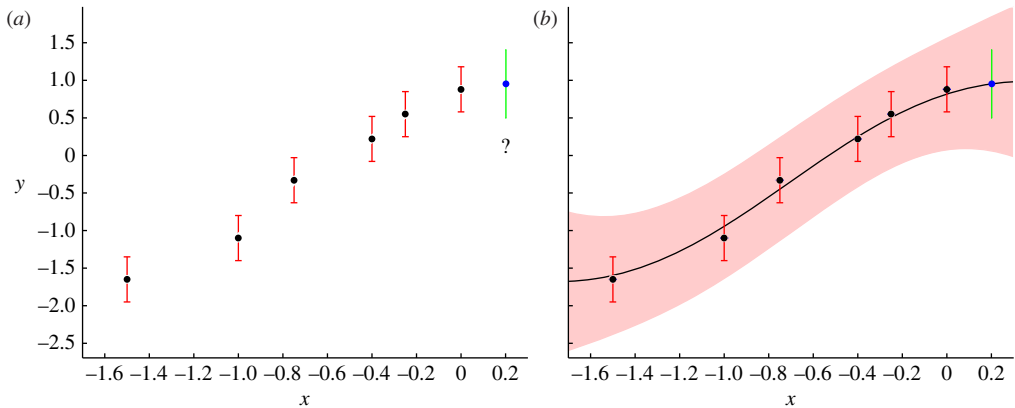


Figure 6. (a) Given six noisy data points (error bars are indicated with vertical lines), we are interested in estimating a seventh at $x_* = 0.2$. (b) The solid line indicates an estimation of y_* for x_* across the range of the plot. Along with the posterior mean, the posterior uncertainty, $\pm 2\sigma$, is shaded. (Online version in colour.)

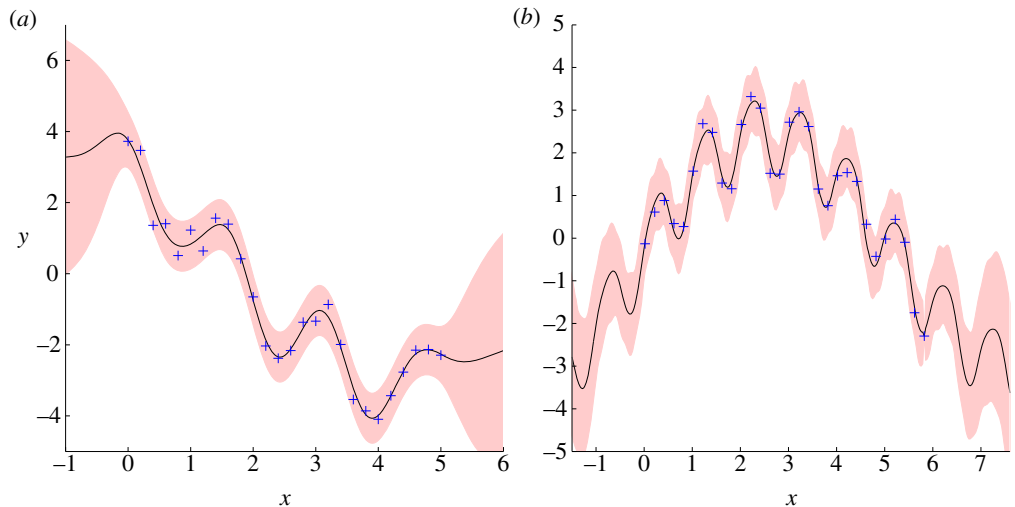


Figure 7. (a) Estimation of y_* (solid line) and $\pm 2\sigma$ posterior deviance for a function with short-term and long-term dynamics, and (b) long-term dynamics and a periodic component. Observations are shown as pluses. As in the previous example, we evaluate the posterior GP over an extended range to show both interpolation and extrapolation. (Online version in colour.)

(b) Sequential modelling and active data selection

We start by considering a simple example, shown in figure 8. Figure 8a shows a set of data points and the GP posterior distribution *excluding* observation of the right-most datum (darker shaded point). Figure 8b depicts the same inference *including* this last datum. We see how the posterior variance shrinks as we make the observation. The previous example showed how making an observation, even of a noisy time series, shrinks our uncertainty associated with beliefs about the function local to the observation. We can see this even more clearly if we successively extrapolate until we see another datum, as shown in figure 9. Rather than observations coming on a fixed time-interval grid, we can imagine a scenario in which observations are costly to acquire, and we wish to find a natural balance between sampling and reducing uncertainty in the functions of interest. This concept leads us naturally in two directions. Firstly, for the active *requesting* of observations when our uncertainty has grown beyond acceptable limits (of course these limits are

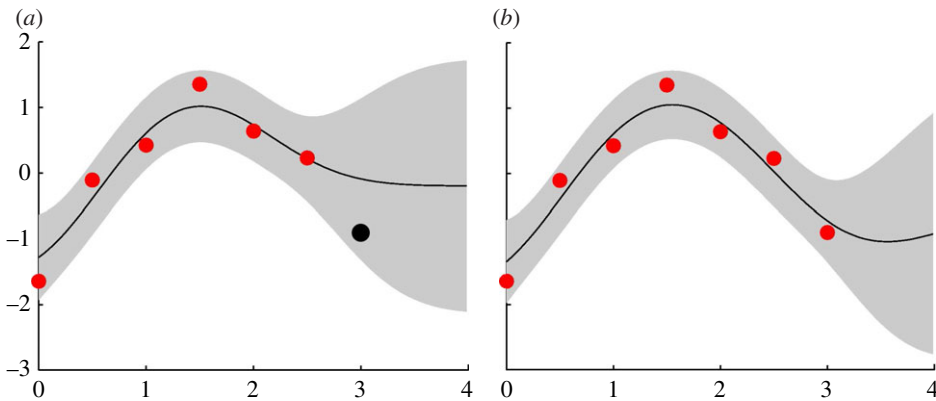


Figure 8. A simple example of a GP applied sequentially. (a) The posterior mean and $\pm 2\sigma$ prior to observing the right-most datum (darker shaded) and (b) after observation. (Online version in colour.)

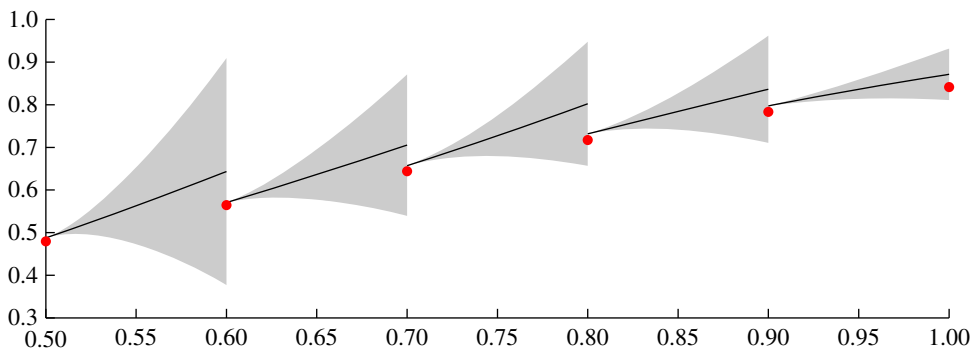


Figure 9. The GP is run sequentially making forecasts until a new datum is observed. Once we make an observation, the posterior uncertainty drops to zero (assuming noiseless observations). (Online version in colour.)

related to the cost of sampling and observation and the manner in which uncertainty in the time series can be balanced against this cost) and secondly to dropping previously observed samples from our model. The computational cost of GPs is dominated by the inversion of a covariance matrix (as in equation (3.9)) and hence scales with the cube of the number of retained samples. This leads to an adaptive *sample retention*. Once more, the balance is problem specific, in that it relies on the trade-off between computational speed and (for example) forecasting uncertainty. The interested reader is pointed to Osborne *et al.* [2] for more detailed discussions. We provide some examples of active data selection in operation in real problem domains later in this study.

(c) Choosing covariance and mean functions

The prior mean of a GP represents whatever we expect for our function before seeing any data. The covariance function of a GP specifies the correlation between any pair of outputs. This can then be used to generate a covariance matrix over our set of observations and predictants. Fortunately, there exist a wide variety of functions that can serve in this purpose [3,4], which can then be combined and modified in a further multitude of ways. This gives us a great deal of flexibility in our modelling of functions, with covariance functions available to model periodicity, delay, noise and long-term drifts and other phenomena.

(i) Covariance functions

In the following section, we briefly describe commonly used kernels. We start with simple white noise, and then consider common *stationary* covariances, both uni- and multi-dimensional. We finish this section with periodic and quasi-periodic kernel functions. The interested reader is referred to Rasmussen & Williams [1] for more details. Although the following list is not exclusive by any means, it provides details for most of the covariance functions suitable for analysis of time series. We note once more that sums (and products) of valid covariance kernels give valid covariance functions (i.e. the resultant covariance matrices are positive (semi-)definite) and so we may entertain with ease multiple explanatory hypotheses. The price we pay lies in the extra complexity of handling the increased number of hyperparameters.

White noise with variance σ^2 is represented by

$$k_{\text{WN}}(x_i, x_j) = \sigma^2 \delta(i, j). \quad (3.12)$$

This kernel allows us to entertain uncertainty in our observed data and is so typically found added to other kernels (as we saw in equation (3.3)).

The SE kernel is given by

$$k_{\text{SE}} = h^2 \exp \left[- \left(\frac{x_i - x_j}{\lambda} \right)^2 \right], \quad (3.13)$$

where h is an output-scale amplitude and λ is an input (length, or time) scale. This gives rather smooth variations with a typical time scale of λ and admits functions drawn from the GP that are infinitely differentiable.

The rational quadratic (RQ) kernel is given by

$$k_{\text{RQ}}(x_i, x_j) = h^2 \left(1 + \frac{(x_i - x_j)^2}{\alpha \lambda^2} \right)^{-\alpha}, \quad (3.14)$$

where α is known as the index. Rasmussen & Williams [1] show that this is equivalent to a scale mixture of SE kernels with different length scales, the latter distributed according to a Beta distribution with parameters α and λ^{-2} . This gives variations with a range of time scales, the distribution peaking around λ but extending to significantly longer period (but remaining rather smooth). When $\alpha \rightarrow \infty$, the RQ kernel reduces to the SE kernel with length scale λ .

The Matérn class of covariance functions is defined by

$$k_{\text{M}}(x_i, x_j) = h^2 \frac{1}{\Gamma(\nu) 2^{\nu-1}} \left(2\sqrt{\nu} \frac{|x_i - x_j|}{\lambda} \right) \mathbb{B}_{\nu} \left(2\sqrt{\nu} \frac{|x_i - x_j|}{\lambda} \right), \quad (3.15)$$

where h is the output scale, λ is the input scale, $\Gamma(\cdot)$ is the standard Gamma function and $\mathbb{B}(\cdot)$ is the modified Bessel function of second order. The additional hyperparameter ν controls the degree of differentiability of the resultant functions modelled by a GP with a Matérn covariance function, such that they are only $(\nu + \frac{1}{2})$ times differentiable. As $\nu \rightarrow \infty$, so the functions become infinitely differentiable and the Matérn kernel becomes the SE one. Taking $\nu = \frac{1}{2}$ gives the exponential kernel

$$k(x_i, x_j) = h^2 \exp \left(- \frac{|x_i - x_j|}{\lambda} \right), \quad (3.16)$$

which results in functions that are only once differentiable, and correspond to the Ornstein–Uhlenbeck process, the continuous time equivalent of a first-order autoregressive model, AR(1). Indeed, as discussed in Rasmussen & Williams [1], time-series models corresponding to AR(p) processes are discrete time equivalents of GP models with Matérn covariance functions with $\nu = p - \frac{1}{2}$.

Multiple inputs and outputs. The simple distance metric, $|x_1 - x_2|$, used thus far clearly allows only for the simplest case of a one-dimensional input x , which we have hitherto tacitly assumed to represent a time measure. In general, however, we assume our input space has finite dimension and write $x^{(e)}$ for the value of the e th element in \mathbf{x} and denote $x_i^{(e)}$ as the value of the e th element at

the i th index point. In such scenarios, we entertain multiple exogenous variables. Fortunately, it is not difficult to extend covariance functions to allow for these multiple input dimensions. Perhaps the simplest approach is to take a covariance function that is the product of one-dimensional covariances over each input (the *product correlation* rule [5]),

$$k(x_i, x_j) = \prod_e k^{(e)}(x_i^{(e)}, x_j^{(e)}), \quad (3.17)$$

where $k^{(e)}$ is a valid covariance function over the e th input. As the product of covariances is a covariance, so equation (3.17) defines a valid covariance over the multi-dimensional input space. We can also introduce distance functions appropriate for multiple inputs, such as the Mahalanobis distance,

$$d^{(M)}(x_i, x_j; \Sigma) = \sqrt{(x_i - x_j)^T \Sigma^{-1} (x_i - x_j)}, \quad (3.18)$$

where Σ is a covariance matrix over the input variable vector \mathbf{x} . Note that this is a matrix which represents hyperparameters of the model, and should not be confused with covariances formed from covariance functions (which are always denoted by \mathbf{K} in this study). If Σ is a diagonal matrix, its role in equation (3.18) is simply to provide an individual input scale $\lambda^e = \sqrt{\Sigma(e, e)}$ for the e th dimension. However, by introducing off-diagonal elements, we can allow for correlations among the input dimensions. To form the multi-dimensional kernel, we simply replace the scaled distance measure $|x_i - x_j|/\lambda$ of, e.g. equation (3.13) with $d^{(M)}(x_1, x_2)$ from equation (3.18) above.

For multi-dimensional outputs, we consider a multi-dimensional space consisting of a set of time series along with a label l , which indexes the time series, and x denoting time. Together, these thence form the two-dimensional set of $[l, x]$. We will then exploit the fact that a product of covariance functions is a covariance function in its own right, and write

$$k([l_m, x_i], [l_n, x_j]) = k_x(x_i, x_j) k_l(l_m, l_n),$$

taking covariance function terms over both time and time-series label. If the number of time series is not too large, we can arbitrarily represent the covariance matrix over the labels, using the spherical decomposition [6]. This allows us to represent any covariance structure over the labels. More details of this approach, which enables the dependencies between time series to be modelled, are found in Roberts *et al.* [7], and we use this as the focus of one of our examples in §5.

Periodic and quasi-periodic kernels. Note that a valid covariance function under any arbitrary (smooth) map remains a valid covariance function [1,8]. For any function $u: x \rightarrow u(x)$, a covariance function $k()$ defined over the range of x gives rise to a valid covariance $k'()$ over the domain of u . Hence, we can use simple, stationary covariances in order to construct more complex (possibly non-stationary) covariances. A particularly relevant example of this,

$$u(x) = (u^{(a)}(x), u^{(b)}(x)) = \left(\cos\left(2\pi \frac{x}{T}\right), \sin\left(2\pi \frac{x}{T}\right) \right), \quad (3.19)$$

allows us to modify our simple covariance functions above to model periodic functions. We can now take this covariance over u as a valid covariance over x . As a result, we have the covariance function, for the example of the SE (3.13),

$$k_{\text{per-SE}}(x_j, x_i; h, w, T) = h^2 \exp\left(-\frac{1}{2w^2} \sin^2\left(\pi \left|\frac{x_j - x_i}{T}\right|\right)\right). \quad (3.20)$$

In this case, the output scale h serves as the amplitude and T is the period. The hyperparameter w is a ‘roughness’ parameter that serves a role similar to the input scale λ in stationary covariances. With this formulation, we can perform inference about functions of arbitrary roughness and with arbitrary period. Indeed a periodic covariance function can be constructed from any kernel involving the squared distance $(x_i - x_j)^2$ by replacing the latter with $\sin^2[\pi(x_i - x_j)/T]$, where T is the period. The length scale w is now relative to the period, and letting $w \rightarrow \infty$ gives sinusoidal variations, while increasingly small values of w give periodic variations with

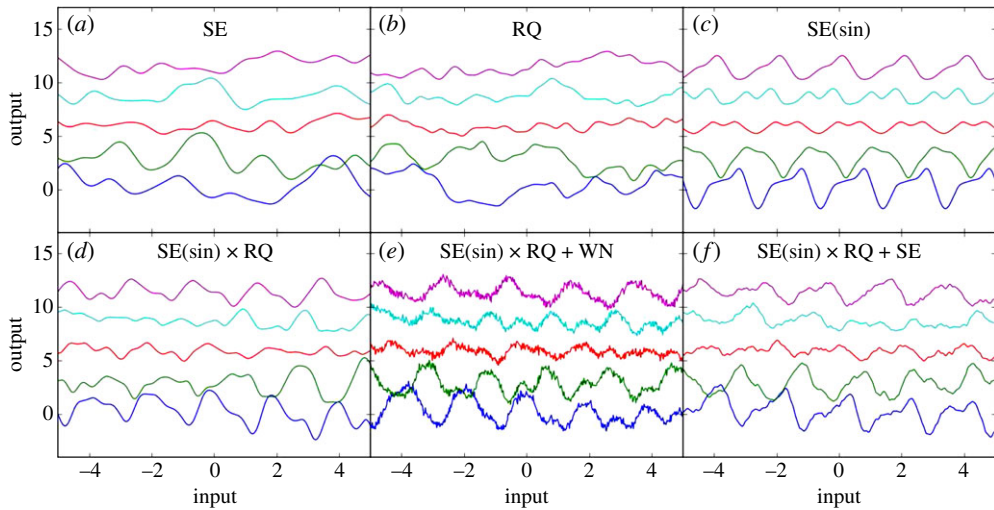


Figure 10. Random draws from GPs with different kernels. (a) Shows the SE kernel (equation (3.13), with $h = 1$, $\lambda = 1$), (b) the RQ (equation (3.14), with $h = 1$, $\lambda = 1$ and $\alpha = 0.5$) and (c) a periodic kernel based on the SE (equation (3.20), with $h = 1$, $T = 2$ and $w = 0.5$). (d) Shows a quasi-periodic kernel constructed by multiplying the periodic kernel of equation (3.13) (with $h = 1$, $T = 2$, $w = 1$) with the RQ kernel of equation (3.14) (with $\lambda = 4$ and $\alpha = 0.5$). (e, f) Show noisy versions of this kernel obtained by adding, respectively, a white noise term (equation (3.13), with $\sigma = 0.2$) and an SE term (equation (3.13), with $h = 0.1$, $\lambda = 0.1$). Each line consists of equally spaced samples over the interval $[-5, 5]$, and is offset from the previous one by 3 for clarity. The random number generated was initiated with the same seed before generating the samples shown in each panel. (Online version in colour.)

increasingly complex harmonic content. Similar periodic functions could be constructed from any kernel. Other periodic functions could also be used, so long as they give rise to a symmetric, positive definite covariance matrix – \sin^2 is merely the simplest.

As described in Rasmussen & Williams [1], valid covariance functions can be constructed by adding or multiplying simpler covariance functions. Thus, we can obtain a quasi-periodic kernel simply by multiplying a periodic kernel with one of the basic stationary kernels described earlier. The latter then specifies the rate of evolution of the periodic signal. For example, we can multiply equation (3.20) with an SE kernel,

$$k_{\text{QP,SE}}(x_i, x_j) = h^2 \exp \left(-\frac{\sin^2[\pi(x_i - x_j)/T]}{2w^2} - \frac{(x_i - x_j)^2}{\lambda^2} \right), \quad (3.21)$$

to model a quasi-periodic signal with a single evolutionary time scale λ .

Examples of functions drawn from these kernels are shown in figure 10. There are many more types of covariance functions in use, including some (such as the Matérn family above) that are better suited to model rougher, less smooth variations. However, the SE and RQ kernels already offer a great degree of freedom with relatively few hyperparameters, and covariance functions based on these are widely used to model time-series data.

Changepoints. We now describe how to construct appropriate covariance functions for functions that experience sudden changes in their characteristics. This section is meant to be expository; the covariance functions we describe are intended as examples rather than an exhaustive list of possibilities. To ease exposition, we assume the (single) input variable of interest, x , represents time. If additional features are available, they may be readily incorporated into the derived covariances [1].

A drastic change in covariance: we start by considering a function of interest as well behaved, except for a drastic change at the point x_c , which separates the function into two regions with associated covariance functions $k_1(\cdot, \cdot; \theta_1)$ before x_c and $k_2(\cdot, \cdot; \theta_2)$ after, where θ_1 and θ_2 represent

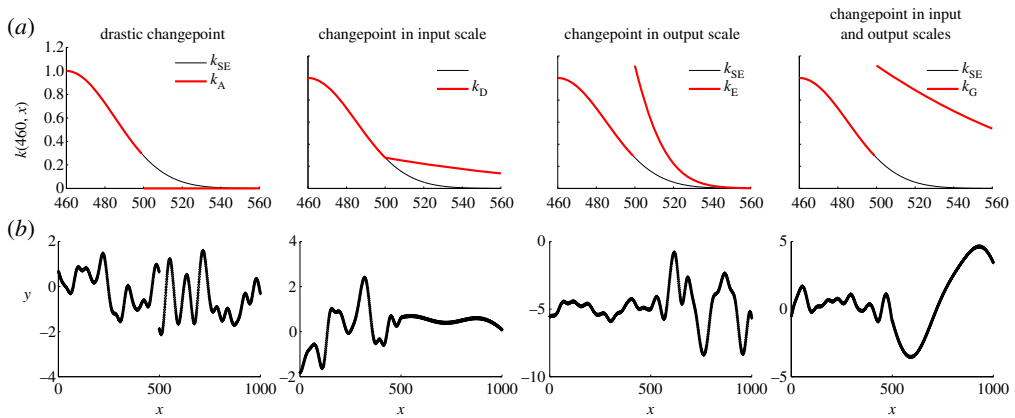


Figure 11. Example covariance functions (a) for the modelling of data with changepoints and associated draws (b) from the resultant GPs, indicating what kind of data that they might be appropriate for. Each changepoint covariance function is drawn as a bold line, with the standard SE kernel shown as k_{SE} for comparison (thin line). For ease of comparison, we fix the location of the changepoint hyperparameter to $x_c = 500$ and plot the functions over the interval from $460 \leq x \leq 560$. (Online version in colour.)

the values of any hyperparameters associated with k_1 and k_2 , respectively. The change is so drastic that the observations before x_c are completely uninformative about the observations after the changepoint. The full set of hyperparameters for this covariance function are hence the hyperparameters of the two covariance functions as well as the location of the changepoint, x_c . This covariance function is easily seen to be semi-positive definite and hence admissible [9,10]. The covariance function, and an example draw from the GP associated with it, are presented in the left-most plots of figure 11.

A drastic change in covariance with constraints: suppose a *continuous function* of interest is best modelled by different covariance functions, before and after a changepoint x_c . The function values after the changepoint are conditionally independent of the function values before, given the value at the changepoint itself. This represents an extension to the drastic covariance described earlier; our two regions can be drastically different, but we can still enforce continuity and smoothness constraints across the boundary between them. We call this covariance function the *continuous conditionally independent* covariance function. This covariance function can be extended to multiple changepoints, boundaries in multi-dimensional spaces, and also to cases where function derivatives are continuous at the changepoint. For proofs and details of this covariance function, the reader is invited to see Osborne *et al.* [10] and Reece *et al.* [11].

A sudden change in input scale: suppose a function of interest is well behaved, except for a drastic change in the input scale λ at time x_c , which separates the function into two regions with different degrees of long-term dependence. We let λ_1 and λ_2 represent the input scale of the function before and after the changepoint at x_c , respectively. The hyperparameters of this covariance consist of the two input scales, λ_1, λ_2 along with a common output scale, h , and the changepoint location, x_c . The second panel in figure 11 shows an example covariance function of this form (figure 11a) and an example function (figure 11b).

A sudden change in output scale: we now consider a function of interest as well behaved, except for a drastic change in the output scale h at time x_c , which separates the function into two regions. As before, we let h_1 and h_2 represent the output scales before and after the changepoint at x_c . The full set of hyperparameters of this model consists of the two output scales, h_1, h_2 , a common input scale, λ , and the location of the changepoint, x_c . The third panel of figure 11 shows an example covariance and associated example function. We note that we may readily combine changes in input and output scale into a single changepoint covariance (an example of which is shown in the right-most plots of figure 11).

A change in observation likelihood: hitherto, we have taken the observation likelihood as being defined by a single GP. We now consider other possible observation models, motivated by fault detection and removal [11,12]. For example, a sensor fault implies that the relationship between the underlying process model and the observed values is temporarily corrupted. In situations where a model of the fault is known, the faulty observations need not be discarded; they may still contain valuable information about the plant process. The interested reader is directed to Reece *et al.* [11,12], in which covariances for biased readings, stuck sensors and sensor drifts are discussed.

(ii) Mean functions

As the mean function will dominate our forecasts in regions far from the data, the choice of the prior mean function can have a profound impact on our predictions and must be chosen with this in mind. In the majority of cases in the literature, we find vague (i.e. high uncertainty) flat mean functions used. This choice is reinforced by considering the prior mean function as the expectation function, prior to any observed data, of our domain beliefs. In the vast majority of situations, the symmetry of our ignorance (i.e. we are equally unsure that a trend is up or down) leads to flat, often zero-offset, mean functions. As a simple example, we may have domain knowledge that our functions have a linear drift term, but we do not know the magnitude or direction. Whatever prior we place over the gradient of the drift will be necessarily symmetric and leads to a zero mean with variance defined by the vagueness of our priors. If we do have such domain knowledge, then we are free to incorporate this into our GP models. For example, consider the case in which we know that the observed time series consists of a deterministic component and an unknown additive component. Draws from our GP are hence

$$\mathbf{y}(\mathbf{x}) \sim \mathcal{N}(\mathbf{m}(\mathbf{x}; \boldsymbol{\theta}_m), \mathbf{K}(\mathbf{x}, \mathbf{x}; \boldsymbol{\theta}_C)), \quad (3.22)$$

in which the mean function, \mathbf{m} , has hyperparameters $\boldsymbol{\theta}_m$ that encode domain knowledge regarding the deterministic component and the covariance matrix \mathbf{K} has hyperparameters $\boldsymbol{\theta}_C$. For example, we may know that our observations are obtained from an underlying exponential decay with an unknown additive function along with coloured noise. Our mean function will hence be of the form $m(x_*) = A \exp(-ax_*)$ where A, a are unknown hyperparameters. Figure 12a shows a standard SE covariance GP used to model a small set of noisy data samples (dots) drawn from a function with an underlying exponential decay. The GP models the observed data well, but long-term predictions are naturally dominated by a flat prior mean function. In figure 12b, a GP with identical covariance is used, but the mean function is that of an exponential decay with unknown hyperparameters. Even a few data points are sufficient for the exponential function hyperparameters to be inferred, leading to long-term forecasts that are dominated by a (albeit uncertain) decay function.

4. Marginalizing hyperparameters

GP models have a number of hyperparameters (owing to both the covariance and mean functions) that we must *marginalize*² in order to perform inference. That is, we must first assign a prior $p(\theta)$ to these hyperparameters, informed by our domain knowledge. For example, in assigning a prior to the period of a tidal signal (as in §5a), we would use a prior that expressed that the most important period was of the order of days, rather than nanoseconds or millenia. In the absence of hard domain knowledge, these priors are chosen to be diffuse: for example, a Gaussian with high variance. Then, the quantity we are interested in is

$$p(y_* | \mathbf{y}) = \frac{\int p(y_* | \mathbf{y}, \theta) p(\mathbf{y} | \theta) p(\theta) d\theta}{\int p(\mathbf{y} | \theta) p(\theta) d\theta}, \quad (4.1)$$

²The process of marginalization refers to ‘integrating out’ uncertainty. For example, given $p(y, \theta) = p(y | \theta) p(\theta)$, we may obtain $p(y)$ by marginalizing over the unknown parameter θ , such that $p(y) = \int p(y | \theta) p(\theta) d\theta$.

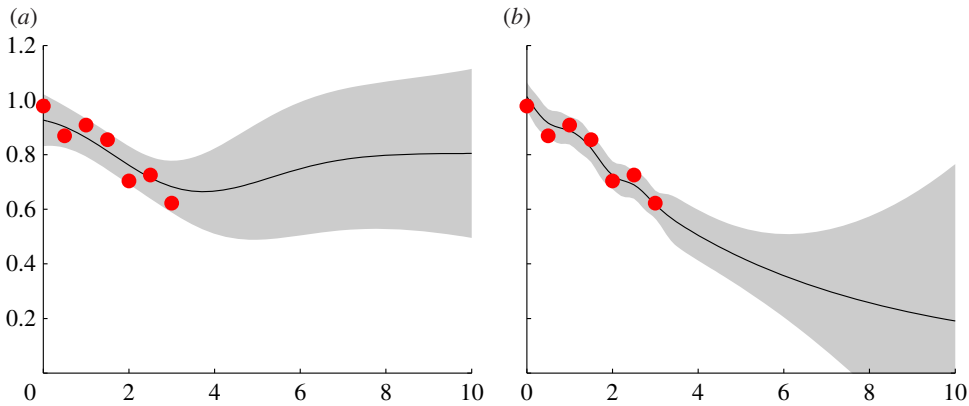


Figure 12. The effect of including a simple mean function. (a) A GP model with a flat prior mean and SE covariance function. The noisy observations are indicated by dots. The posterior from the GP is shown, along with $\pm 2\sigma$. (b) The same covariance function is used, but now the mean function has extra hyperparameters corresponding to an exponential decay with unknown time constant and scale. We see that the long-term forecasts in this example encode our prior belief in the decay function. (Online version in colour.)

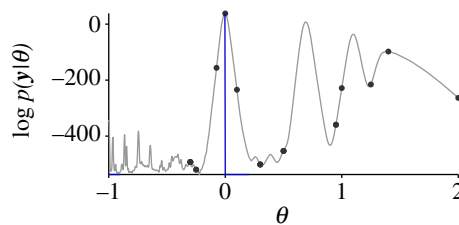


Figure 13. Samples (dots) obtained by optimizing the log-likelihood (grey curve) using a global optimizer, and the maximum-likelihood approximation (vertical line) of the likelihood surface. (Online version in colour.)

which requires two integrals to be evaluated. These are both typically non-analytic, owing to the complex form of the likelihood $p(\mathbf{y}|\theta)$ when considered as a function of hyperparameters θ . As such, we are forced to resort to approximate techniques.

Approximating an integral requires two problems to be solved. First, we need to make observations of the integrand, to explore it, and then those observations need to be used to construct an estimate for the integral. There are a number of approaches to both problems.

Optimizing an integrand (figure 13) is one fairly effective means of exploring it: we will take samples around the maxima of the integrand, which are likely to describe the majority of the mass comprising the integral. A local optimizer, such as a gradient ascent algorithm, will sample the integrand around the peak local to the start point, giving us information pertinent to at least that part of the integrand. If we use a global optimizer, our attempts to find the global extremum will ultimately result in all the integrand being explored, as desired.

Maximizing an integrand is most common when performing *maximum likelihood*. The integrands in (4.1) are proportional to the likelihood $p(\mathbf{y}|\theta)$: if the prior $p(\theta)$ is relatively flat, the likelihood will explain most of the variation of the integrands as a function of θ . Maximizing the likelihood hence gives a reasonable means of integrand exploration, as above. Maximum likelihood, however, specifies a generally unreasonable means of integral estimation: the likelihood is approximated as a Dirac delta function located at the θ that maximized the likelihood. As per figure 13, this completely ignores the width of the integrands, leading to potentially problematic features [13]. This approximation finds use when the likelihood is very peaked, as is the case when we have a great deal of data.

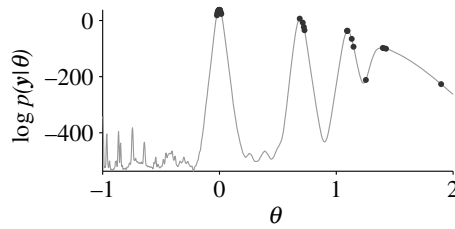


Figure 14. Samples obtained by taking draws from the posterior using a Markov chain Monte Carlo method.

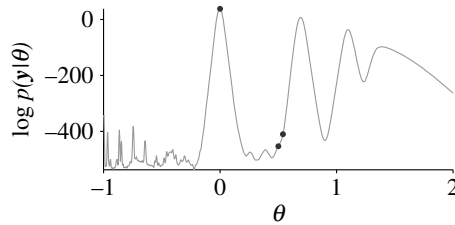


Figure 15. A set of samples that would lead to unsatisfactory behaviour from simple Monte Carlo.

A slightly more sophisticated approach to integral estimation is to take a *Laplace approximation*, which fits a Gaussian around the maximum-likelihood peak. This gives at least some representation of the width of the integrands. Yet further sophistication is displayed by the methods of *variational Bayes* [14], which treat the fitting of probability distributions to the problematic terms in our integrands as an optimization problem.

Monte Carlo techniques represent a very popular means of exploring an integrand. *Simple Monte Carlo* techniques draw random samples from the prior $p(\phi)$, to which our integrands are proportional. Note that (4.1) can be rewritten as

$$p(y_*|\mathbf{y}) = \int p(y_*|\mathbf{y}, \theta) p(\theta|\mathbf{y}) \, d\theta. \quad (4.2)$$

More sophisticated *Markov chain Monte Carlo* techniques [15] attempt to generate samples from the hyperparameter posterior

$$p(\theta|\mathbf{y}) = \frac{p(\mathbf{y}|\theta) p(\theta)}{\int p(\mathbf{y}|\theta) p(\theta) \, d\theta}, \quad (4.3)$$

to which (4.2) is proportional (figure 14 illustrates samples drawn using such a method). Sampling in this way ensures that we have many samples where the prior/posterior is large, and hence, where our integrands are likely to be large. This is a particular concern for multi-dimensional integrals, where the problem is complicated by the ‘curse of dimensionality’ [16]. Essentially, the volume of space that could potentially be explored is exponential in its dimension. However, a probability distribution, which must always have a total probability mass of one, will be highly concentrated in this space, ensuring our samples are likewise concentrated is a great boon. Moreover, Monte Carlo sampling ensures a non-zero probability of obtaining samples from any region where the prior is non-zero. This means that we can achieve some measure of broader exploration of our integrands.

Monte Carlo, does not, however, provide a very satisfactory means of integral estimation: it simply approximates the integral as the average over the obtained samples. As discussed by O’Hagan [17], this ignores the information content contained in the locations of the samples, leading to unsatisfactory behaviour. For example, imagine that we had three samples, two of which were identical: $\theta_1 = \theta_2$. In this case, the identical value will receive two-thirds of the weight, whereas the equally useful other value will receive only one-third. This is illustrated in figure 15.

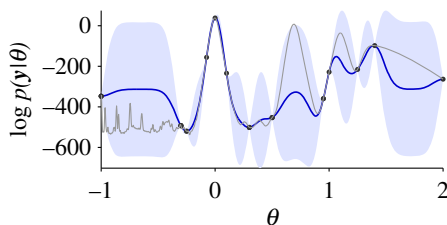


Figure 16. Bayesian quadrature fits a GP to the integrand, and thereby performs inference about the integral. (Online version in colour.)

In an attempt to address these issues, Bayesian quadrature [18,19] provides a model-based means of integral estimation. This approach assumes GPs over the integrands, using the obtained samples to determine a distribution for the integrals (figure 16). This probabilistic approach means that we can use the obtained variance in the integral as a measure of our confidence in its estimate. Of course, we still need to determine the hyperparameters for the GPs over the integrands. This problem is solved by adopting simple covariance functions for these GPs and using maximum likelihood to fit their hyperparameters (the maximum-likelihood output scale even has a closed-form solution). This renders the approach computationally tractable, to complement its superior accuracy.

In the examples to follow, we will exclusively use Bayesian quadrature to marginalize the hyperparameters of our GP models. Where desired, similar techniques are also used to calculate the posteriors for such hyperparameters

$$p(\theta|\mathbf{y}) = \frac{p(\mathbf{y}|\theta)p(\theta)}{\int p(\mathbf{y}|\theta)p(\theta)d\theta}. \quad (4.4)$$

Where the posterior for a hyperparameter is highly concentrated around a particular value, we will informally describe the hyperparameter as having been *learned* as having that value.

5. Examples

In the following examples, we briefly illustrate the GP approach to practical time-series analysis, highlighting the use of a variety of covariance and mean functions.

(a) Multi-dimensional weather sensor data

The first example we provide is based on real-time data that are collected by a set of weather, sea state and environment sensors on the south coast of the UK (see Roberts *et al.* [7] for more details). The network (Bramblemet) consists of four sensors (named Bramblemet, Sotonmet, Cambermet and Chimet), each of which measures a range of environmental variables (including wind speed and direction, air temperature, sea temperature and tide height) and makes up-to-date sensor measurements. We have two data streams for each variable at our disposal. The first is the real-time, but sporadic, measurements of the environmental variables; it is these that are presented as a multi-dimensional time series to the GP. Second we have access, retrospectively, to finer-grained data. We use this latter dataset for assessment only.

Figure 17 illustrates the efficacy of our GP prediction for a tide height dataset. In order to manage the four outputs of our tide function (one for each sensor), we rewrite so that we have a single output and inputs t , time, and l , a sensor label, as discussed in §3a and in §3c.

Note that our covariance over time is the sum of a periodic term and a *disturbance* term. Both are of the Matérn form with $\nu = \frac{5}{2}$. This form is a consequence of our expectation that the tides would be well modelled by the superposition of a simple periodic signal and an

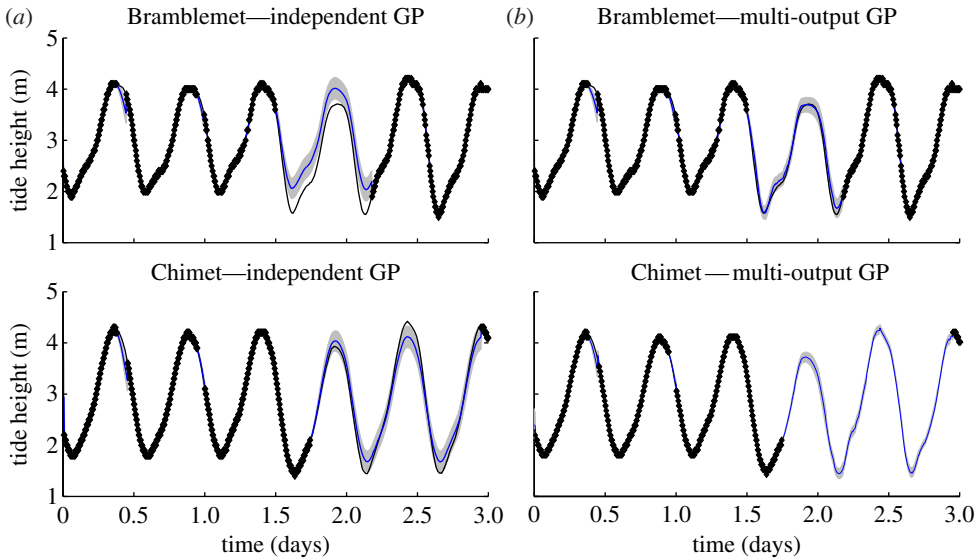


Figure 17. Prediction and regression of tide height data for (a) independent and (b) multi-output GPs. (Online version in colour.)

occasional disturbance signal due to exceptional conditions. Of course, for a better fit over the course of, say, a year, it would be possible to additionally incorporate longer term drifts and periods.

The period T of the periodic covariance term was unsurprisingly learnt as being about half a day, whereas for the disturbance term, the time scale w was found to be about two and a half hours. Note that this latter result is concordant with our expectations for the time scales of the weather events we intend our disturbance term to model.

Our algorithm learned that all four sensors were very strongly correlated, with spherical decomposition of the inferred correlation elements all very close to one. The hyperparameter matrix Σ of equation (3.18) (which defines relationships between variables) additionally gives an appropriate length scale for each sensor. From this inference, we determined weather events to have induced changes in tide height of the order of 20 per cent.

We also make allowances for the prospect of relative latency among the sensors by incorporating delay variables, introduced by a vector of delays in time observations [7]. We found that the tide signals at the Cambermet and Chimet stations were delayed by about 10 minutes relative to the other two. This makes physical sense—the Bramblemet and Sotonmet stations are located to the west of the Cambermet and Chimet stations, and the timing of high tide increases from west to east within the English channel.

Note the performance of our multi-output GP formalism when the Bramblemet sensor drops out at $t = 1.45$ days. In this case, the independent GP quite reasonably predicts that the tide will repeat the same periodic signal it has observed in the past. However, the GP can achieve better results if it is allowed to benefit from the knowledge of the other sensors' readings during this interval of missing data. Thus, in the case of the multi-output GP, by $t = 1.45$ days, the GP has successfully determined that the sensors are all very strongly correlated. Hence, when it sees an unexpected low tide in the Chimet sensor data (caused in this case by the strong northerly wind), these correlations lead it to infer a similarly low tide in the Bramblemet reading. Hence, the multi-output GP produces significantly more accurate predictions during the missing data interval, with associated smaller error bars. Exactly the same effect is seen in the later predictions of the Chimet tide height, where the multi-output GP predictions use observations from the other sensors to better predict the high tide height at $t = 2.45$ days.

Table 1. Predictive performances for 5 day Bramblemet tide height dataset. We note the superior performance of the GP compared with a more standard Kalman filter model. Error metrics shown are root mean square error (r.m.s.e) and normalized mean square error (n.m.s.e.), which is presented on a logarithmic, decibel scale.

algorithm	r.m.s.e (m)	n.m.s.e. (dB)
naive	7.5×10^{-1}	− 2.1
Kalman filter	1.7×10^{-1}	−15.2
independent GPs	8.7×10^{-2}	−20.3
multi-output GP	3.8×10^{-2}	−27.6

Note also that there are two brief intervals of missing data for all sensors just after both of the first two peak tides. During the second interval, the GP’s predictions for the tide are notably better than for the first—the greater quantity of data it has observed allows it to produce more accurate predictions. With time, the GP is able to build successively better models for the series.

The predictive performances for our various algorithms over this dataset can be found in table 1. For the Kalman filter comparison, a history length of 16 observations was used to generate each prediction because this gave rise to the best predictive ability for the Kalman model on out-of-sample data. However, note that our multi-output GP, which exploits correlations between the sensors, and the periodicity in each individual sensors’ measurements, significantly outperforms both the Kalman filter and the independent GP [7]. The naive result is obtained by repeating the last observed sensor value as a forecast and is included as a baseline only.

(b) Active data selection

We now demonstrate our active data selection algorithm. Using the fine-grained data (downloaded directly from the Bramblemet weather sensors), we can simulate how our GP would have chosen its observations had it been in control. Results from the active selection of observations from all four tide sensors are displayed in figure 18. Again, these plots depict dynamic choices; at time t , the GP must decide when next to observe, and from which sensor, given knowledge only of the observations recorded prior to t , in an attempt to maintain the uncertainty in tide height below 10 cm. The covariance function used was that described in the previous example, namely a sum of two $\nu = \frac{5}{2}$ Matérn covariance functions, one stationary and the other of periodic form. Consider first the case shown in figure 18*a*, in which separate independent GPs are used to represent each sensor. Note that a large number of observations are taken initially as the dynamics of the sensor readings are learnt, followed by a low but constant rate of observation. By contrast, for the multi-output case shown in figure 18*b*, the GP is allowed to explicitly represent correlations and delays between the sensors. As mentioned earlier, this dataset is notable for the slight delay of the tide heights at the Chimet and Cambermet sensors relative to the Sotonmet and Bramblemet sensors, due to the nature of tidal flows in the area. Note that after an initial learning phase as the dynamics, correlations and delays are inferred, the GP chooses to sample predominantly from the undelayed Sotonmet and Bramblemet sensors. The dynamics of the tide height at the Sotonmet sensor are more complex than the other sensors owing to the existence of a ‘young flood stand’ and a ‘double high tide’ in Southampton. For this reason, the GP selects Sotonmet as the most informative sensor and samples it most often. Despite no observations of the Chimet sensor being made within the time span plotted, the resulting predictions remain remarkably accurate. Consequently, only 119 observations are required to keep the uncertainty below the specified tolerance, whereas 358 observations were required in the independent case. This represents another clear demonstration of how our prediction is able to benefit from the readings of multiple sensors.

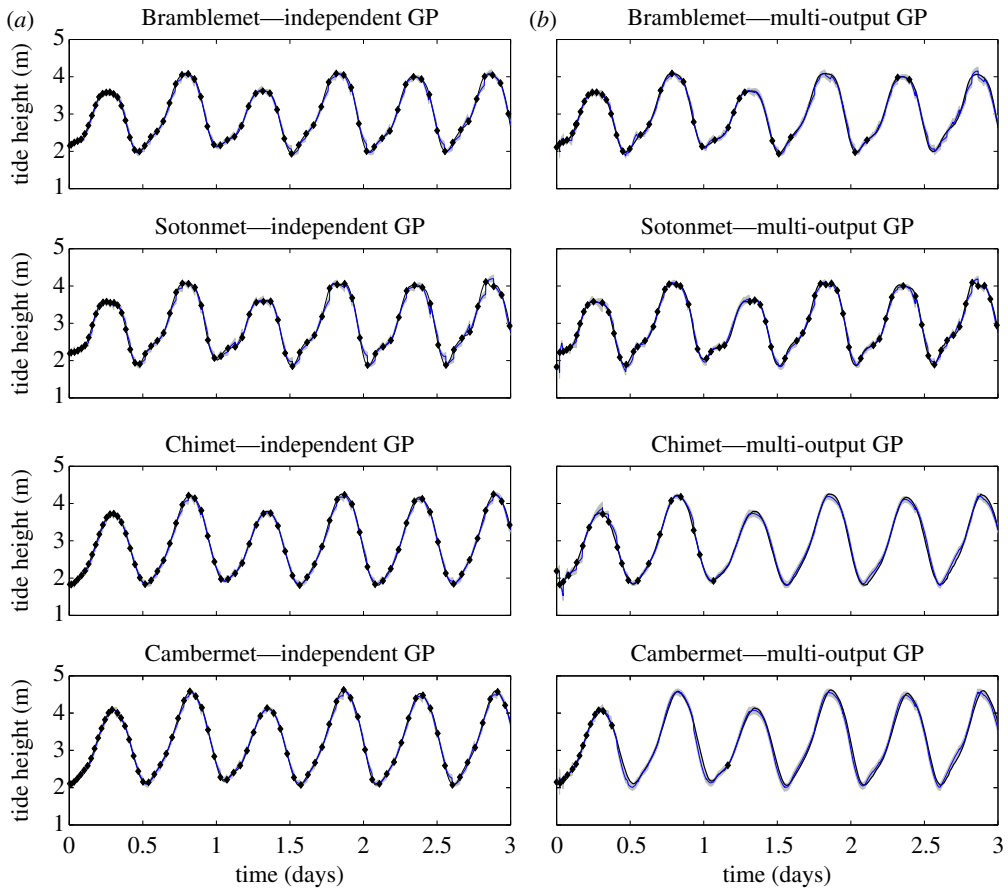


Figure 18. (a) Comparison of active sampling of tide data using independent and (b) multi-output GPs. Note that, in the case of multi-output GPs, one sensor reading (Sotonmet) slightly leads the other readings and is hence sampled much more frequently. In some cases, such as the Cambermet readings, only occasional samples are taken, yet the GP forecasts are excellent. (Online version in colour.)

(c) Changepoint detection

In Garnett *et al.* [9] and Osborne *et al.* [10], a fully Bayesian framework was introduced for performing sequential time-series prediction in the presence of changepoints. The position of a particular changepoint becomes a hyperparameter of the model that is marginalized using Bayesian quadrature. If the locations of changepoints in the data are of interest, the full posterior distribution of these hyperparameters can be obtained given the data. The result is a robust time-series prediction algorithm that makes well-informed predictions, even in the presence of sudden changes in the data. If desired, the algorithm additionally performs changepoint and fault detection as a natural by-product of the prediction process. In this section, we briefly present some exemplar datasets and the associated changepoint inference.

(i) Nile dataset

We first consider a canonical changepoint dataset, the minimum water levels of the Nile river during the period AD 622–1284 [20]. Several authors have found evidence supporting a change in input scale for this data around the year AD 722 [21]. The conjectured reason for this changepoint is the construction in AD 715 of a new device (a ‘nilometer’) on the island of Roda, which affected the nature and accuracy of the measurements.

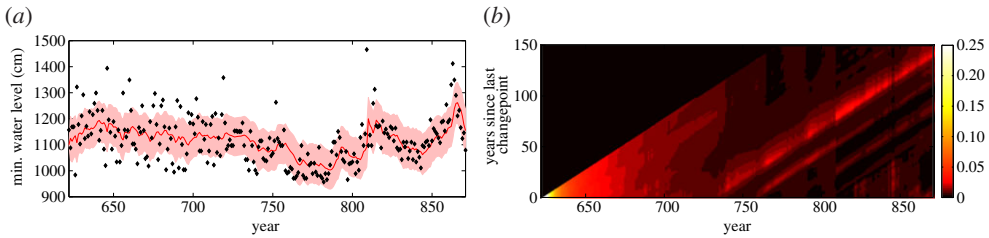


Figure 19. Prediction for the Nile dataset using input-scale changepoint covariance (a) and the corresponding posterior distribution for time since the changepoint (b). (Online version in colour.)

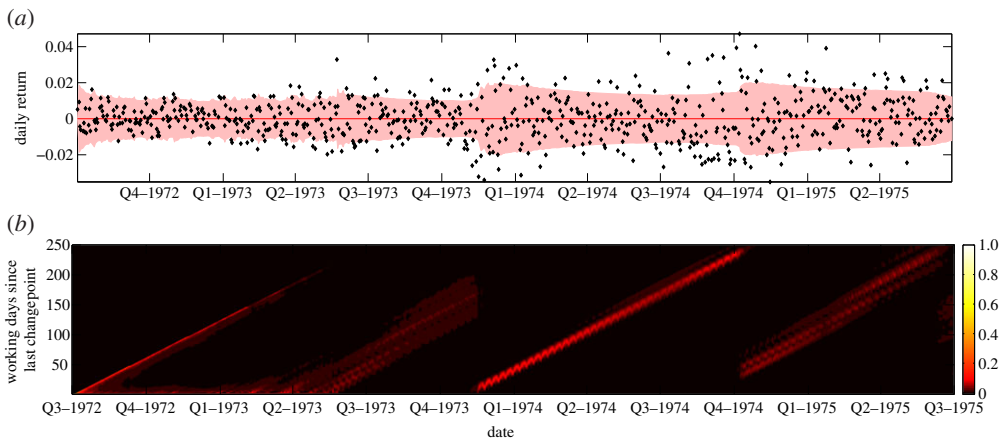


Figure 20. Online (sequential) one-step predictions (a) and posterior for the location of changepoint for the Dow–Jones industrial average data using an output-scale changepoint covariance (b). (Online version in colour.)

We performed one-step (next datum) lookahead prediction on this dataset using the input-scale changepoint covariance discussed earlier. The results can be seen in figure 19. Figure 19a shows our one-step predictions on the dataset, including the mean and $\pm\sigma$ error bars. Figure 19b shows the posterior distribution of the number of years since the last changepoint. A changepoint around AD 720–722 is clearly visible and agrees with previous results [21].

(ii) 1972–1975 Dow–Jones industrial average

As a second canonical changepoint dataset, we present the series of daily returns of the Dow–Jones industrial average between 3 July 1972 and 30 June 1975 [22]. This period included a number of newsworthy events that had significant macroeconomic influences, as reflected in the Dow–Jones returns.

We performed sequential one-step (next datum) prediction on this data using a GP with a diagonal covariance that assumed all measurements were independent and identically distributed (as under the efficient market hypothesis, returns should be uncorrelated). However, the variance of these observations was assumed to undergo changes, and as such, we used a covariance that incorporated such changes in output scale. We had three hyperparameters to marginalize: the variance before the changepoint, the variance after the changepoint and, finally, the location of that changepoint.

Our results are plotted in figure 20. Our model clearly identifies the important changepoints that likely correspond to the commencement of the Organization of the Petroleum Exporting

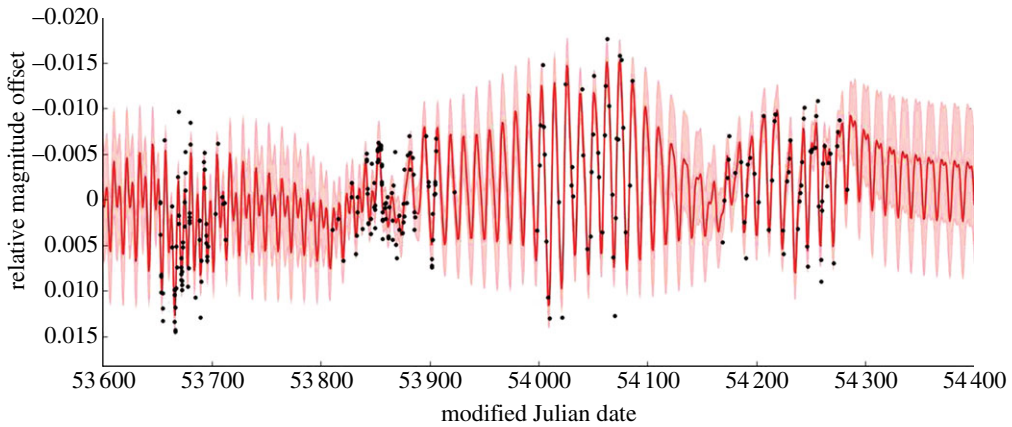


Figure 21. Predictive distribution for a quasi-periodic GP model using a mixed SE and RQ kernel, trained and conditioned on observations made with the 0.8 m Automated Patrol Telescope [24] using the Strömgren b and y filters. The dots represent the observations, the line is the mean of the predictive posterior distribution and the shaded region encompasses the $\pm\sigma$ interval. (Online version in colour.)

Countries embargo on 19 October 1973, and the resignation of Richard Nixon as President of the USA on 9 August 1974. A weaker changepoint is identified early in 1973, which Adams & MacKay [22] speculate is due to the beginning of the Watergate scandal.

(d) Quasi-periodic modelling of stellar light curves

Many Sun-like stars display quasi-periodic brightness variations on time scales of days to weeks, with amplitudes ranging from a few parts per million to a few per cent. These variations are caused by the evolution and rotational modulation of magnetically active regions, which are typically fainter than the surrounding photosphere. In this case, we may expect a range of both periodic covariance scales w and evolutionary time scales λ , corresponding to different active region sizes and lifetimes, respectively. This can be achieved by replacing one or both of the SE kernels in equation (3.13) by RQ kernels (equation (3.14)). Finally, we can also allow for short-term irregular variability or correlated observational noise by including a separate, additive SE or RQ kernel. For example, Pont *et al.* [23] used a GP with such quasi-periodic kernels to model the total irradiance variations of the Sun in order to predict its radial velocity variations.

In figure 21, we show the results of a quasi-periodic GP regression to photometric observations of the well-known planet host star HD 189733, taken from Henry & Winn [24]. The kernel used consists of a periodic SE component (equation (3.21)) multiplied by an RQ term (equation (3.14)) to allow for a range of evolutionary time scales, plus an additive white noise term (equation (3.12)). Inference over the hyperparameters of interest yielded expected values of $h = 6.68$ mmag, $T = 11.86$ days, $w = 0.91$, $\alpha = 0.23$, $\lambda = 17.81$ days and $\sigma = 2.1$ mmag, where σ is the amplitude of the white noise term. Our period is in excellent agreement with Henry & Winn [24]. The relatively long periodic length scale w indicates that the variations are dominated by a small number of fairly large active regions. The evolutionary term has a relatively short time scale, λ , but a shallow index, α , which is consistent with the notion that the active regions on this star evolve relatively fast and/or that, as in the Sun, active regions located at different latitudes have different rotation rates (known as differential rotation).

(e) Modelling light curves of transiting exoplanets

One of the most successful ways of discovering and characterizing extra-solar planets (i.e. planets not in our solar system) is through observing transit light curves. A transit occurs when a

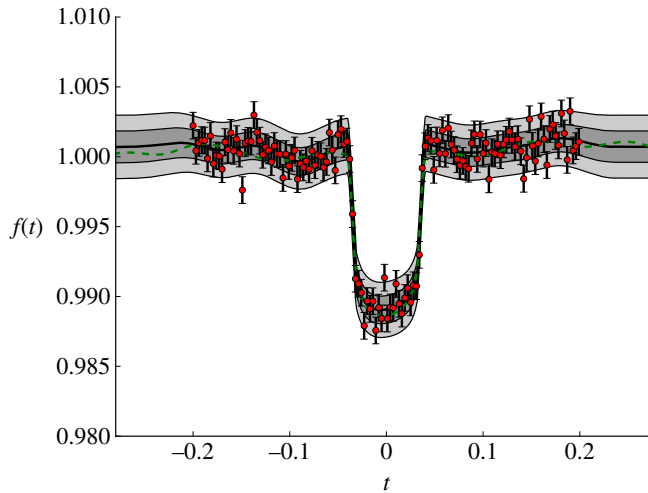


Figure 22. As an example of a complex mean function, we here model data from an exoplanet transit light curve. The data are fitted with a GP with an exoplanet transit mean function and a squared exponential covariance kernel to model the correlated noise process and the effects of external state variables. The shaded regions are at $\pm 1, 2\sigma$ from the posterior mean. (Online version in colour.)

planet periodically passes between its host star and the Earth, blocking a portion of the stellar light, and produces a characteristic dip in the light curve. From this transit, we can measure such physical parameters as the planet-to-star radius ratio and the inclination of the orbit. While transit light curves are readily described by a deterministic parametric function, real observations are corrupted by systematic noise in the detector, external state variables (such as the temperature of the detector, orbital phase, position of the host star on the charge-coupled device array, etc.), as well as the underlying flux variability of the host star. As it is not possible to produce a deterministic model to account for all these systematics, a GP may be used to place a distribution over possible artefact functions, modelling correlated noise as well as subtle changes in observed light curves due to external state variables. We hence encode the transit curve as the mean function of a GP. The covariance function has inputs given by time and external state variables (hence, this is a multi-input, single-output model). By integrating out our uncertainty (see §4) in the hyperparameters of the GP (which model all the systematic artefacts and noise processes), we can gain much more realistic inference of probability distribution of the transit function parameters (the hyperparameters of the mean function). For a detailed discussion of the application of GPs to transit light curves, see Gibson *et al.* [25], in which the instrumental systematics are represented by a GP with an SE kernel (equation (3.13)) and input parameters representing the external state variables. Robust inference of transit parameters is required to perform detailed studies of transiting systems, including the search for atomic and molecular signatures in the atmospheres of exoplanets. Figure 22 shows this GP model fitting to the time series of observations. More details are found in Gibson *et al.* [25].

6. Conclusion

In this paper, we have presented a brief outline of the conceptual and mathematical basis of GP modelling of time series. As ever, a practical implementation of the ideas concerned requires jumping algorithmic rather than theoretical hurdles, which we do not discuss here because of space constraints. Some introductory code may be found at ftp://ftp.robots.ox.ac.uk/pub/outgoing/mebden/misc/GP_tut.zip and more general code can be downloaded from <http://www.gaussianprocess.org/gpml>. Space has not permitted discussion

of exciting recent trends in GP modelling that allow for more explicit incorporation of differential equations governing the system dynamics (either observed or not), such as *latent force models* [26]. Further extensions, using GPs as building blocks in more complex probabilistic models, are of course possible, and recent research has also highlighted the use of GPs for numerical integration, global optimization, mixture-of-experts models, unsupervised learning models and much more.

The authors thank Alex Rogers, Roman Garnett, Richard Mann, Tom Evans, Mark Smith and Chris Hart. Part of this work was funded by the UK research councils, whose support is gratefully acknowledged.

References

1. Rasmussen CE, Williams CKI. 2006 *Gaussian processes for machine learning*. Cambridge, MA: MIT Press.
2. Osborne MA, Rogers A, Ramchurn S, Roberts SJ, Jennings NR. 2008 Towards real-time information processing of sensor network data using computationally efficient multi-output Gaussian processes. In *Int. Conf. on Information Processing in Sensor Networks (IPSN 2008)*, April 2008, pp. 109–120. IEEE Computer Society.
3. Abrahamsen P. 1997 A review of Gaussian random fields and correlation functions, 2nd edn. Technical Report no. 917, Norwegian Computing Center, Box 114, Blindern, 0314 Oslo, Norway.
4. Stein ML. 2005 Space-time covariance functions. *J. Am. Stat. Assoc.* **100**, 310–322. (doi:10.1198/016214504000000854)
5. Sasena MJ. 2002 Flexibility and efficiency enhancements for constrained global design optimization with Kriging approximations. PhD thesis, University of Michigan, Ann Arbor, Michigan.
6. Pinheiro J, Bates D. 1996 Unconstrained parameterizations for variance–covariance matrices. *Stat. Comput.* **6**, 289–296. (doi:10.1007/BF00140873)
7. Osborne M., Roberts S, Rogers A, Jennings N. In press. Real-time information processing of environmental sensor network data. *Trans. Sensor Networks* **9**.
8. MacKay DJC. 1998 Introduction to Gaussian processes. In *Neural networks and machine learning* (ed. CM Bishop), pp. 84–92. Berlin, Germany: Springer.
9. Garnett R, Osborne MA, Roberts SJ. 2009 Sequential Bayesian prediction in the presence of changepoints. In *Proc. 26th Annual Int. Conf. Machine Learning, Montreal, June 2009*, pp. 345–352. New York, NY: ACM.
10. Osborne M, Reece S, Rogers A, Roberts S, Garnett R. 2010 Sequential Bayesian prediction in the presence of changepoints and faults. *Comp. J.* **53**, 1430–1446.
11. Reece S, Garnett R, Osborne MA, Roberts SJ. 2009 Anomaly detection and removal using non-stationary Gaussian processes. Technical report PARG-09-01, University of Oxford, Oxford, UK.
12. Reece S, Claxton C, Nicholson D, Roberts SJ. 2009 Multi-sensor fault recovery in the presence of known and unknown fault types. In *Proc. 12th Int. Conf. on Information Fusion (FUSION 2009)*, Seattle, WA, pp. 1696–1703. IEEE Computer Society.
13. MacKay DJC. 2002 *Information theory, inference and learning algorithms*. Cambridge, UK: Cambridge University Press.
14. Bishop C. 2006 *Pattern recognition and machine learning*. Berlin, Germany: Springer.
15. Chen J, Gupta AK. 2000 *Parametric statistical change point analysis*. Berlin, Germany: Birkhäuser.
16. Duda RO, Hart PE, Stork DG. 2000 *Pattern classification*, 2nd edn. London, UK: Wiley-Interscience.
17. O'Hagan A. 1987 Monte Carlo is fundamentally unsound. *Statistician* **36**, 247–249. (doi:10.2307/2348519)
18. O'Hagan A. 1991 Bayes–Hermite quadrature. *J. Stat. Plan. Inference* **29**, 245–260. (doi:10.1016/0378-3758(91)90002-V)
19. Rasmussen CE, Ghahramani Z. 2003 Bayesian Monte Carlo. In *Advances in neural information processing systems*, vol. 15, pp. 489–496. Cambridge, MA: MIT Press.
20. Whitcher B, Byers SD, Guttorp P, Percival DB. 2002 Testing for homogeneity of variance in time series: long memory, wavelets and the Nile River. *Water Resour. Res.* **38**, 10–1029. (doi:10.1029/2001WR000509)

21. Ray BK, Tsay RS. 2002 Bayesian methods for change-point detection in long-range dependent processes. *J. Time Series Anal.* **23**, 687–705. (doi:10.1111/1467-9892.00286)
22. Adams RP, MacKay DJC. 2007 Bayesian online changepoint detection. (<http://arxiv.org/abs/0710.3742>)
23. Pont F, Aigrain S, Zucker S. 2011 A simple method to estimate radial velocity variations due to stellar activity using photometry. *Mon. Not. R. Astron. Soc.* **419**, 3147.
24. Henry GW, Winn JN. 2008 The rotation period of the planet-hosting star HD 189733. *Astron. J.* **135**, 68. (doi:10.1088/0004-6256/135/1/68)
25. Gibson NP, Aigrain S, Roberts S, Evans T, Osborne M, Pont F. In press. A Gaussian process framework for modelling instrumental systematics: application to transmission spectroscopy. *Mon. Not. R. Astron. Soc.*
26. Luengo D, Lawrence N, Álvarez M. 2009 Latent force models. In *Proc. 12th Int. Workshop on Artificial Intelligence and Statistics, Clearwater Beach, FL, 16–18 April 2009*, pp. 9–16. JMLR Workshop and Conference Proceedings, vol. 5. MIT Press.