

Real-Time Information Processing of Environmental Sensor Network Data using Bayesian Gaussian Processes

MICHAEL A. OSBORNE, University of Oxford
 STEPHEN J. ROBERTS, University of Oxford
 ALEX ROGERS, University of Southampton
 NICHOLAS R. JENNINGS, University of Southampton

In this paper, we consider the problem faced by a sensor network operator who must infer, in real-time, the value of some environmental parameter that is being monitored at discrete points in space and time by a sensor network. We describe a powerful and generic approach built upon an efficient multi-output Gaussian process that facilitates this information acquisition and processing. Our algorithm allows effective inference even with minimal domain knowledge, and we further introduce a formulation of Bayesian Monte Carlo to permit the principled management of the hyperparameters introduced by our flexible models. We demonstrate how our methods can be applied in cases where the data is delayed, intermittently missing, censored and/or correlated. We validate our approach using data collected from three networks of weather sensors and show that it yields better inference performance than both conventional independent Gaussian processes and the Kalman filter. Finally, we show that our formalism efficiently re-uses previous computations by following an online update procedure as new data sequentially arrives, and that this results in a four-fold increase in computational speed in the largest cases considered.

Categories and Subject Descriptors: C.2.4 [Computer Communication Networks]: Distributed Systems – Distributed Applications

General Terms: Algorithms, Design, Performance

Additional Key Words and Phrases: Learning of models from data, Gaussian processes, information processing, adaptive sampling

ACM Reference Format:

Osborne, M. A., Roberts, S. J., Rogers, A. and Jennings, N. R. 2012. Real-time information processing using Bayesian Gaussian processes. *ACM Trans. Sensor. Netw.* 9, 1, Article 1 (January 2012), 32 pages.
 DOI = 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

1. INTRODUCTION

Sensor networks have recently generated a great deal of research interest within the computer and physical sciences, and their use for the scientific monitoring of remote and hostile environments is increasingly common-place. While early sensor networks were a simple evolution of existing automated data loggers that collected data for later off-line scientific analysis, more recent sensor networks now make their data available

A preliminary version of the work presented in this paper appeared as: Osborne, M. A., Rogers, A., Ramchurn, S., Roberts, S. J. and Jennings, N. R. (2008) Towards Real-Time Information Processing of Sensor Network Data using Computationally Efficient Multi-output Gaussian Processes. In: *International Conference on Information Processing in Sensor Networks (IPSN 2008)*, St. Louis, Missouri, USA. pp. 109-120.

Author's addresses: M. A. Osborne and S. J. Roberts, Department of Engineering Science, University of Oxford, Oxford, UK; A. Rogers and N. R. Jennings, Electronics and Computer Science, University of Southampton, Southampton, UK.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2012 ACM 1539-9087/2012/01-ART1 \$10.00

DOI 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

in real-time through the internet, and increasingly perform some form of real-time data processing or aggregation to provide useful summary information to users [Hart and Martinez 2006].

Providing real-time access to sensor data in this way presents many novel challenges; not least the need for self-describing data formats, and standard protocols such that the existence and capabilities of sensors can be advertised to potential users. However, more significantly for us, many of the information processing tasks that would previously have been performed off-line by the expert owner or single user of an environmental sensor network (such as detecting faulty sensors, performing ‘data cleaning’ to remove erroneous readings, fusing noisy measurements from several sensors, and deciding how frequently data should be collected), must now be performed autonomously in real-time. Furthermore, since such information processing is likely to be performed within centralised sensor repositories — of which Weather Underground (<http://www.wunderground.com>), pachube (<http://www.pachube.com/>) and Microsoft Research’s SensorMap (<http://atom.research.microsoft.com/sensormap/>) are embryonic examples — and will be applied to open sensor networks where additional sensors may be deployed at any time, and existing sensors may be removed, repositioned or updated after deployment (such as the rooftop weather sensors within the Weather Underground), these information processing tasks may have to be performed with only limited knowledge of the precise location, reliability, and accuracy of each sensor.

Now, many of the information processing tasks described above have previously been tackled by applying principled Bayesian methodologies from the academic literature of geospatial statistics and machine learning: specifically, kriging [Cressie 1991] and Gaussian processes [Rasmussen and Williams 2006]. However, due to the computational complexity of these approaches, to date they have largely been used off-line in order to analyse and re-design existing sensor networks (e.g. to reduce maintenance costs by removing the least informative sensors from an existing sensor network [Fuentes et al. 2007], or to find the optimum placement of a small number of sensors, after a trial deployment of a larger number has collected data indicating their spatial correlation [Krause et al. 2006]). Alternatively, they have been applied to single sensors and have ignored the correlations that exist between different sensors within the network [Kho et al. 2009; Padhy et al. 2010]. Thus, there is a clear need for more computationally efficient algorithms, in order that this information processing can be performed at scale in real-time.

Against this background, this paper describes our work developing just such an algorithm. More specifically, we present a novel iterative formulation of a multi-output Gaussian process (GP) that uses a computationally efficient implementation of *Bayesian Monte Carlo* to marginalise hyperparameters, efficiently re-uses previous computations by following an online update procedure as new data sequentially arrives, and uses a principled ‘windowing’ of data in order to maintain a reasonably sized data set. We use this GP to build a probabilistic model of the environmental variables being measured by sensors within the network, that is tolerant to data that may be missing, delayed, censored and/or correlated. This model allows us to then perform information processing tasks including: modelling the accuracy of the sensor readings, predicting the value of missing sensor readings, predicting how the monitored environmental variables will evolve in the near future, and performing active sampling by deciding when and from which sensor to acquire readings. We validate our approach using data collected from three networks of weather sensors and show that it yields better inference performance than both conventional independent single-output Gaussian processes, that model each sensor independently, and the Kalman filter. Finally, we analyse the computational speed of our algorithm. We show that our formalism efficiently re-uses previous computations by following an online update procedure as

new data sequentially arrives, and this results in a four-fold increase in computational speed in the largest cases considered. Furthermore, it enables all of the information processing tasks described in this paper to be performed in real-time.

The remainder of this paper is organised as follows: Section 2 describes the information processing problem that we face. Section 3 presents our Gaussian process formulation, and Section 4 describes the three sensor networks used to validate it. In Section 5 we present experimental results using data from these networks. In Section 6 we present results on the computational cost of our algorithm, and we discuss its trial deployment in Section 7. Related work is discussed in Section 8, and finally, we conclude in Section 9.

2. THE INFORMATION PROCESSING PROBLEM

As discussed above, we require that our algorithm be able to autonomously perform data acquisition and information processing despite having only limited specific knowledge of each of the sensors in its local neighbourhood (e.g. their precise location, reliability, and accuracy). To this end, we require that it explicitly represents:

- (1) The uncertainty in the estimated values of environmental variables being measured, noting that sensor readings will always incorporate some degree of measurement noise.
- (2) The correlations or delays that exist between sensor readings; sensors that are close to one another, or in similar environments, will tend to make similar readings, while many physical processes involving moving fields (such as the movement of weather fronts) will induce delays between sensors.

We then require that it uses this representation in order to:

- (1) Perform regression and prediction of environmental variables; that is, interpolate between sensor readings to predict variables at missing sensors (i.e. sensors that have failed or are unavailable through network outages), and perform short term prediction in order to support decision making.
- (2) Perform efficient active sampling by selecting when to take a reading, and which sensor to read from, such that the minimum number of sensor readings are used to maintain the estimated uncertainty in environmental variables below a specified threshold (or similarly, to minimise uncertainty given a constrained number of sensor readings). Such constraints may reflect the computational limitations of the processor on which the algorithm is running, or alternatively, where the algorithm is actually controlling the sensors within the network, it may reflect the constrained power consumption of the sensors themselves.

More specifically, the problem that we face can be cast as a multivariate regression and decision problem in which we have $l = 1, \dots, L$ environmental variables $x_l \in \mathbb{R}$ of interest (such as air temperature, wind speed or direction specified at different sensor locations). We assume a set of N potentially noisy sensor readings, $\{[l_1, t_1], z_1], \dots, [l_N, t_N], z_N]\}$, in which we, for example, observe the value z_1 for the l_1^{th} variable at time t_1 , whose true unknown value is y_1 . Where convenient, we may group the inputs as $x = [l, t]$. Note that we do not require that all the variables are observed at the same time, nor do we impose any discretisation of our observations into regularly spaced time steps. We define our vector of observations as $z_d \triangleq [z_1, \dots, z_N]$ of variables labelled by $l_d \triangleq [l_1, \dots, l_N]$ at times $t_d \triangleq [t_1, \dots, t_N]$. Given this data, we are interested in inferring the vector of values y_* for any other vector of variables labelled by l_* at times t_* . We might, for example, wish to predict the values of variables at

locations we do not currently have readings for (perhaps due to sensor failure), or to forecast future values for a variable at a particular sensor.

In the sections to follow, we'll describe an algorithm built around Gaussian processes to address this problem. We exploit whatever prior expectations we may have for the variation of the environmental variables over time and sensor, by building a model that can, for example, express that a particular variable may have some (unknown) periodicity over time and that readings from the various sensors may be correlated with one another. Such strong expectations, however, are not necessary; our algorithm is flexible enough to learn a broad range of models. This is achieved by using a means of approximate integration to marginalise over the class of models that are compatible with our observations. In essence, we sample a number of reasonable models and weight the predictions made under each potential model by its likelihood given the data. Those predictions are made in a computationally efficient, iterative fashion as we receive ever more data, and given our principled Bayesian framework, provide a measure of the uncertainty in our estimates. By managing these uncertainties, our algorithm is also able to select for itself the most informative observations, recognising that observations may be costly and must be selected intelligently.

3. GAUSSIAN PROCESSES

Multivariate regression problems of the form described above are increasingly being addressed using Gaussian processes (GPs). These represent a powerful way to perform Bayesian inference about functions, and here, we consider our environmental variables to be just such a function [Rasmussen and Williams 2006]. This function takes as inputs the variable label and time pair x and produces as output the variable's value y . In this work, we will assume that our inputs are always known (e.g. our data is time-stamped), and will incorporate them into our background knowledge, or context, I . A GP is then a generalised multivariate Gaussian prior distribution over the (potentially infinite number of) outputs of this function:

$$\begin{aligned} p(\mathbf{y} \mid \boldsymbol{\mu}, \mathbf{K}, I) &\triangleq \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \mathbf{K}) \\ &\triangleq \frac{1}{\sqrt{\det 2\pi\mathbf{K}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \mathbf{K}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right). \end{aligned}$$

It is specified by prior mean and covariance functions, which generate $\boldsymbol{\mu}$ and \mathbf{K} . The multivariate Gaussian distribution is qualified for this role due to the fact that both its marginal probabilities and conditional probabilities are themselves Gaussian. This allows us to produce analytic posterior distributions for outputs of interest, conditioned on whatever sensor readings have been observed. Furthermore, this posterior distribution will have both a predictive mean and a variance to explicitly represent our uncertainty.

While the fundamental theory of GPs is well established (see Rasmussen and Williams [2006] for example), there is much scope for the development of computationally efficient implementations especially for applications such as sensor networks which involve large numbers of sensors collecting large quantities of timestamped data. To this end, in this work we present a novel on-line formalism of a multi-dimensional GP that allows us to model the correlations between sensor readings, and to update this model on-line as new observations are sequentially available. In the next sections we describe the covariance functions that we use to represent correlations and delays between sensor readings, the *Bayesian Monte Carlo* method that we use to marginalise the hyperparameters of these covariance functions, and how we efficiently update the model as new data is received, by reusing the results of previous computations, and applying a principled 'windowing' of our data series.

3.1. Mean and Covariance Functions

The prior mean of a GP represents whatever we expect for our function before seeing any data. Throughout this paper, we will take this as a function constant in time for each sensor, such that $\mu([l, t]) = \mu_l$. 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 [Abrahamsen 1997; Stein 2005], 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.

Examples of covariance functions are given by the Matérn class [Rasmussen and Williams 2006], given by

$$K_{\text{Mtn}}(t, t'; h, w, \nu) \triangleq h^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu} r \right)^\nu \mathfrak{K}_\nu \left(\sqrt{2\nu} r \right) \quad (1)$$

where $r = \left| \frac{t-t'}{w} \right|$, $\nu > 0$ is a smoothness hyperparameter (larger ν implies smoother functions) and \mathfrak{K}_ν is the modified Bessel function. Fortunately, (1) simplifies for half-integer ν , to give, for the example of $\nu = 5/2$

$$K_{\text{Mtn}}(t, t'; h, w, \nu = 5/2) \triangleq h^2 \left(1 + \sqrt{5}r + \frac{5r^2}{3} \right) \exp \left(-\sqrt{5}r \right), \quad (2)$$

and its modification for periodic signals

$$K_{\text{per-Mtn}}(t, t'; h, w, \nu = 5/2) \triangleq h^2 \left(1 + \sqrt{5}r_{\text{per}} + \frac{5r_{\text{per}}^2}{3} \right) \exp \left(-\sqrt{5}r_{\text{per}} \right), \quad (3)$$

where $r_{\text{per}} = \sin \pi \left| \frac{t-t'}{w} \right|$. This modification of r to accommodate periodicity is readily made for other covariance functions.

There are many ways to construct a covariance over a multi-dimensional input space, such as our two-dimensional input space $x = [l, t]$. We will typically exploit the fact that a product of covariances is a covariance in its own right, and write

$$K([l, t], [l', t']) = K_t(t, t') K_l(l, l'),$$

taking independent covariance terms over time and sensor label. If the number of sensors is not too large, we can arbitrarily represent the covariance matrix K_l over sensor labels using the spherical decomposition [Pinheiro and Bates 1996]. This allows us to arbitrarily represent any possible covariance over sensor labels. For a covariance matrix W , we write

$$W = S \text{diag}(\tau). \quad (4)$$

Here S is an upper triangular matrix, whose n th column contains the spherical coordinates in \mathbb{R}^n of a point on the hypersphere \mathbb{S}^{n-1} , followed by the requisite number of zeros. As an example, S for a four dimensional space is

$$S = \begin{bmatrix} 1 & \cos \theta_1 & \cos \theta_2 & \cos \theta_4 \\ 0 & \sin \theta_1 & \sin \theta_2 \cos \theta_3 & \sin \theta_4 \cos \theta_5 \\ 0 & 0 & \sin \theta_2 \sin \theta_3 & \sin \theta_4 \sin \theta_5 \cos \theta_6 \\ 0 & 0 & 0 & \sin \theta_4 \sin \theta_5 \sin \theta_6 \end{bmatrix}. \quad (5)$$

This form ensures that $S^T S$ has ones across its diagonal and hence all other entries may be thought of as akin to correlation coefficients, lying between -1 and 1 . Mean-

while, τ is a vector of the scales for each dimension. If W is the covariance matrix parameterising the Mahalanobis distance, τ represents a vector of input scales for each input, and the off-diagonal elements of $S^T S$ express the correlation amongst inputs. If instead we have used W as a parameterisation of $K_1([1, \dots, L], [1, \dots, L])$, τ gives a vector of output scales corresponding to the different values $l = [1, \dots, L]$. $S^T S$ gives the correlation coefficients for each pair of values l might take.

These intuitive connections provide the motivation for the use of this parameterisation. Finally, note that the total number of hyperparameters required by the spherical parameterisation is $\frac{1}{2}n(n+1)$ if the side length of W is n . Note that this is the same number as would be required to parameterise the non-zero elements of an upper triangular Cholesky factor $R = \text{chol } W$ directly. The spherical parameterisation hence requires a large number of hyperparameters, but allows us to express any possible covariance matrix.

Finally, to represent measurement noise, we further extend the covariance function to

$$V([l, t], [l', t']) \triangleq K([l, t], [l', t']) + \sigma^2 \delta([l, t] - [l', t']),$$

where $\delta(-)$ is the Kronecker delta and σ^2 represents the variance of additive Gaussian noise.

The flexibility of our model comes at the cost of the introduction of a number of hyperparameters, which we collectively denote as ϕ . These include correlation hyperparameters, along with others such as the periods and amplitudes of each covariance term (i.e. w and h) and the noise deviation σ . The constant prior means μ_1, \dots, μ_L are also included as additional hyperparameters. Taking these hyperparameters as given and using the properties of the Gaussian distribution, we are able to write our predictive equations as

$$p(\mathbf{x}_* | \mathbf{y}_d, \phi, I) = \mathbf{N}(\mathbf{x}_*; \mathbf{m}_*, \mathbf{C}_*), \quad (6)$$

where, collecting our inputs as $\mathbf{x}_* \triangleq [\mathbf{l}_*, \mathbf{t}_*]$ and $\mathbf{x}_d \triangleq [\mathbf{l}_d, \mathbf{t}_d]$, we have:

$$\mathbf{m}_* = \boldsymbol{\mu}(\mathbf{x}_*; \phi) + \mathbf{K}(\mathbf{x}_*, \mathbf{x}_d; \phi) \mathbf{V}(\mathbf{x}_d, \mathbf{x}_d; \phi)^{-1} (\mathbf{y}_d - \boldsymbol{\mu}(\mathbf{x}_d; \phi)) \quad (7)$$

$$\mathbf{C}_* = \mathbf{K}(\mathbf{x}_*, \mathbf{x}_*; \phi) - \mathbf{K}(\mathbf{x}_*, \mathbf{x}_d; \phi) \mathbf{V}(\mathbf{x}_d, \mathbf{x}_d; \phi)^{-1} \mathbf{K}(\mathbf{x}_d, \mathbf{x}_*; \phi). \quad (8)$$

For notational convenience, we will often drop the dependence upon ϕ from our covariance and mean functions.

3.2. Marginalisation

Of course, it is rare that we can be certain a priori about the values of our hyperparameters. One common approach to dealing with this uncertainty is to explore the possible space of hyperparameters (using a local search method such as gradient ascent) in order to select the set that exhibit the greatest log likelihood [Rasmussen and Williams 2006]. However, this approach fails to represent the uncertainty in the value of these hyperparameters, and thus, it fails to represent the true uncertainty in the final predictions of the environmental variables being measured.

To correctly represent this uncertainty we must marginalise over ϕ , and thus, rather than equation (6), we must consider

$$p(\mathbf{y}_* | \mathbf{z}_d, I) = \frac{\int d\phi p(\mathbf{y}_* | \mathbf{z}_d, \phi, I) p(\mathbf{z}_d | \phi, I) p(\phi | I)}{\int d\phi p(\mathbf{z}_d | \phi, I) p(\phi | I)}. \quad (9)$$

Unfortunately, both our likelihood $p(\mathbf{z}_d | \phi, I)$ and predictions $p(\mathbf{y}_* | \mathbf{z}_d, \phi, I)$ exhibit non-trivial dependence upon ϕ and so our integrals are non-analytic. As such, we resort

to quadrature, which inevitably involves evaluating the two quantities

$$\begin{aligned} q(\phi) &\triangleq p(\mathbf{y}_* | \mathbf{z}_d, \phi, I) \\ r(\phi) &\triangleq p(\mathbf{z}_d | \phi, I) \end{aligned} \quad (10)$$

at a set of sample points $\phi_s = \{\phi_1, \dots, \phi_\eta\}$, giving $\mathbf{q}_s \triangleq q(\phi_s)$ and $\mathbf{r}_s \triangleq r(\phi_s)$. Of course, this evaluation is a computationally expensive operation. Defining the vector of all possible function inputs as ϕ , we clearly can't afford to evaluate $\mathbf{q} \triangleq q(\phi)$ or $\mathbf{r} \triangleq r(\phi)$. Note that the more complex our model, and hence the greater the number of hyperparameters, the higher the dimension of the hyperparameter space we must sample in. As such, the complexity of models we can practically consider is limited by the curse of dimensionality. We can view our sparse sampling as introducing a form of uncertainty about the functions q and r , which we can again address using Bayesian probability theory.

To this end, we apply *Bayesian Monte Carlo*, and thus, assign further GP priors to q and r [Rasmussen and Ghahramani 2003]. This choice is motivated by the fact that variables over which we have a multivariate Gaussian distribution are joint Gaussian with any projections of those variables. As such, given that integration is a projection, we can use our computed samples \mathbf{q}_s in order to perform Gaussian process regression about the value of integrals over $q(\phi)$, and similarly for r . Note that the quantity we wish to perform inference about,

$$\psi \triangleq p(\mathbf{y}_* | \mathbf{q}, \mathbf{r}, \mathbf{z}_d, I) = \frac{\int d\phi_* q(\phi_*) r(\phi_*) p(\phi_* | I)}{\int d\phi_* r(\phi_*) p(\phi_* | I)}, \quad (11)$$

possesses richer structure than that previously considered using Bayesian Monte Carlo techniques. In our case, $r(\phi)$ appears in both our numerator and denominator integrals, introducing correlations between the values we estimate for them. The correlation structure of this system is illustrated in figure 1.

In considering any problem of inference, we need to be clear about both what information we have and which uncertain variables we are interested in. In our case, both function values, \mathbf{q}_s and \mathbf{r}_s , and their locations, ϕ_s , represent valuable pieces of knowledge¹. As with our convention above, we will take knowledge of sample locations ϕ_s to be implicit within I . We respectively define $\mathbf{m}^{(q)}$ and $\mathbf{m}^{(r)}$ as the means for q and r conditioned on \mathbf{q}_s and \mathbf{r}_s from (7), $\mathbf{C}^{(q)}$ and $\mathbf{C}^{(r)}$ the similarly conditional covariances from (8). The ultimate quantity of our interest is then

$$\begin{aligned} p(\mathbf{y}_* | \mathbf{q}_s, \mathbf{r}_s, \mathbf{z}_d, I) &= \iiint p(\mathbf{y}_* | \mathbf{q}, \mathbf{r}, \mathbf{z}_d, I) p(\psi | \mathbf{q}, \mathbf{r}, I) p(\mathbf{q} | \mathbf{q}_s, I) p(\mathbf{r} | \mathbf{r}_s, I) d\psi d\mathbf{q} d\mathbf{r} \\ &= \iiint \psi \delta\left(\psi - \frac{\int d\phi_* q_* r_* p(\phi_* | I)}{\int d\phi_* r_* p(\phi_* | I)}\right) \mathbf{N}(\mathbf{q}; \mathbf{m}^{(q)}, \mathbf{C}^{(q)}) \mathbf{N}(\mathbf{r}; \mathbf{m}^{(r)}, \mathbf{C}^{(r)}) d\psi d\mathbf{q} d\mathbf{r} \\ &= \int \frac{\int d\phi_* m_*^{(q)} r_* p(\phi_* | I)}{\int d\phi_* r_* p(\phi_* | I)} \mathbf{N}(\mathbf{r}; \mathbf{m}^{(r)}, \mathbf{C}^{(r)}) d\mathbf{r}. \end{aligned} \quad (12)$$

Here, unfortunately, our integration over \mathbf{r} becomes nonanalytic. However, we can employ a maximum *a posteriori* approximation by assuming

$$\mathbf{N}(\mathbf{r}; \mathbf{m}^{(r)}, \mathbf{C}^{(r)}) \simeq \delta(\mathbf{r} - \mathbf{m}^{(r)}).$$

¹As discussed by O'Hagan [1987], traditional, frequentist Monte Carlo effectively ignores the information content of ϕ_s , leading to several unsatisfactory features.

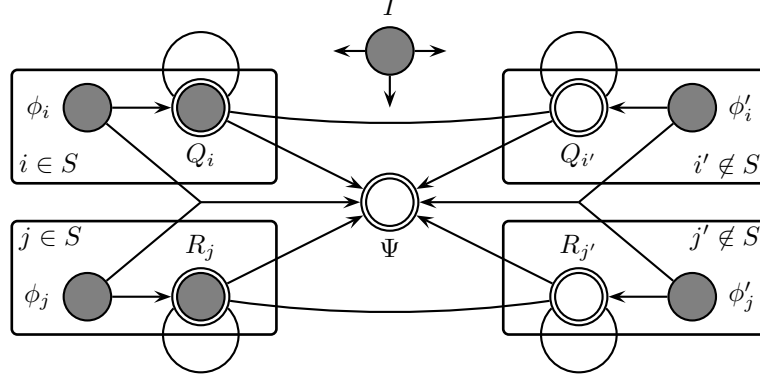


Fig. 1: Bayesian network for marginalising hyperparameters using Bayesian Monte Carlo. Shaded nodes are known and double-circled nodes are deterministic given all their parents. All Q nodes are correlated with one another, as are all R nodes, and the context I is correlated with all nodes.

Essentially, we assume that our uncertainty in the likelihood function is small and trust that our more correct treatment of the uncertainty in the prediction function q is sufficient. Given that we are likely to be uncertain about $r(\phi)$ at exactly the same ϕ we are uncertain about $q(\phi)$, those points far removed from our samples ϕ_s , this approximation is not unreasonable. Before we can state our consequent results, to each of our hyperparameters we assign a Gaussian prior distribution (or if our hyperparameter is restricted to the positive reals, we instead assign a Gaussian distribution to its log) given by

$$p(\phi | I) \triangleq \mathbf{N}(\phi; \boldsymbol{\nu}, \lambda^T \lambda) . \quad (13)$$

Note that the intuitive spherical parameterisation (5) assists with the elicitation of priors over its hyperparameters. We then assign a *squared exponential* covariance function for the GP over both q and r given by

$$K(\phi, \phi') \triangleq \mathbf{N}(\phi; \phi', \mathbf{w}^T \mathbf{w}) . \quad (14)$$

Finally, using the further definition for $i, j \in \mathcal{I}_s$, that

$$\mathfrak{N}_s(i, j) \triangleq \mathbf{N} \left(\begin{bmatrix} \phi_i \\ \phi_j \end{bmatrix}; \begin{bmatrix} \boldsymbol{\nu} \\ \boldsymbol{\nu} \end{bmatrix}, \begin{bmatrix} \lambda^T \lambda + \mathbf{w}^T \mathbf{w} & \lambda^T \lambda \\ \lambda^T \lambda & \lambda^T \lambda + \mathbf{w}^T \mathbf{w} \end{bmatrix} \right) , \quad (15)$$

our approximation gives us

$$p(\mathbf{y}_* | \mathbf{y}_d, I) \simeq \mathbf{q}_s^T \boldsymbol{\rho} , \quad (16)$$

where the weights of this linear combination are

$$\boldsymbol{\rho} \triangleq \frac{\mathbf{K}(\phi_s, \phi_s)^{-1} \mathfrak{N}_s \mathbf{K}(\phi_s, \phi_s)^{-1} \mathbf{r}_s}{\mathbf{1}_{s,1}^T \mathbf{K}(\phi_s, \phi_s)^{-1} \mathfrak{N}_s \mathbf{K}(\phi_s, \phi_s)^{-1} \mathbf{r}_s} , \quad (17)$$

and $\mathbf{1}_{s,1}$ is a vector containing only ones of dimensions equal to \mathbf{q}_s . With a GP on $p(\mathbf{y}_* | \phi, I)$, each $q_i = p(\mathbf{y}_* | \mathbf{z}_d, \phi_i, I)$ will be a slightly different Gaussian. Hence we effectively approximate $p(\mathbf{y}_* | \mathbf{z}_d, I)$ as a Gaussian (process) mixture; Bayesian Monte Carlo returns a weighted sum of our predictions evaluated at a sample set of hyperparameters. The assignment of these weights is informed by the best use of all pertinent

information. As such, it avoids the risk of overfitting that occurs when applying a less principled technique such as maximum likelihood [MacKay 2002].

3.3. Censored Observations

In the work above, we have assumed that observations of our variables of interest were corrupted by simple Gaussian noise. However, in many contexts, we instead observe *censored* observations. That is, we might observe that a variable was above or below certain thresholds, but no more. Examples are rich within the weather sensor networks considered. Float sensors are prone to becoming lodged on sensor posts, reporting only that the water level is below that at which it is stuck. Conversely, wind sensors may reach their maximum reading in very strong winds, reporting that the wind is greater than this limit, but providing no more information. Furthermore, readings are often rounded to the nearest integer, meaning that if we observe a reading of x , we can say only that the true value was between $x - 0.5$ and $x + 0.5$. We can readily extend our inference framework to allow for such a noise model.

More precisely, we assume that we actually observe bounds b_c that constrain Gaussian-noise corrupted versions z_c of the underlying variables of interest y_c at x_c . This framework allows for imprecise censored observations. Note that the noise variance for censored observations may differ from the noise variance associated with other observations. Conditioned on a combination of censored and un-censored observations, the distribution for our variables of interest is

$$p(y_* | z_d, b_c, I) = \frac{\int d\phi \int_{b_c} dz_c p(y_* | z_d, z_c, \phi, I) p(z_c | z_d, \phi, I) p(z_d | \phi, I) p(\phi | I)}{\int d\phi \int_{b_c} dz_c p(z_c | z_d, \phi, I) p(z_d | \phi, I) p(\phi | I)}.$$

While we cannot determine this full, non-Gaussian distribution easily, we can determine the analytic forms for its mean and covariance. We use the abbreviations $m_{c|d,\phi} \triangleq m(z_c | z_d, \phi, I)$ and $C_{c|d,\phi} \triangleq C(z_c | z_d, \phi, I)$. To reflect the influence of our censored observations, our likelihoods become

$$p(z_d, b_c | \phi, I) = N(z_d; m(z_d | \phi, I), C(z_d | \phi, I)) \int_{b_c} dz_c N(z_c; m_{c|d,\phi}, C_{c|d,\phi}),$$

giving the new weights over hyperparameter samples $\rho^{(cd)}$. We can then write our predictive mean as

$$m(y_* | z_d, b_c, I) = \sum_{i \in s} \rho_i^{(cd)} \left(\mu(x_*, \phi_i) + K(x_*, [x_c, x_d]; \phi_i) \right. \\ \left. V([x_c, x_d], [x_c, x_d]; \phi_i)^{-1} \left[\frac{\int_{b_c} dz_c z_c N(z_c; m_{c|d,\phi_i}, C_{c|d,\phi_i})}{\int_{b_c} dz_c N(z_c; m_{c|d,\phi_i}, C_{c|d,\phi_i})} - \mu(y_c; \phi_i) \right] \right),$$

noting that a censored observation is intuitively treated as an uncensored observation equal to the conditional mean of the GP over the bounded region. We have also the predictive covariance

$$C(y_{st} | z_d, b_c, I) = \sum_{i \in s} \rho_i^{(cd)} \left(K(x_*, x_*; \phi_i) - \right. \\ \left. K(x_*, [x_c, x_d]; \phi_i) V([x_c, x_d], [x_c, x_d]; \phi_i)^{-1} K([x_c, x_d], y_*; \phi_i) \right).$$

We now have the problem of determining the integrals

$$\int_{b_c} dz_c N(z_c; \mathbf{m}_{c|d,\phi}, \mathbf{C}_{c|d,\phi}) \quad \text{and} \quad \int_{b_c} dz_c z_c N(z_c; \mathbf{m}_{c|d,\phi}, \mathbf{C}_{c|d,\phi}) ,$$

which are non-analytic. With only a single censored observation, we can use standard functions for such integrals. Assuming that our censored observations are independent of one another (conditioned on the other observations z_d) would also allow us the use of such results, but is usually an unreasonable approximation. Fortunately, there also exists an efficient Monte Carlo algorithm for the purpose of resolving such integrals [Genz 1992]. It is this algorithm that we will choose for the appropriate empirical tests to follow.

3.4. Efficient Implementation

Now the implementation of equations (7) and (8) involves the inverse of a matrix, and the most stable way to implement this is through the use of the Cholesky decomposition, $\mathbf{R}(x_d, x_d)$, of $\mathbf{V}(x_d, x_d)$, rather than inverting the matrix directly. Performing this Cholesky decomposition represents the most computationally expensive operation we must perform; its cost scaling as $O(N^3)$ in the number of data points N . However, as discussed earlier, we do not intend to use our GP with a fixed set of data, but rather, within an on-line algorithm that receives new observations over time. As such, we must be able to iteratively update our predictions in as little time as possible. Fortunately, we can do so by exploiting the special structure of our problem. When we receive new data, our \mathbf{V} matrix is changed only in the addition of a few new rows and columns. Hence most of the work that went into computing its Cholesky decomposition at the last iteration can be recycled to produce the new Cholesky decomposition (see Appendix A.1 for details of this operation). Another problematic calculation required by (7) and (8) is the computation of the data-dependent term $\mathbf{R}(x_d, x_d)^{-1}(y_d - \mu(x_d))$, in which $y_d - \mu(x_d)$ is also only changing due to the addition of new rows. As such, efficient updating rules are also available for this term (see Appendix A.2). Finally, we are also able to perform an efficient update to the likelihood² of our hyperparameters, as per Appendix A.3. As per (17), the updated sample likelihoods r_s then give rise to updated weights ρ over our hyperparameter samples. With these efficient rules, we are able to reduce the overall cost of an update from $O(N^3)$ to $O(N^2)$.

However, we can further increase the efficiency of our updates by making a judicious assumption. In particular, experience shows that our GP requires only a very small number of recent observations in order to produce good estimates. Indeed, most covariance functions have very light tails such that only points within a few multiples of the time scale are at all relevant to the point of interest. Hence we seek sensible ways of discarding information once it has been rendered ‘stale’, to reduce both memory usage and computational requirements.

One pre-eminently reasonable measure of the value of data is the uncertainty we still possess after learning it. In particular, we are interested in how uncertain we are about x_* ; as given by the covariance of our Gaussian mixture equation (16). Our approach is thus to drop our oldest data points (those which our covariances typically deem least relevant to the current predictant) until this uncertainty exceeds some predetermined threshold.

Just as we were able to efficiently update our Cholesky factor upon the receipt of new data, so we can downdate to remove data (see Appendix A.4 for the details of this operation). Unfortunately, there is no such efficient downdate for the data-dependent term

²Note that likelihoods are commonly smaller than machine precision and so for numerical reasons we prefer to instead work with log-likelihoods, exponentiating as required.

$\mathbf{R}(x_d, x_d)^{-1}(\mathbf{y}_d - \mu(x_d))$. Naturally there is no computational benefit here to be gained in downdating our hyperparameter likelihoods, whose dimension remains constant regardless of the size of the dataset. Hence our likelihoods are never downdated, and we retain all knowledge of otherwise dropped data that is pertinent to our marginalisation of hyperparameters.

Our proposal allows us to rapidly remove unwanted data, compute our uncertainty about \mathbf{y}_* , and then repeat as required; the GP will retain only as much data as necessary to achieve a pre-specified degree of accuracy. This allows a principled way of ‘windowing’ our data series.

Finally, we turn to the implementation of our marginalisation procedure. Essentially, our approach is to maintain an ensemble of GPs, one for each hyperparameter sample, each of which we update and downdate according to the proposals above. Note that the computations for each hyperparameter sample can be treated independently, lending our algorithm to parallelisation.

The predictions of each sample are then weighted and combined according to equation (16). Note that the only computations whose computational cost grows at greater than a quadratic rate in the number of samples, η , are the Cholesky decomposition and multiplication of covariance matrices in equation (17), and these scale rather poorly as $O(\eta^3)$. To address this problem, we take our Gaussian priors for each different hyperparameter $\phi_{(e)} \in \phi$ as independent. We further take a covariance structure given by the product of terms over each hyperparameter, the common *product correlation rule* (e.g. Sasena [2002])

$$K(\phi, \phi') = \prod_e K_e(\phi_{(e)}, \phi'_{(e)}). \quad (18)$$

If we additionally consider a simple grid of samples, such that ϕ_s is the tensor product of a set of samples $\phi_{(e),s}$ over each hyperparameter, then the problematic term in equation (16) reduces to the Kronecker product of the equivalent term over each individual hyperparameter:

$$\begin{aligned} & \mathbf{K}(\phi_s, \phi_s)^{-1} \mathfrak{N}_s \mathbf{K}(\phi_s, \phi_s)^{-1} \\ &= \mathbf{K}(\phi_{(1),s}, \phi_{(1),s})^{-1} \mathfrak{N}_s(\phi_{(1),s}, \phi_{(1),s}) \mathbf{K}(\phi_{(1),s}, \phi_{(1),s})^{-1} \\ & \quad \otimes \mathbf{K}(\phi_{(2),s}, \phi_{(2),s})^{-1} \mathfrak{N}_s(\phi_{(2),s}, \phi_{(2),s}) \mathbf{K}(\phi_{(2),s}, \phi_{(2),s})^{-1} \\ & \quad \otimes \dots \end{aligned} \quad (19)$$

This means that we only have to perform the expensive Cholesky factorisation and multiplication with matrices whose size equals the number of samples for each hyperparameter, rather than on a matrix of size equal to the total number of hyperparameter samples. This hence represents an effective way to avoid the ‘curse of dimensionality’.

Note that we can also leverage some of the benefits of maximum likelihood to our own advantage. In particular, given an initial grid of samples, we can perform regular updates to their positions by performing single-stage gradient ascent. That is, the value of $\phi_{(e),i}$, the i th sample in the e th hyperparameter, is adjusted according to the average derivative with respect to $\phi_{(e),i}$ of the likelihood of all samples that share $\phi_{(e),i}$. Shifting a hyperparameter sample does, however, require the fresh Cholesky decomposition of the covariance matrix associated with that sample, along with the re-computation of the data-dependent term, likelihood and the weights term (19). As such, we perform such sample shifting relatively rarely to avoid too excessive a computational burden.

In general, limiting the number of hyperparameter samples is of critical importance to achieving practical computation. As such, we should exploit any and all prior in-

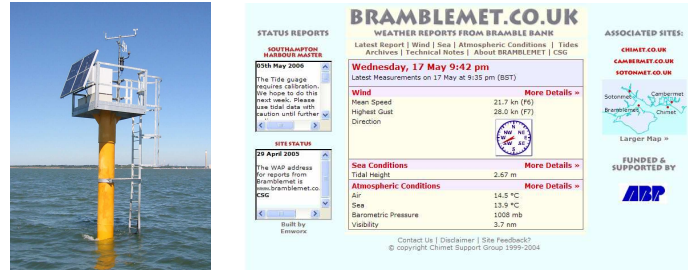


Fig. 2: The Bramble Bank weather station and web site.

formation that we possess about the system to limit the volume of hyperparameter space that our GP is required to explore online. For example, an informative prior expressing that a tidal period is likely to be around half a day will greatly reduce the number of samples required for this hyperparameter. Similarly, an offline analysis of any available training data will return sharply peaked posteriors over our hyperparameters that will further restrict the required volume to be searched over on-line. For example, in the results that we present in Section 5, we represent the tidal period hyperparameter with only a single sample on-line, so certain does training data make us of its value. Finally, a simpler and less flexible covariance model, with fewer hyperparameters, could be chosen if computational limitations become particularly severe.

Applied together, these features provide us with an efficient on-line algorithm that can be applied in real-time as data is sequentially collected from the sensor network. Pseudo-code for our approach is displayed in Algorithm 1.

3.5. Active Data Selection

Finally, in addition to the regression and prediction problem described in Section 2, we are able to use the same algorithm to perform active data selection. This is a decision problem concerning which observations should be taken. In this, we once again take a utility that is a function of the uncertainty in our predictions. We specify a utility of negative infinity if our uncertainty about any variable is greater than a pre-specified threshold, and a fixed negative utility is assigned as the cost of an observation (in general, this cost could be different for different sensors). Note that the uncertainty increases monotonically in the absence of new data, and shrinks in the presence of an observation. Hence our algorithm is simply induced to make a reading whenever the uncertainty grows beyond a pre-specified threshold.

Our algorithm can also decide which observation to make at this time, by determining which sensor will allow it the longest period of grace until it would be forced to observe again. This clearly minimises the number of costly observations. Note that the uncertainty of a GP, as given by (8), is actually dependent only on the location of the observation, not its actual value. Hence the uncertainty we imagine remaining after taking an observation from a sensor can be quickly determined without having to speculate about what data we might possibly collect. However, this is true only so long as we do not consider the impact of new observations on our hyperparameter sample weights (17). Our approach is to take the model, in particular the weights over samples, as fixed, and investigate only how different schedules of observations affect our predictions within it. With this proviso, we are guaranteed to maintain our uncertainty below a specified threshold, while taking as few observations as possible.

ALGORITHM 1: GP-TRACK($\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \delta, \text{MEAN}(\cdot; \cdot), \text{COV}(\cdot, \cdot; \cdot), \phi_s$)

```

// Returns the sequential lookahead predictions made with our GP algorithm.
//    $\{(x_t, z_t)\}_{t=1}^T$ : All data.
//    $\delta$ : Lookahead.
//    $\text{MEAN}(x, \phi)$ : Mean function ( $\mu$ ), taking as arguments an input  $x$  and
//                 hyperparameter vector  $\phi$ .
//    $\text{COV}(x, x'; \phi)$ : Covariance function ( $K$ ), taking as arguments two inputs
//                  $x$  and  $x'$  and hyperparameter vector  $\phi$ .
//    $\phi_s \triangleq \{\phi_j\}_{j=1}^\eta$ : Hyperparameter samples.

1:  $\mathbf{x}_d \leftarrow \{\}$ 
2:  $\mathbf{z}_d \leftarrow \{\}$ 
3: for  $i = 1$  to  $\eta$  do
4:    $GP[i] \leftarrow \text{INITIALISE-GP}(\text{MEAN}, \text{COV}, \phi_i)$  // Supply  $GP[i]$  with the mean function
    $\text{MEAN}(x, \phi_i)$  and covariance function  $\text{COV}(x, x'; \phi_i)$ .
5: end for
6:  $\text{weights-term} \leftarrow \text{CALCULATE-WEIGHTS-TERM}(\phi_s)$  // Calculate the weights term from (19),
    $\mathbf{K}(\phi_s, \phi_s)^{-1} \mathfrak{N}_s \mathbf{K}(\phi_s, \phi_s)^{-1}$ .
7: for  $t = 1$  to  $T - \delta$  do
8:    $(\mathbf{x}_d, \mathbf{z}_d, \text{drop-inds}, \text{add-ind}) \leftarrow \text{INCREMENT-DATA}((\mathbf{x}_d, \mathbf{z}_d), (x_t, z_t), GP)$ 
   // Add new datum (index  $\text{add-ind}$ ), and drop old data (indices  $\text{drop-inds}$ ) as necessary (as
   // per Section 3.4).
9:   for  $i = 1$  to  $\eta$  do
10:     $GP[i] \leftarrow \text{DOWNDATe-GP}(GP[i], \text{drop-inds})$  // DOWNDATeS a GP as per Section 3.4,
    revising covariance matrix and data-dependent term to allow for dropped data.
11:     $GP[i] \leftarrow \text{UPDATE-GP}(GP[i], (\mathbf{x}_d, \mathbf{z}_d), \text{add-ind})$  // UPDATeS a GP as per Section 3.4,
    revising covariance matrix, data-dependent term and likelihoods to allow for added
    data.
12:     $(m_{t+\delta, i}, C_{t+\delta, i}) \leftarrow \text{MAKE-PREDICTIONS}(x_{t+\delta}, GP[i])$  // Compute the posterior mean
    and variances for predictant  $y_\star = y_{t+\delta}$  using (7) and (8).
13:   end for
14:    $\rho \leftarrow \text{DETERMINE-WEIGHTS}(\text{weights-term}, GP[1 : \eta])$  // Using (17).
15:    $(m'_{t+\delta}, C'_{t+\delta}) \leftarrow \text{COMBINE-PREDICTIONS}(\{\rho_i, m_{t+\delta, i}, C_{t+\delta, i}\}_{i=1}^\eta)$ 
   // Gaussians are combined in a weighted mixture, using (16).
16:   Return:  $(m'_{t+\delta}, C'_{t+\delta})$ 
17:   if IS-TIME-TO-MOVE( $t$ ) // Infrequently true then
18:      $(\phi_s, \text{move-inds}) \leftarrow \text{MANAGE-HYPERSAMPLES}(GP[1 : \eta], \phi_s, t)$ 
     // Move hyperparameter samples  $\text{move-inds}$ , as per Section 3.4.
19:     for  $i \in \text{move-inds}$  do
20:        $GP[i] \leftarrow \text{OVERWRITE-GP}(GP[i], (\mathbf{x}_d, \mathbf{z}_d), \phi_i)$  // Recompute all quantities given the
       new hyperparameter sample.
21:     end for
22:      $\text{weights-term} \leftarrow \text{CALCULATE-WEIGHTS-TERM}(\phi_s)$ 
23:   end if
24: end for

```

4. EVALUATION DATASETS

In order to empirically evaluate the information processing algorithm described in the previous section, we have used data from three networks of weather sensors of different sizes (specifically, 4, 16 and 670 sensors) that provide a range of different sensor measurements. These networks are the Bramblemet sensor network, the Wannengrat Alpine Observatory, and the UK Met Office MIDAS land surface weather stations. The use of these weather sensor networks is attractive since they exhibit challenging correlations and delays whose physical processes are well understood. However, we stress

that the approach that we describe in this paper is general, and can equally well be applied to any time series that exhibits complex correlations and delays.

4.1. Bramblemet Weather Sensor Network

The Bramblemet network 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 available through separate web pages (see Figure 2 and <http://www.bramblemet.co.uk/>). To facilitate the autonomous collection of sensor data by our information processing algorithm, we have worked with the sensor network operators (Associated British Ports) and supplemented each sensor web page with machine readable RDF data (see Figure 3 for an example of this format – current sensor data in this format is available at <http://www.bramblemet.co.uk/bra.rdf>). This format is attractive as it represents a fundamental element of the semantic web, and there exist a number of software tools to parse, store and query it. More importantly, it allows the sensor data to be precisely defined through standard ontologies [Lassila and Swick 1999]. For example, linking the predicate *geo:lat* to the ontology available at http://www.w3.org/2003/01/geo/wgs84_pos# precisely defines the value “50.79472” as representing a latitude in the WGS84 geodetic reference datum. While ontologies for sensor data have yet to be standardised, a number of candidates exist (see [Barnaghi et al. 2009] for an example based upon Sensor Web Enablement (SWE) and SensorML data component models).

The use of the Bramblemet sensor network is particularly attractive since the network is subject to real network outages that generate instances of missing sensor readings on which we can evaluate our information processing algorithms. Furthermore, by working with the sensor network operators we have been able to recover the missing data from these periods from the sensors’ local caches, and thus, perform quantitative comparisons between our information processing algorithm’s predictions and the actual data. Furthermore, within the Bramblemet network sensors, the sensor measurements are rounded prior to their transmission over the wireless link. This results in censored observations, and thus, again by accessing the raw data from the sensors’ local caches, we are also able to compare the ability of our information processing algorithms to handle these censored observations.

However, the Bramblemet sensor network only consists of four sensors, and thus, in order to compare the performance of our approach on larger networks, we turn to the Wannengrat Alpine observatory and the UK Met Office MIDAS land surface stations.

4.2. Wannengrat Alpine Observatory

The Wannengrat site consists of over twenty SensorScope stations to measure local weather conditions³. This wireless sensor network is deployed with the aim of understanding the complex meteorological processes occurring on the snowpack. For this network, we simply use the data from online data repository, and have had no direct interaction with the sensor network operators.

4.3. MIDAS Land Surface Weather Stations

Finally, we use data from the UK Met Office MIDAS land surface stations data [British Atmospheric Data Centre 2009], focusing on the network’s readings of maximum daily temperature at 670 sensors distributed across the UK. Our experiments on the data collected from these several hundred sensors serve as a demonstration of the ability of our algorithm to cope with larger sensor networks.

³See <http://www.swiss-experiment.ch/index.php/Wannengrat:Sensorscopedeployment>.

```

<sit:Location rdf:about="%sit;bramblemet"
  rdfs:label="Bramble Bank"
  geo:lat="50.79472"
  geo:lng="-1.2875"
  sit:altitude="1">
</sit:Location>

<sit:Sensor rdf:about="%sit;bramblemet/windspeed"
  rdfs:label="Wind speed">
  <sit:sensorType rdf:resource="%sit;windspeed"/>
  <sit:location rdf:resource="%sit;bramblemet"/>
</sit:Sensor>

<sit:SensorType rdf:about="%sit;windspeed"
  rdfs:label="Wind speed">
</sit:SensorType>

<sit:Unit rdf:about="%sit;knots"
  rdfs:label="Knots"
  sit:unitAbbr="kn">
</sit:Unit>

<sit:Reading
  rdf:about="%sit;bramblemet/windspeed/reading/1234"
  rdfs:value="9.3"
  sit:datetime="2007-10-25T21:55:00">
  <sit:sensor rdf:resource="%sit;bramblemet/windspeed"/>
  <sit:unit rdf:resource="%sit;knots"/>
</sit:Reading>

```

Fig. 3: Example RDF data from the Bramblemet sensor.

5. EMPIRICAL EVALUATION

In this section we empirically evaluate our information processing algorithm on real weather data collected from the sensor networks described above. We will compare the use of our GP formalism built around a multiple-sensor covariance against a number of benchmarks:

- **Conventional independent GP** in which each environmental variable is modelled separately (i.e. correlations between these parameters are ignored).
- **Kalman filter** in which a maximum a posteriori (MAP) recursion is used to govern the adaptations of the hyper-parameters (namely the state and observation noise processes) by sequentially using the maximum-likelihood formulation first proposed by Jazwinski [1970]. One step-ahead forecasts are successively iterated through the algorithm to provide multi-step forecasts of arbitrary length. Missing values, if they occur, are accommodated for using the scheme advocated by Shumway and Stoffer [2000]. The algorithm operates, for multiple time-series forecasting, using a p -th order lookback into the data, including cross-coupling between time-series, making it a recursive multivariate autoregressive, $MAR(p)$, process.

- **Naïve algorithm** in which the prediction of the variable at the next time step is simply equal to that of the most recent observation at that sensor.

The first benchmark illustrates the effectiveness of our Gaussian process formalism that expresses correlations between different sensors. The second benchmark is a state-of-the-art alternative, that has previously been used for tracking environmental sensor data [Bertino et al. 2003]. Finally, the naïve third benchmark provides a lower bound on the performance of a prediction algorithm.

For the purposes of illustration, we will plot the sensor readings acquired by our algorithm (shown as markers), the mean and standard deviation of the GP prediction (shown as a solid line with plus or minus a single standard deviation shown as shading), and the true fine-grained sensor readings (shown as bold) that were downloaded directly from the sensor (rather than through the web site) after the event. Note that we present a limited number of sensors for reasons of space, but we use readings from all relevant sensors in order to perform inference.

Where appropriate, predictive quality will be evaluated using the standard root mean squared error (RMSE), defined as

$$\text{RMSE}(\mathbf{m}, \mathbf{y}) = \sqrt{\frac{1}{N} \sum_{i=1}^N (m_i - y_i)^2}.$$

Here \mathbf{m} is a vector of N estimates for the vector of variables \mathbf{y} . The estimates we use are the means for \mathbf{y} conditioned on some set of observations \mathbf{z} (as dependent on the application). We also provide the normed mean squared errors (NMSEs), defined as

$$\text{NMSE}(\mathbf{m}, \mathbf{y}) = 10 \log_{10} \frac{1/N \sum_{i=1}^N (m_i - y_i)^2}{1/N \left(\sum_{i=1}^N y_i^2 \right) - \left(1/N \sum_{i=1}^N y_i \right)^2},$$

the ratio of the mean squared error and the variance of our data, expressed in units of decibels. The NMSE provides a measure of the accuracy in our predictions that is relatively insensitive to the scale of variation (and units) of the relevant variable.

It is commonly the case that after a significant volume of data, the majority of the weight will be concentrated at a single trial hyperparameter sample – the posterior distribution for the hyperparameter $p(\phi | \mathbf{z}_d, I)$ will be close to a delta function at that value. In this case, we will say, informally, that we have *learned* the hyperparameter to have that single value. Where this appears to be the case, performing a local optimisation of the likelihood $p(\mathbf{z}_d | \phi, I)$ can help to find the point of highest posterior density.

In the subsections that follow we first evaluate our approach on tide height, air temperature and rounded air pressure data from the Bramblemet weather sensor network, and then discuss the prediction of maximum daily air temperature over the larger MIDAS land surface weather station network. In each case, we describe the covariance function used, and present quantitative comparisons of prediction quality. We then go on to discuss our approach to active data selection using the tide and wind speed sensors of the Bramblemet weather sensor network, and the air temperature sensors of the Wannangrat Alpine Observatory.

5.1. Prediction of Tide Heights

Figure 4 illustrates the efficacy our GP prediction for a tide height dataset. That is, at time t , Figure 4 depicts the posterior distribution of the GP, conditioned on all observations prior to and inclusive of t . 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

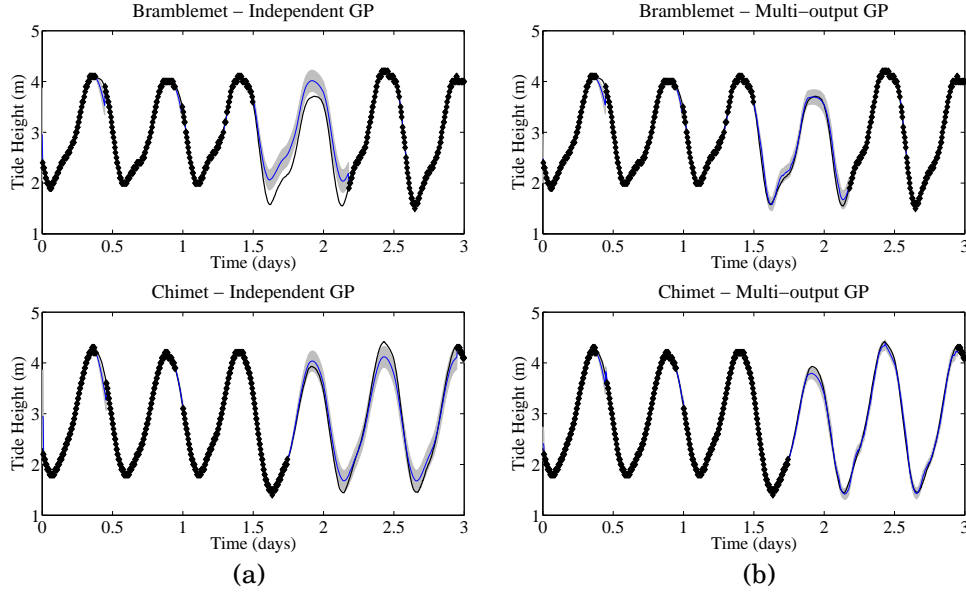


Fig. 4: Prediction and regression of tide height data for (a) independent and (b) multi-output Gaussian processes.

t , time, and l , a sensor label. As such, our covariance function is of the form

$$K([t, l], [t', l']) = W(l, l') \left(K_{\text{per-Mtn}}(t - \Delta_l, t' - \Delta_{l'}; h_P = 1, w_P, \nu = \frac{5}{2}) + K_{\text{Mtn}}(t - \Delta_l, t' - \Delta_{l'}; h_D, w_D, \nu = \frac{5}{2}) \right), \quad (20)$$

where W is a covariance matrix parameterised using the spherical parameterisation (5). We used two different GP models: one in which the parameters of the spherical decomposition were inferred from the data; and one for which the matrix S was taken as an identity matrix. In the former case, we consequently have a multi-output formalism, in the latter, independent GPs for each sensor.

Note that our covariance over time is the sum of a periodic term and a *disturbance* term. Both are of the Matérn form with $\nu = \frac{5}{2}$, (3) and (2) respectively. 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 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 w_P of the first term was unsurprisingly learnt as being about half a day, whereas for the disturbance term the time scale w_D 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 $S^T S$ (where S is from (5)) containing elements all very close to one. W additionally gives an appropriate length scale for each sensor. Over this data set, the Chimet sensor was found to have a length scale of 1.4m, with the remainder possessing scales of close to 1m. Note that h_P and h_D in (20) are dimensionless ratios, serving as weightings between our various covariance terms. Hence we can set $h_D = 1$ without loss of gener-

Table I: Predictive performances for five-day Bramblemet tide height dataset.

Algorithm	RMSE (m)	NMSE (dB)
Naïve	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

ality, upon which we learn that h_P is around 0.2. Hence weather events were observed to have induced changes in tide height on the order of 20%.

We also make allowances for the prospect of relative latency amongst the sensors by incorporating delay variables, represented by the vector Δ . 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.

Only the correlations amongst sensors and their individual delays were sampled over, all other hyperparameters for our covariance were learned by the GP prior to the depicted data set. Such hyperparameters were then specified with single samples at the learned values.

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.

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 I. For the Kalman filter, a history length of 16 observations was used to generate each prediction, since this gave rise to the best predictive ability. 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.

5.2. Prediction of Air Temperatures

We next tested on air temperature data. Figure 5 demonstrates our algorithm's 5 minute lookahead predictive performance, where 5 minutes was also the usual (al-

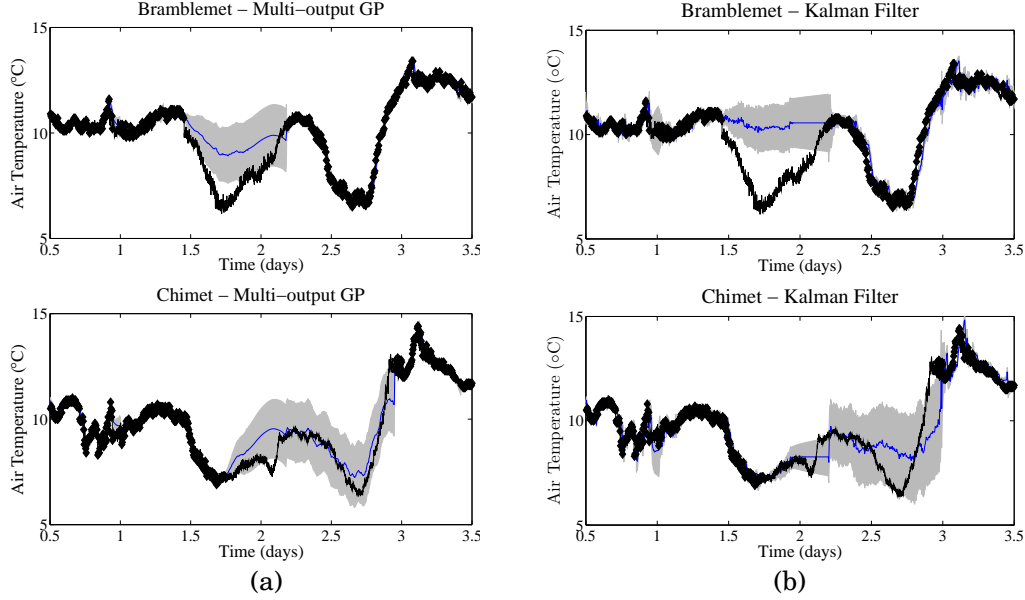


Fig. 5: One-step lookahead prediction of air temperature data for (a) a multi-output Gaussian process and (b) a Kalman filter.

though not exclusive) time increment from one datum to the next. That is, Figure 5 depicts the posterior distribution at time t conditioned on all observations prior to and inclusive of $t - 5$ mins. Note that the Cambermet station does not record air temperature, and is hence discluded from consideration here. Taking the covariance function used for tide height data shown in equation (20) as a model, we assumed the covariance function

$$\begin{aligned}
 K([t, l], [t', l']) = & W(l, l') \\
 & \left(K_{\text{per-Mtn}}(t - \Delta_l, t' - \Delta_{l'}; h_P = 1, w_P, \nu = \frac{5}{2}) + \right. \\
 & \left. K_{\text{Mtn}}(t - \Delta_l, t' - \Delta_{l'}; h_T, w_T, \nu = \frac{5}{2}) \right) + \\
 & \delta(l - l') K_{\text{Mtn}}(t, t'; h_D, w_D, \nu = \frac{3}{2}),
 \end{aligned}$$

Our covariance over time consisted of three terms. The first was periodic, whose period w_P was learned as 1 day, describing nightly dips in temperature. The second term describes longer-term trends in temperature. Its time scale w_T was learned as 3 days, with a relative magnitude h_T just over twice that of the periodic term, h_P . From this we can induce that the periodic term represents only a relatively minor component of this model. Finally, the third term represents higher frequency disturbances, with a time scale learned as 4 hours. Note that this final term is not correlated amongst sensors as were the other two; these disturbances appear to be local to each sensor. The individual output scales for each sensor given by W were learned as 2°C , 6°C and 2.6°C respectively. These give the scales of the periodic fluctuations and longer-term drifts in the air temperature at each sensor. For comparison, the scale of the disturbances h_D were learned as 1°C . As for the tide heights example, only the sensor correlations and delays were sampled over for the depicted predictions.

Table II: Predictive performances for five-day Bramblemet air temperature dataset.

Algorithm	RMSE (°C)	NMSE (dB)
Naïve	0.84	-10.0
Kalman filter	0.98	-9.0
Independent GPs	0.82	-10.5
Multi-output GP	0.45	-15.4

The predictive performances for this dataset are shown in Table II. Note that in this case a history of 15 observations proved optimal for the Kalman filter. However, the Kalman filter performs particularly badly, due to the significant lag (compared to the real temperature measurement) that the filter introduces. Again, the multi-output GP significantly outperforms the other benchmarks.

5.3. Prediction of Air Pressures

Figure 6 demonstrates regression and prediction over the heavily rounded air pressure observations from the Bramblemet sensor alone. We used a GP model with the Matérn covariance of the non-periodic form shown in equation (2). We further made use of the methods described in Section 3.3 for managing the rounded observations (to the nearest Pa) that we receive from the sensor. Here we see a dramatic improvement over Kalman filter prediction (again using a history of 16 observations to generate the predictions) as shown in Table III. The Kalman filter behaves poorly as it is unable to both capture the static behaviour when different raw measurements are rounded to the same integer value, and still respond to the sudden jumps when a change in this integer value occurs. In contrast, within our GP approach, we explicitly model the censored observations that result from this rounding, and thus, the prediction performance is significantly better.

5.4. Prediction of MIDAS data

We now turn to a larger sensor network and demonstrate the suitability of our prediction algorithm for such networks. In particular, we apply our algorithm to data from January 1990 collected by the 670 sensors that comprised UK Met Office MIDAS land surface stations at the time. We arbitrarily selected maximum air temperature as the environmental variable of interest for our experiments, and aimed to perform one-step lookahead prediction. Specifically, we aim to predict the value of the next observation to be received by the network, which may be an observation from a day equal to that of the current observation, but from a sensor different to that most recently observed. The use of a model over both space and time is integral to success on this task.

The dataset comes equipped with latitude and longitude measurements for each sensor, which we can use to express the covariance between any pair of sensors. That is, we take the covariance

$$K([t, l], [t', l']) = K_{\text{Mtn}}(l, l'; h, w_l, \nu = \frac{3}{2}) K_{\text{Mtn}}(t, t'; h = 1, w_t, \nu = \frac{3}{2}), \quad (21)$$

where for the covariance between l and l' we replace r with the great-circle distance between the positions of sensors l and l' , divided by isotropic spatial scale hyperparameter w_l . Note that our approach here is an alternative to the arbitrary parameterisation presented in equation (5). This parameterisation requires a hyperparameter for every distinct pair of sensors and is hence impractical for larger numbers of sensors.

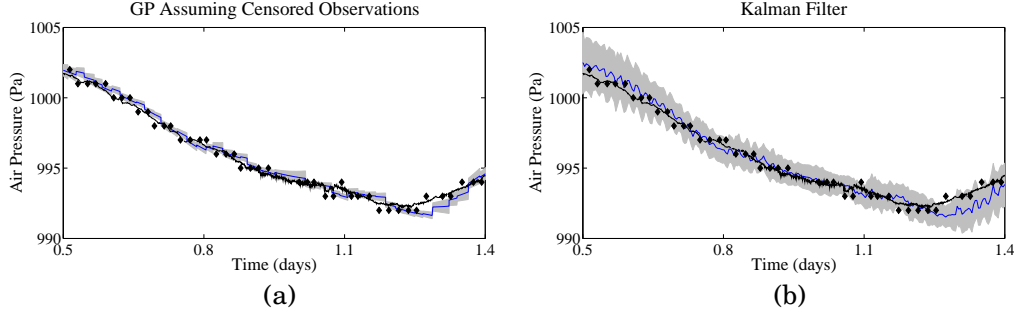


Fig. 6: Prediction and regression of air pressure data for (a) a Gaussian process employing censored observations and (b) a Kalman filter.

Table III: Predictive performances for one-day Bramblemet air pressure dataset. Note that with data from only one sensor, our multi-output GP formalism is not applicable.

Algorithm	RMSE (Pa)	NMSE (dB)
Naïve	3.61	-4.75
Kalman filter	3.29	-5.55
GP for censored data	0.43	-20.49

Table IV: Predictive performances for MIDAS maximum air temperature data.

Algorithm	RMSE (°C)	NMSE (dB)
Naïve	3.93	3.61
Basis functions	2.23	-1.32
Independent GPs	2.50	-0.31
Multi-output GP	1.53	-4.57

The multi-output GP model built using (21) is compared against an independent GPs model in which no spatial inference is performed.

Results from our tests on this dataset are displayed in Table IV. Note that the GP is easily able to outperform its naïve competitor. We also performed a comparison against a method built around basis functions (a Kalman filter is incapable of the simultaneous spatio-temporal inference required for this particular problem). A set of Gaussian basis functions, ϕ_i , were centred at the sensor locations for the last 30 observations and used to infer the value at the next sensor on the list via a linear basis model of the form $y = \sum_i w_i \phi_i$. The factor weights were inferred using a Bayesian linear regression model with shrinkage priors over all the parameters. The inference was performed via variational Bayes [Bishop 1995]. The multi-output GP demonstrated significant improvement over this alternative method.

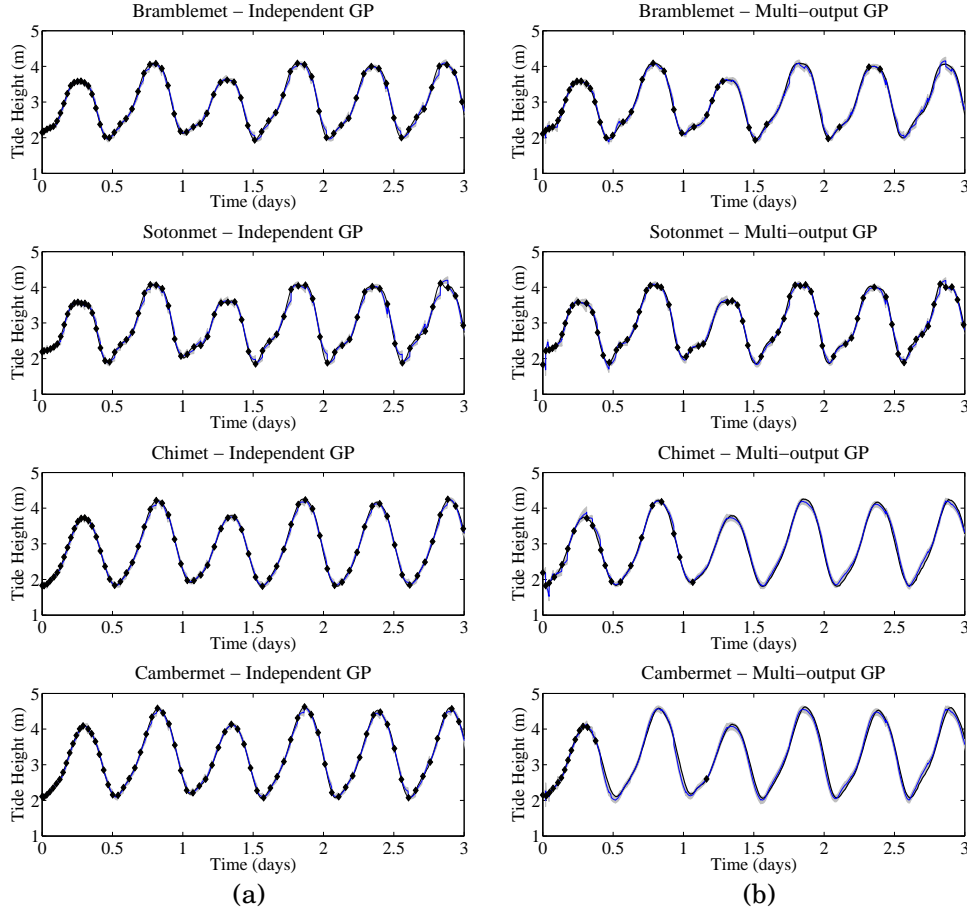


Fig. 7: Comparison of active sampling of tide data using (a) independent and (b) multi-output Gaussian processes.

5.5. Active Data Selection of Tide Heights

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 the four tide sensors are displayed in Figure 7. 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 10cm. The covariance function used was that described in (20).

Consider first the case shown in Figure 7(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. In contrast, for the multi-output case shown in Figure 7(b), the GP is allowed to explicitly represent correlations and delays between the sensors. As mentioned above, this data set is notable for the slight delay of the tide heights at the Chimet and Cambermet sensors relative to the Sotonmet and Bramblemet sen-

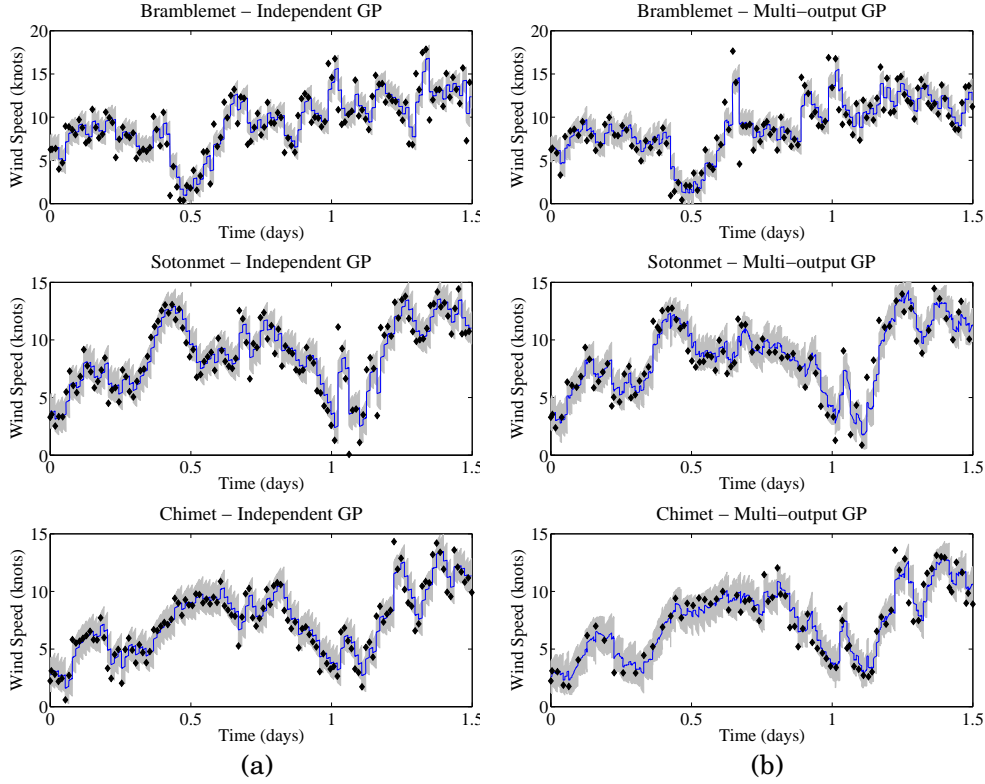


Fig. 8: Comparison of active sampling of wind speed using (a) independent and (b) multi-output Gaussian processes.

sors, 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 sensors⁴. 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.

5.6. Active Data Selection of Wind Speeds

Figure 8 shows similar results for the wind speed measurements from three of the four sensors (the Cambermet sensor being faulty during this period) where the goal was to maintain the uncertainty in wind speed below 1.5 knots. By analogy to the covariance function over tides (20), we used the following covariance function

In this case, for purposes of clarity, the fine-grained data is not shown on the plot. Note that the measurement noise is much greater in this case, and this is reflected in the uncertainty in the GP predictions. Furthermore, note that while the Sotonmet and

⁴The dynamics of the tide height at the Sotonmet sensor are more complex than the other sensors due 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.

Chimet sensors exhibit a noticeable correlation, Bramblemet appears to be relatively uncorrelated with both. This observation is reflected in the sampling that the GP performs. The independent GPs sample the Bramblemet, Sotonmet and Chimet sensors 126, 120 and 121 times respectively, while over the same period, our multi-output GP samples the same sensors 115, 88 and 81 times. Our multi-output GP learns on-line that the wind speed measurements of the Sotonmet and Chimet sensors are correlated, and then exploits this correlation in order to reduce the number of times that these sensors are sampled (inferring the wind speed at one location from observations of another). However, there is little or no correlation between the Bramblemet sensor and the other sensors, and thus, our multi-output GP samples Bramblemet almost as often as the independent GPs.

5.7. Active Data Selection of Air Temperature

To investigate the practicality of our active data selection algorithm on larger sensor networks, we apply it to data from 16 air temperature sensors from the Wannengrat Alpine Observatory. As in the case of the MIDAS data described above we use the latitude and longitude co-ordinates of these sensors within the covariance function shown in equation (21).

$$K([t, l], [t', l']) = W(l, l') K_{\text{per-Mtn}}(t - \Delta_l, t' - \Delta_{l'}; h_P = 1, w_P, \nu = \frac{1}{2}) + \delta(l - l') K_{\text{Mtn}}(t, t'; h_D, w_D, \nu = \frac{3}{2}).$$

Figure 9 shows the detailed predictions and samples for five of the sensors, and a summary plot showing the sample times of all 16 sensors. Note that the sensor sampling rate is relatively low at first, however, as the volatile fluctuations in temperature begin to occur at about $t = 0.7$ days, the sampling frequency automatically increases to ensure that the prediction uncertainty is maintained below the specified threshold.

6. COMPUTATION TIME

As described earlier, a key requirement of our algorithm is computational efficiency, in order that it can be used to represent multiple correlated sensors, and hence, used for real-time information processing. Here we consider the computation times involved in producing the results presented in the previous section. We should note at this point that we are not comparing the computation time of our approach against the Kalman filter. The computation performed by the Kalman filter is far less intensive than that of the Gaussian process (in our empirical evaluations, the update time of the Kalman filter was typically < 0.01 seconds). However, as demonstrated in the previous section, the Kalman filter is unable to represent the complex correlations and periodicities seen in the sensor data, and thus, it provides significantly worse predictions than our approach. Rather, in this section, we are evaluating the computation time required in order to update the Gaussian process as a new observation is received. This computation time represents the cost of computing the Cholesky factor \mathbf{R} of \mathbf{V} , the associated data-dependent term $\mathbf{R}(x_d, x_d)^{-1}(\mathbf{y}_d - \boldsymbol{\mu}(x_d))$ in equations (7) and (8), and the hyperparameter sample likelihoods and subsequently the weights of equation (16). Once these terms have been calculated, by either method, making predictions at any point in time is extremely fast, requiring simply the adding of another element to x_* . The conventional approach to performing this update would be to simply recompute the relevant quantities from scratch upon receipt of new data (using the highly optimised LAPACK routines within Matlab - see <http://www.netlib.org/lapack/> for more details of this linear algebra package). In contrast, in our formalism, we make use of

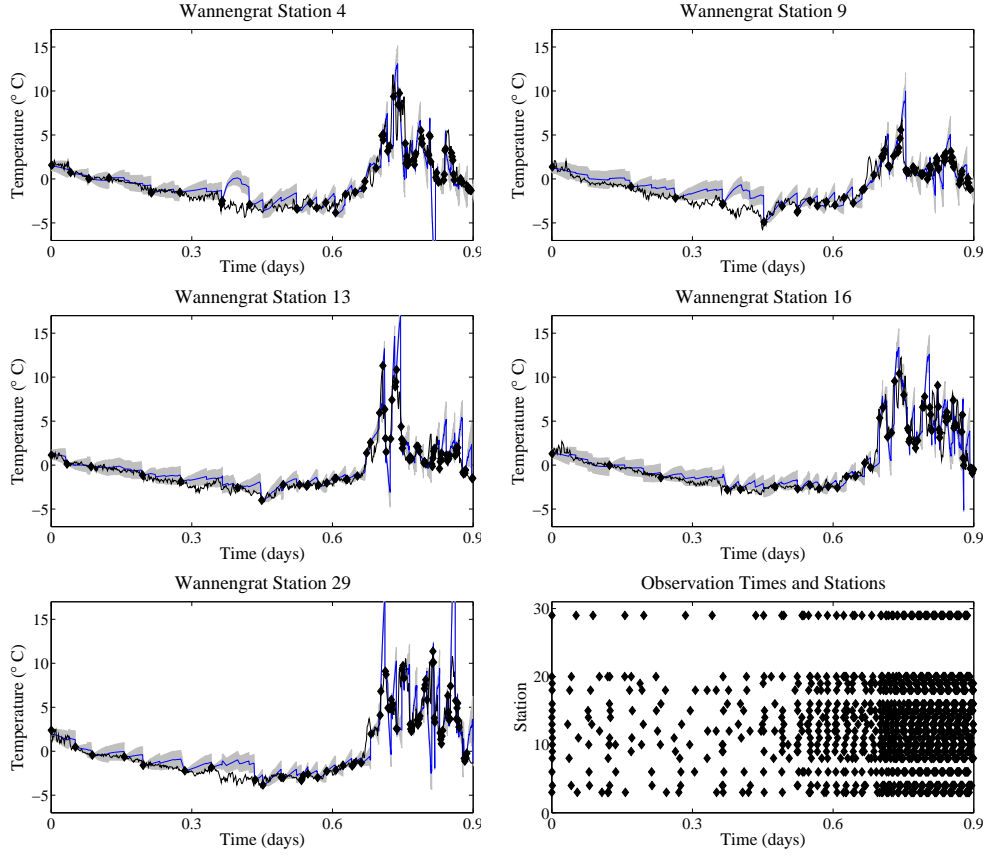


Fig. 9: Active sampling of ambient temperatures at 16 Wannengrat sensor stations.

Table V: The required computation time (seconds) per GP update, with (a) and without (b) our efficient update rules, over N the number of stored data points and η the number of hyperparameter samples.

η	N		
	10	100	500
1	< 0.01	< 0.01	< 0.01
10	0.01	0.01	0.08
100	0.11	0.13	0.74
1000	1.30	1.39	4.51

(a)

η	N		
	10	100	500
1	< 0.01	< 0.01	0.06
10	0.01	0.02	0.34
100	0.09	0.17	3.32
1000	0.92	1.71	17.34

(b)

the computationally efficient online rules described in Section 3.4 to perform an online update of previously calculated results.

Table VI: The numbers of stored data points, N , and hyperparameter samples, η , used for the various datasets presented in this paper.

Dataset	N	η
Bramblemet Tide Heights	200	1620
Bramblemet Air Temperatures	500	972
Bramblemet Air Pressures	16	125
MIDAS Air Temperature	1500	105
Bramblemet Wind Speeds	454	675
Wannengrat Air Temperature	500	125

Thus, against this background, Table V shows the computation time required using MATLAB on an Intel Core i7 Processor 860 (2.80GHz, 8MB) with 8GB of RAM, with displayed results the average of 10 trials. In more detail, Table V(a) the results when our computationally efficient online rules are used, and Table V(b) shows the results when the conventional approach is used.

Note that with our efficient updates, we expect the cost of computation to grow as $O(N^2)$ in the number of stored data points. Our proposed algorithm will automatically determine the quantity of data to store in order to achieve the desired level of accuracy. In the problems we have studied, a few hundred points were typically sufficient. The largest number we required was 1500, for the MIDAS air temperature data. Note also that the cost of computing equation (19) will grow in the cube of the number of samples in each hyperparameter. However, we consider only a fixed set of samples in each hyperparameter, and thus, equation (19) need only be computed once, off-line. In this case, our on-line costs are limited by the multiplication of that term by the likelihoods r_s to give the weights of equation (16), and this only grows as $O(\eta^2)$. Furthermore, note that this cost is independent of how the η samples are distributed amongst the hyperparameters. For completeness, Table VI shows the number of data points and the number of hyperparameter samples used for each of the empirical evaluations.

In practice, we see that our efficient updates result in substantial improvements in computational speed once the quantity of data N exceeds the threshold required to offset the small computational overheads introduced by using our approach instead of the already highly optimised LAPACK routines. For the largest cases investigated here we see an approximate four-fold increase in computation speed. For our weather datasets, samples are taken no more frequently than once a minute. Thus all of the information processing that we have described in this paper can easily be performed in real-time, and the increase in computational speed seen in our approach can be used to increase the number of data points considered or the number of hyperparameters samples used; both of which will likely increase the predictive power of the Gaussian process.

7. TRIAL DEPLOYMENT

The results in the previous section indicate that the Gaussian process information processing approach developed here is relatively computationally expensive, and thus, is unlikely to be implementable on very low-power embedded devices. Due to this, and the fact that modelling the correlations between any two sensors requires the full history of observations from each sensor, it is most appropriate that our methodology

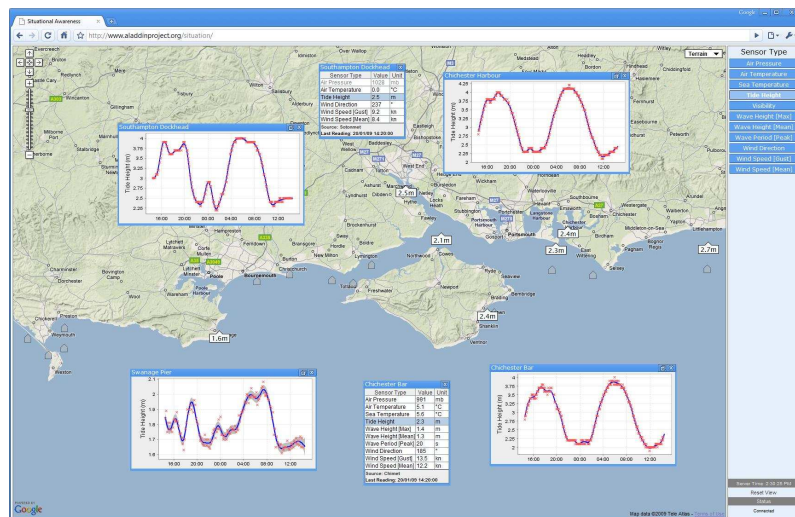


Fig. 10: Server-based implementation of our information processing algorithm accessed through a web based Google Maps interface (available at <http://www.aladdinproject.org/situation/>).

is deployed at the sensor data repository where the sensor data is finally stored⁵, or on the sink node of the sensor network (assuming that this node has more significant computational resources).

As a concrete example of the first setting, we have implemented a server-side application that collects the RDF data from the 4 Bramblemet sensors and from 30 additional sensors that form part of the Channel Coastal Observatory (located on the south coast of England). This application parses the RDF from the Bramblemet sensors, and uses an API to collect it from the Channel Coastal Observatory sensors, before storing it in RDF format using Jena (see <http://jena.sourceforge.net/>). The sensor data is then presented in tabular and graphical forms using a Google Maps interface. A live version of this system is available at <http://www.aladdinproject.org/situation/> and is shown in Figure 10.

The application collects data from the sensors every fifteen minutes and applies the Gaussian process prediction presented here across a number of different sensor readings (including tide height, air temperature, and wind speed). Taking tide height as an example, the multi-output GP is used on the 4 Bramblemet sensors, and the independent GP prediction on the other 30 sensors, to predict the tide height over the next fifteen minutes. This computation takes less than one minute to complete, and enables the system to present graphs that update every fifteen seconds ensuring that they are always showing the most up-to-date estimate of the current tide height. If sensor readings are unavailable, due to sensor failure or communication loss, then the prediction process continues, filling in the missing interval with the results of the GP prediction.

Now, the sensors deployed within the Bramblemet and Channel Coastal Observatory networks are relatively large devices, having solar panels and batteries to ensure that

⁵Current examples of such repositories include Weather Underground (<http://www.wunderground.com>), pachube (<http://www.pachube.com/>) and Microsoft Research's SensorMap (<http://atom.research.microsoft.com/sensormap/>).

they can continue to operate indefinitely. In other application domains, in order to minimise the energy consumed by the sensors, it might be necessary to use the output of the GP prediction to perform active sensing (as described in Sections 3.5 and 5.5) whereby the server sends signals back to individual sensors telling them when they should make another measurement.

For example, we are currently exploring the application of this approach within water distribution networks. In this setting, battery powered sensors use GPRS modems to report mains water pressure at regular intervals to a server. These sensor readings are then used to actively manage the mains water pressure in order to satisfy minimum water pressure requirements, while also minimising leakage losses (which are typically proportional to the pressure in the water main). The sensors sleep between taking readings to conserve battery power and maximise the useful life of the sensor (or equivalently, to minimise the frequency at which the sensors must be revisited to change their batteries). Using the Gaussian process algorithms presented here, it is possible to explicitly model the correlations between the sensors, and then choose the interval at which each sensor should make a reading. In this way, the minimum number of sensor readings can be taken whilst ensuring that the level of uncertainty regarding the true water pressure is not exceeded.

Similarly, we are also exploring the application of the Gaussian process algorithm on the sink node of a body sensor network. In this setting, small battery powered sensors are deployed on the body to record physiological measurements (such as skin temperature, conductivity, and pulse rate) which are transmitted to the sink node (in this case a smart phone running a dedicated application) where they are stored for later upload. Again, by modelling correlations explicitly using the multi-output GP described here, the sampling rate of the individual sensors can be minimised, extending their battery life without maximising the amount of information collected.

8. RELATED WORK

Gaussian process regression has a long history of use within geophysics and geospatial statistics (where the process is known as kriging [Cressie 1991]), but has only recently been applied within sensor networks. For example, Gaussian processes have been used to represent spatial correlations between sensors so that the near-optimal sensor placement (in terms of mutual information) can be determined, for modelling wireless propagation between sensor nodes subject to censored observations of received signal strength [Ertin 2007], and in the form of multi-variate Gaussian distributions to represent correlations between different sensors and sensor types for energy efficient querying of a sensor network [Deshpande et al. 2004]. They have also been used to represent temporal correlations between the readings from a single sensor in order to perform adaptive sampling — maximising the information gain subject to a sampling constraint [Kho et al. 2009].

Our work differs in that we use GPs to represent temporal correlations, and represent correlations and delays between sensors with additional hyperparameters. It is thus closely related to other work using GPs to perform regression over multiple responses [Boyle and Frean 2005; Teh et al. 2005]. However, our focus is to derive a computationally efficient algorithm, and thus, we use a number of novel computational techniques to allow the re-use of previous calculations as new sensor observations are made. We additionally use a novel Bayesian Monte Carlo technique to marginalise the hyperparameters that describe the correlations and delays between sensors. Finally, we use the variance of the GP's predictions in order to perform active data selection. Likewise, our work differs from that previously using censored sensor readings within a GP framework [Ertin 2007], since our work proposes a principled Bayesian Monte Carlo method for adapting our models to the data, where Ertin's model assumes that

the hyperparameters are known a priori, and hence, does not consider how likelihoods should be computed in this context. Furthermore, in our work, Monte Carlo techniques are also used to evaluate our other integrals, rather than taking Laplace approximations, allowing more accurate representation of the correlation amongst censored readings.

Finally, our approach has several advantages relative to sequential state-space models such as the Kalman filter [Girard et al. 2003; Jazwinski 1970] which have also been applied in environmental settings for tracking sensor readings [Bertino et al. 2003]. Firstly, these state-space models require the discretisation of the time input, representing a discarding of potentially valuable information. Secondly, their sequential nature means they must necessarily perform difficult iterations in order to manage missing or late data, or to produce long-range forecasts. In our GP approach, what observations we have are readily managed, regardless of when they were made. Equally, the computation cost of all our predictions is identical, irrespective of the time or place we wish to make them about. Finally, a sequential framework requires an explicit specification of a transition model. In our approach, we are able to learn a model from data even if our prior knowledge is negligible, and the benefits of our approach are empirically supported by the results presented in Section 5.

9. CONCLUSIONS

In this paper we addressed the need for algorithms capable of performing real-time information processing of sensor network data, and we presented a novel computationally efficient formalism of a multi-output Gaussian process. Using weather data collected from three sensor networks, we demonstrated that this formalism allows an end-user to make effective use of sensor data even in the presence of network outages and sensor failures (recovering missing data by making predictions based on previous sensor readings and the current readings of sensors that are functioning), and to automatically determine the sampling rate of each individual sensor in order to ensure that the uncertainty in the environmental parameter being measured stays within a pre-specified limit. Furthermore, we showed that our formalism that efficiently re-uses previous computations by following an online update procedure as new data sequentially arrives results a significant increase in computation speed compared to the highly optimised LAPACK linear algebra package.

Our future work in this area consists of three areas. First, we would like to investigate the use of schemes more sophisticated than the maximum likelihood approach for the movement of the locations of our set of hyperparameter samples. As the posterior distributions of these hyperparameters become more sharply peaked, we might aim to reduce the number of samples to further increase the computational efficiency of our algorithm.

Second, we intend to investigate the use of correlations between different sensor types (rather than between different sensors of the same type as presented here) to perform regression and prediction within our weather sensor network. Our formalism is in no way restricted to identical sensors, and it is expected that many sensors will exhibit correlations in this setting. For example, both wind speed and wind direction, and air temperature and air pressure are likely to exhibit correlations as weather fronts move over the sensor network.

Finally, we would like to use the probabilistic model that the GP builds to automatically handle faulty and unreliable sensors within the network. Note that we do not wish to simply detect these faulty sensors, but would like to simultaneously perform both detection and prediction, despite the presence of these failures. Our preliminary work in this area indicates that nonstationary covariance functions that model such changes can be introduced to the formalism described here to achieve this, and have

already been shown to work effectively on real-world sensor data derived from the sensor networks described in this paper [Garnett et al. 2009].

Acknowledgments

This research was undertaken as part of the ALADDIN (Autonomous Learning Agents for Decentralised Data and Information Networks) project, jointly funded by a BAE Systems and EPSRC strategic partnership (EP/C548051/1), and was also part of the Systems Engineering for Autonomous Systems (SEAS) Defence Technology Centre established by the UK Ministry of Defence. We would like to thank B. Blaydes of the Bramblemet/Chimet Support Group, and W. Heaps of Associated British Ports (ABP) for allowing us access to the weather sensor network, hosting our RDF data on the sensor web sites, and for providing raw sensor data as required. We also thank the UK Met office for the use of the MIDAS dataset.

References

- ABRAHAMSEN, P. 1997. A review of Gaussian random fields and correlation functions. Tech. Rep. 917, Norwegian Computing Center, Box 114, Blindern, N-0314 Oslo, Norway. 2nd edition.
- BARNAGHI, P., MEISSNER, S., PRESSER, P., AND MOESSNER, K. 2009. Sense and sens'ability: Semantic data modelling for sensor networks. In *Proceedings of the ICT Mobile and Wireless Communications Summit*. Santander, Spain.
- BERTINO, L., EVENSEN, G., AND VACKERNAGEL, H. 2003. Sequential data assimilation techniques in oceanography. *International Statistical Review* 71, 223–242.
- BISHOP, C. 1995. *Neural networks for pattern recognition*. Oxford university press.
- BOYLE, P. AND FREAN, M. 2005. Dependent Gaussian processes. In *Advances in Neural Information Processing Systems 17*. The MIT Press, 217–224.
- BRITISH ATMOSPHERIC DATA CENTRE. 2009. UK Meteorological Office MIDAS Land Surface Stations data (1853-current). Available from: <http://badc.nerc.ac.uk/data/ukmo-midas>.
- CRESSIE, N. A. C. 1991. *Statistics for spatial data*. John Wiley & Sons.
- DESHPANDE, A., GUESTRIN, C., MADDEN, S., HELLERSTEIN, J., AND HONG, W. 2004. Model-driven data acquisition in sensor networks. In *Proceedings of the Thirtieth International Conference on Very Large Data Bases (VLDB 2004)*. Toronto, Canada, 588–599.
- ERTIN, E. 2007. Gaussian Process Models for Censored Sensor Readings. In *Proceedings of the 14th IEEE/SP Workshop on Statistical Signal Processing*. Madison, Wisconsin, USA, 665–669.
- FUENTES, M., CHAUDHURI, A., AND HOLLAND, D. H. 2007. Bayesian entropy for spatial sampling design of environmental data. *Environmental and Ecological Statistics* 14, 323–340.
- GARNETT, R., OSBORNE, M. A., AND ROBERTS, S. J. 2009. Sequential bayesian prediction in the presence of changepoints. In *Proceedings of the 26th International Conference on Machine Learning*. Montreal, Canada.
- GENZ, A. 1992. Numerical Computation of Multivariate Normal Probabilities. *Journal of Computational and Graphical Statistics* 1, 2, 141–149.
- GIRARD, A., RASMUSSEN, C., CANDELA, J., AND MURRAY-SMITH, R. 2003. Gaussian process priors with uncertain inputs – application to multiple-step ahead time series forecasting. In *Advances in Neural Information Processing Systems 16*. MIT Press, Cambridge, MA, 545–552.
- HART, J. K. AND MARTINEZ, K. 2006. Environmental Sensor Networks: A revolution in the earth system science? *Earth-Science Reviews* 78, 177–191.
- JAZWINSKI, A. 1970. *Stochastic processes and filtering theory*. Academic Press New York.
- KHO, J., ROGERS, A., AND JENNINGS, N. R. 2009. Decentralized control of adaptive sampling in wireless sensor networks. *ACM Transactions on Sensor Networks* 5, 3, 1–35.
- KRAUSE, A., GUESTRIN, C., GUPTA, A., AND KLEINBERG, J. 2006. Near-optimal sensor placements: Maximizing information while minimizing communication cost. In *Proceedings of the Fifth International Conference on Information Processing in Sensor Networks*. Nashville, Tennessee, USA, 2–10.
- LASSILA, O. AND SWICK, R. R. 1999. Resource description framework (rdf) model and syntax specification. Available at <http://www.w3.org/TR/1999/REC-rdf-syntax-19990222/>.
- MACKAY, D. J. C. 2002. *Information Theory, Inference & Learning Algorithms*. Cambridge University Press.
- O'HAGAN, A. 1987. Monte Carlo is fundamentally unsound. *The Statistician* 36, 247–249.

- PADHY, P., DASH, R. K., MARTINEZ, K., AND JENNINGS, N. R. 2010. A utility-based adaptive sensing and multi-hop communication protocol for wireless sensor networks. *ACM Transactions on Sensor Networks* 6, 3. In print.
- PINHEIRO, J. AND BATES, D. 1996. Unconstrained parameterizations for variance-covariance matrices. *Statistics and Computing* 6, 289–296.
- RASMUSSEN, C. E. AND GHAHRAMANI, Z. 2003. Bayesian Monte Carlo. In *Advances in Neural Information Processing Systems 15*. The MIT Press, 489–496.
- RASMUSSEN, C. E. AND WILLIAMS, C. K. I. 2006. *Gaussian Processes for Machine Learning*. MIT Press.
- SASENA, M. J. 2002. Flexibility and Efficiency Enhancements for Constrained Global Design Optimization with Kriging Approximations. Ph.D. thesis, University of Michigan.
- SHUMWAY, R. AND STOFFER, D. 2000. *Time Series Analysis and its applications*. Springer.
- STEIN, M. 2005. Space-Time Covariance Functions. *Journal of the American Statistical Association* 100, 469, 310–322.
- TEH, Y. W., SEEGER, M., AND JORDAN, M. I. 2005. Semiparametric latent factor models. In *Proceedings of the Conference on Artificial Intelligence and Statistics*. Barbados, 333–340.
- THE MATHWORKS. 2007. MATLAB R2007a. Natick, MA.

A. APPENDIX

A.1. Cholesky Factor Update

We have a positive definite matrix, represented in block form as $\begin{bmatrix} V_{1,1} & V_{1,3} \\ V_{1,3}^T & V_{3,3} \end{bmatrix}$ and its Cholesky factor, $\begin{bmatrix} R_{1,1} & R_{1,3} \\ 0 & R_{3,3} \end{bmatrix}$. Given a new positive definite matrix, which differs from the old only in the insertion of some new rows and columns, $\begin{bmatrix} V_{1,1} & V_{1,2} & V_{1,3} \\ V_{1,2}^T & V_{2,2} & V_{2,3} \\ V_{1,3}^T & V_{2,3}^T & V_{3,3} \end{bmatrix}$, we wish to efficiently determine its Cholesky factor, $\begin{bmatrix} S_{1,1} & S_{1,2} & S_{1,3} \\ 0 & S_{2,2} & S_{2,3} \\ 0 & 0 & S_{3,3} \end{bmatrix}$. For \mathbf{A} triangular, we define $\mathbf{x} = \mathbf{A} \setminus \mathbf{b}$ as the solution to the equations $\mathbf{A} \mathbf{x} = \mathbf{b}$ as found by the use of backwards or forwards substitution. The following rules are readily obtained

$$S_{1,1} = R_{1,1} \quad (22)$$

$$S_{1,2} = R_{1,1}^T \setminus V_{1,2} \quad (23)$$

$$S_{1,3} = R_{1,3} \quad (24)$$

$$S_{2,2} = \text{chol}(V_{2,2} - S_{1,2}^T S_{1,2}) \quad (25)$$

$$S_{2,3} = S_{2,2}^T \setminus (V_{2,3} - S_{1,2}^T S_{1,3}) \quad (26)$$

$$S_{3,3} = \text{chol}(R_{3,3}^T R_{3,3} - S_{2,3}^T S_{2,3}) . \quad (27)$$

By setting the appropriate row and column dimensions (to zero if necessary), this allows us to efficiently determine the Cholesky factor given the insertion of rows and columns in any position.

A.2. Data Term Update

We have all terms defined in Section A.1, in addition to the data

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix}$$

and the product

$$\begin{bmatrix} C_1 \\ C_3 \end{bmatrix} \triangleq \begin{bmatrix} R_{1,1} & R_{1,3} \\ 0 & R_{3,3} \end{bmatrix}^{-T} \begin{bmatrix} Y_1 \\ Y_3 \end{bmatrix} .$$

To efficiently determine the updated product

$$\begin{bmatrix} D_1 \\ D_2 \\ D_3 \end{bmatrix} \triangleq \begin{bmatrix} S_{1,1} & S_{1,2} & S_{1,3} \\ 0 & S_{2,2} & S_{2,3} \\ 0 & 0 & S_{3,3} \end{bmatrix}^{-T} \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix},$$

we have

$$D_1 = C_1 \quad (28)$$

$$D_2 = S_{2,2}^{-T} (Y_2 - S_{1,2}^T C_1) \quad (29)$$

$$D_3 = S_{3,3}^{-T} (R_{3,3}^T C_3 - S_{2,3}^T D_2). \quad (30)$$

A.3. Log-Gaussian Update

We have all terms defined in Sections A.1 and A.2, in addition to

$$\begin{aligned} K &= \log \mathbf{N} \left(\begin{bmatrix} Y_1 \\ Y_3 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} V_{1,1} & V_{1,3} \\ V_{1,3}^T & V_{3,3} \end{bmatrix} \right) \\ &= -\frac{1}{2} \text{tr} \log(\sqrt{2\pi} R_{1,1}) - \frac{1}{2} \text{tr} \log(\sqrt{2\pi} R_{3,3}) - \frac{1}{2} C_1^T C_1 - \frac{1}{2} C_3^T C_3. \end{aligned}$$

We can then calculate the updated term

$$\begin{aligned} L &= \log \mathbf{N} \left(\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} V_{1,1} & V_{1,2} & V_{1,3} \\ V_{1,2}^T & V_{2,2} & V_{2,3} \\ V_{1,3}^T & V_{2,3}^T & V_{3,3} \end{bmatrix} \right) \\ &= -\frac{1}{2} \text{tr} \log(\sqrt{2\pi} S_{1,1}) - \frac{1}{2} \text{tr} \log(\sqrt{2\pi} S_{2,2}) - \frac{1}{2} \text{tr} \log(\sqrt{2\pi} S_{3,3}) \\ &\quad - \frac{1}{2} D_1^T D_1 - \frac{1}{2} D_2^T D_2 - \frac{1}{2} D_3^T D_3 \\ &= K + \frac{1}{2} \text{tr} \log(\sqrt{2\pi} R_{3,3}) - \frac{1}{2} \text{tr} \log(\sqrt{2\pi} S_{2,2}) - \frac{1}{2} \text{tr} \log(\sqrt{2\pi} S_{3,3}) \\ &\quad + \frac{1}{2} C_3^T C_3 - \frac{1}{2} D_2^T D_2 - \frac{1}{2} D_3^T D_3 \end{aligned} \quad (31)$$

A.4. Cholesky Factor Downdate

We have a positive definite matrix, represented in block form as $\begin{bmatrix} V_{1,1} & V_{1,2} & V_{1,3} \\ V_{1,2}^T & V_{2,2} & V_{2,3} \\ V_{1,3}^T & V_{2,3}^T & V_{3,3} \end{bmatrix}$ and its

Cholesky factor, $\begin{bmatrix} S_{1,1} & S_{1,2} & S_{1,3} \\ 0 & S_{2,2} & S_{2,3} \\ 0 & 0 & S_{3,3} \end{bmatrix}$. Given a new positive definite matrix, which differs from the

old only in the deletion of some new rows and columns, $\begin{bmatrix} V_{1,1} & V_{1,3} \\ V_{1,3}^T & V_{3,3} \end{bmatrix}$, we wish to efficiently deter-

mine its Cholesky factor $\begin{bmatrix} R_{1,1} & R_{1,3} \\ 0 & R_{3,3} \end{bmatrix}$. The following rules are readily obtained

$$R_{1,1} = S_{1,1} \quad (32)$$

$$R_{1,3} = S_{1,3} \quad (33)$$

$$R_{3,3} = \text{chol}(S_{2,3}^T S_{2,3} + S_{3,3}^T S_{3,3}). \quad (34)$$

Note that the special structure of equation (34) can be exploited for the efficient resolution of the required Cholesky operation, as, for example, in the MATLAB function cholupdate [The MathWorks 2007]. By setting the appropriate row and column dimensions (to zero if necessary), this allows us to efficiently determine the Cholesky factor given the deletion of rows and columns in any position.