

---

# A Machine-Learning Method for Pattern Detection, Cleaning, and Prediction of Environmental Data

---

Anonymous Author 1  
Unknown Institution 1

Anonymous Author 2  
Unknown Institution 2

Anonymous Author 3  
Unknown Institution 3

## Abstract

Many signals of interest are corrupted by faults of an unknown type. We propose an approach that uses Gaussian processes and a general “fault bucket” to capture *a priori* uncharacterised faults, along with an approximate method for marginalising the potential faultiness of all observations. This gives rise to an efficient, flexible algorithm for the detection and automatic correction of faults. Our method is deployed in the domain of water monitoring and management, where it is able to solve several fault detection, correction, and prediction problems. The method works well despite the fact that the data is plagued with numerous difficulties, including missing observations, multiple discontinuities, non-stationarity, nonlinearity and many unanticipated types of fault.

## 1 Introduction

Water sustainability is one of the greatest challenges that humankind faces. It is also a problem to which the machine-learning community can make positive, significant contributions. Water sustainability begins with proper water monitoring, which requires the analysis and interpretation of vast amounts of environmental data (WAGNER and US GEOLOGICAL SURVEY, 2006). In this paper, we attack the problem of fault detection, correction, and prediction in water monitoring signals. Here measurements are often corrupted in non-trivial ways by various intermittent faulty sensing and communication mechanisms, giving rise to outliers, telemetry spikes, missing data, drift, and multiple unanticipated exogenous disturbances (see Figure 1).

Further, signals are not well-modelled by simple parametric approaches, such as linear or Markovian models. Despite the enormous importance of such monitoring, appropriate machine-learning techniques are yet to be deployed for this purpose. In particular, there is a clear need for flexible algorithms, able to cope with signals and faults of many different types without placing a significant model-building burden upon users. Such algorithms must also be able to run reliably in real-time on incoming data. These techniques will enable us to provide operators with high-level summaries for better decision support and, in the future, to increase the level of automation and efficiency in water-management systems.

The collection of literature on fault- (also known as novelty-, anomaly- and one-class-) detection is vast (ECIOLOZA, *et al.* (2007); DE FREITAS, *et al.* (1996); ISERMANN (2005); DING (2008); MARKOU and SINGH (2003); CHANDOLA, *et al.* (2009); KHAN and MADDEN (2010); DERESZYNSKI and DIETTERICH (2011)). Unfortunately, the problems solved by most of these techniques are of very different character to our own, rendering such techniques inapplicable. Further, after much experimentation with those methods that are applicable (some of the results appear in our experiments section) it became clear that off-the-shelf techniques could not satisfy our requirements for reliable water monitoring. This was predominately due to excessively restrictive assumptions (e.g., that signals were linear, Markov or Gaussian), and/or a failure to produce reasonable uncertainty estimates. Green-tech areas, including environmental monitoring and energy-demand prediction, are still far from full automation; the provision of uncertainty estimates is necessary to allow human operators to make appropriate decisions. For this reason, we focus on developing probabilistic nonlinear models of the signal. In addition to providing posterior probabilities of observation faultiness, we are able to perform effective prediction for the latent process even in the presence of faults.

Our proposed method will rely on Gaussian processes (GPs) due to their flexibility and widely demonstrated

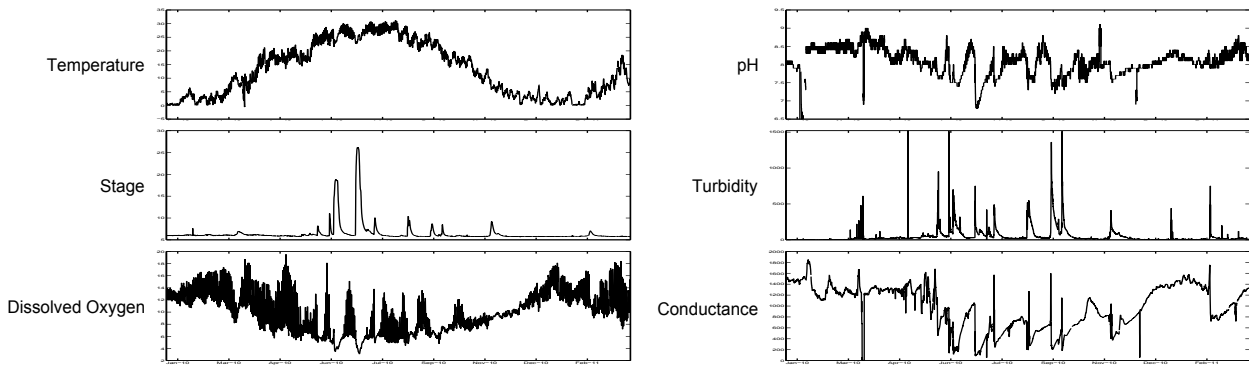


Figure 1: 16 months worth of data from six representative signals in water quality monitoring, which corresponds to approximately 11 000 measurements per series. These signals are highly nonlinear, demonstrating periodicity at different scales, intermittent pulses, and changes in dynamics. Not only do these signals exhibit a wide range of dynamics, but different signals of the same measurement type can also differ drastically if they are taken in different regions.

effectiveness at modeling nonlinear distributions. GPs have been used previously for fault detection in (ECIOLAZA, *et al.*, 2007), but in a very different context, unsuitable for our problem. Previous work along similar lines has approached this problem by creating observation models that specify the anticipated potential fault types *a priori* (GARNETT, *et al.*, 2010), but this is usually an unreasonable assumption in highly variable or poorly understood environments. In our proposed “fault bucket” approach, each point is considered to have been generated from either a nominal or generic faulty process; we do not require the specification of precise fault models. In this way, our model can simultaneously identify anomalies and robustly make predictions in the presence of sensor faults. The result is a fast and efficient method for data-stream prediction that can manage a wide range of faults without requiring significant domain-specific knowledge.

## 2 Gaussian Processes

Gaussian processes provide a simple, flexible framework for performing Bayesian inference about functions (RASMUSSEN and WILLIAMS, 2006). A Gaussian process is a distribution on the functions  $f: \mathcal{X} \rightarrow \mathbb{R}$  (on an arbitrary domain  $\mathcal{X}$ ) with the property that the distribution of the function values at a finite subset of points  $F \subseteq \mathcal{X}$  are multivariate Gaussian distributed. A Gaussian process is completely defined by its first two moments: a mean function  $\mu: \mathcal{X} \rightarrow \mathbb{R}$  and a symmetric positive semidefinite covariance function  $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . The mean function describes the overall trend of the function and is typically set to a constant for convenience. The covariance function describes how function values are correlated as a function of their locations in the domain, thereby encapsulating information about the

overall shape and behavior of the signal. Many covariance functions are available to model a wide variety of anticipated signals.

Suppose we have chosen a Gaussian process prior distribution on the function  $f: \mathcal{X} \rightarrow \mathbb{R}$ , and a set of input points  $\mathbf{x}$ , the prior distribution on  $\mathbf{f} = f(\mathbf{x})$  is

$$p(\mathbf{f} | \mathbf{x}, \theta) = \mathcal{N}(\mathbf{f}; \mu(\mathbf{x}; \theta), K(\mathbf{x}, \mathbf{x}; \theta)),$$

where  $K(\mathbf{x}, \mathbf{x}; \theta)$  is the Gram matrix of the points  $\mathbf{x}$ , and  $\theta$  is a vector containing any parameters required of  $\mu$  and  $K$ , which form hyperparameters of the model.

Exact measurements of the latent function are typically not available. Let  $y(x)$  represent the value of an observation of the signal at  $x$  and  $f(x)$  represent the value of the unknown true latent signal at that point. When the observation mechanism is not expected to experience faults, the usual noise model used is

$$p(y | f, x, \sigma_n^2) = \mathcal{N}(y; f, \sigma_n^2), \quad (1)$$

which represents additive i.i.d. Gaussian observation noise with variance  $\sigma_n^2$ . Note that this model is inappropriate when sensors can experience faults, which corrupt the relationship between  $y$  and  $f$ .

With the observation model above, given a set of observations  $\mathcal{D} = \{(x, y(x))\} = (\mathbf{x}, \mathbf{y})$ , the posterior distribution (itself given by a Gaussian Process) of  $f_\star = f(x_\star)$  given these data is

$$p(f_\star | \mathbf{y}, \theta) = \mathcal{N}(f_\star; m(f_\star | \mathbf{y}, \theta), C(f_\star | \mathbf{y}, \theta)),$$

where the posterior mean and covariance are

$$\begin{aligned} m(f_\star | \mathbf{y}, \theta) &= \mu(x_\star; \theta) + K(x_\star, \mathbf{x}; \theta) V^{-1} (\mathbf{y} - \mu(\mathbf{x}; \theta)) \\ C(f_\star | \mathbf{y}, \theta) &= K(x_\star, x_\star; \theta) - K(x_\star, \mathbf{x}; \theta) V^{-1} K(\mathbf{x}, x_\star; \theta), \end{aligned}$$

for  $V = K(\mathbf{x}, \mathbf{x}; \theta) + \sigma_n^2 \mathbf{I}$ .

We now make some definitions for the sake of readability. Henceforth, we assume that our observations  $\mathbf{y}$  have already been scaled by the subtraction of the prior mean  $\mu(\mathbf{x}; \theta)$ . We will also make use of the covariance matrix shorthand  $K_{m,n} = K(\mathbf{x}_m, \mathbf{x}_n)$ . Finally, for now, we'll drop the explicit dependence of our probabilities on the hyperparameters  $\theta$  (it will be implicitly assumed that all quantities are conditioned on knowledge of them) and will return to them later. Similarly, we drop the dependence of our probabilities on the values of inputs  $x$ , which we assume are always known.

### 3 Fault Bucket

We propose an algorithm that is designed to deal with faults of many different, unspecified types. We use a sequential scheme, applicable for ordered data such as time series, partitioning the data available at any point into old and new halves. We then approximately marginalise the faultiness of old observations, storing and then updating our results for future use. This gives rise to an efficient and fast algorithm. In order to effect our scheme, we make four key approximations:

1. **Fault bucket:** Faulty observations are assumed to be generated from a Gaussian noise distribution with a very wide variance.
2. **Single-Gaussian marginal:** A mixture of Gaussians, weighted by the posterior probabilities of faultiness of old data, is approximated as a single moment-matched Gaussian.
3. **Old/new noise independence:** We assume that noise contributions are independent, and that the contributions for new data are independent of old observations.
4. **Affine precision:** The precision matrix over both old and new halves is assumed to be affine in the precision matrix over the old half.

Approximation 1 represents the state-of-the-art DERESZYNSKI and DIETTERICH (2011). However, using it alone will not give an algorithm that can scale to the real-time problems we consider. Our novel approximations 2-4 permit very fast, fault-tolerant inference. We will detail and justify these approximations further below.

Our single, catch-all, "fault bucket" is expressed by approximation 1. It is built upon the expectation that points that are more likely to have been generated by noise with wide variance than under the normal predictive model of the GP can reasonably be assumed

to be corrupted in some way, assuming we have a good understanding of the latent process. It is hoped that a very broad class of faults can be captured in this way. To formalise this idea, we choose an observation noise distribution to replace (1) that models the noise as independent but not identically distributed with separate variances for the non-fault and fault cases:

$$\begin{aligned} p(y|f, x, \neg \text{fault}, \sigma_n^2) &= \mathcal{N}(y; f, \sigma_n^2) \\ p(y|f, x, \text{fault}, \sigma_f^2) &= \mathcal{N}(y; f, \sigma_f^2), \end{aligned} \quad (2)$$

where  $\text{fault} \in \{0, 1\}$  is a binary indicator of whether the observation  $y(x)$  was faulty and  $\sigma_f > \sigma_n$  is the standard deviation around the mean of faulty measurements. The values of both  $\sigma_n$  and  $\sigma_f$  form hyperparameters of our model and are hence included in  $\theta$ .

Of course, *a priori*, we do not know whether an observation will be faulty. Unfortunately, managing our uncertainty about the faultiness of all available observations is a challenging task. With  $N$  observations, there are  $2^N$  possible assignments of faultiness; it is infeasible to consider them all.

Our solution is founded upon approximation 2. For time series, the value to be predicted  $f_*$  typically lies in the future, and old observations are typically less pertinent for this task than new ones. We hence approximately marginalise the faultiness of old observations, representing the mixture of different Gaussian predictions (each given by a different combination of faultiness) as a single Gaussian. We prefer this approximate marginalisation over faultiness to heuristics that would designate all observations as either faulty or not—we acknowledge our uncertainty about faultiness.

More formally, imagine that we have partitioned our observations  $\mathcal{D}_{a,b}$  into a set of old observations  $\mathcal{D}_a = (\mathbf{x}_a, \mathbf{y}_a)$  and a set of new observations  $\mathcal{D}_b = (\mathbf{x}_b, \mathbf{y}_b)$ . Define  $\sigma_a$  to be the (unknown) vector of all noise variances at observations  $\mathbf{y}_a$ , and define  $\sigma_b$  similarly. Because we have to sum over all possible values for these vectors, we will index the possible values of  $\sigma_a$  by  $i$  (each given by a different combination of faultiness over  $\mathcal{D}_a$ ) and the values of  $\sigma_b$  similarly by  $j$ . We now define the covariances  $V_a^i = K_{a,a} + \text{diag } \sigma_a^i$ ,  $V_b^j = K_{b,b} + \text{diag } \sigma_b^j$  and  $V_{a,b}^{i,j} = K_{\{a,b\},\{a,b\}} + \text{diag}\{\sigma_a^i, \sigma_b^j\}$ , where  $\text{diag } \sigma$  is the diagonal matrix with diagonal  $\sigma$ .

To initialise our algorithm, imagine that  $a$  identifies a small set of data, such that we can readily compute the likelihood of our hyperparameters

$$p(\mathbf{y}_a) = \sum_i p(\mathbf{y}_a | \sigma_a^i) p(\sigma_a^i) = \sum_i \mathcal{N}(\mathbf{y}_a; 0, V_a^i) p(\sigma_a^i) \quad (3)$$

and hence the hyperparameter posterior,  $p(\sigma_a | \mathbf{y}_a)$ . This distribution specifies the probability of our observations  $\mathcal{D}_a$  being faulty; for a single observation  $\mathcal{D}_a$ ,

$p(\text{fault}(\mathcal{D}_a) | \mathbf{y}_a) = p(\sigma_a = \sigma_f | \mathbf{y}_a)$ . If we were to perform predictions for some  $f_\star$  using  $\mathcal{D}_a$  alone, we would need to evaluate

$$\begin{aligned} p(f_\star | \mathbf{y}_a) &= \sum_i p(\sigma_a^i | \mathbf{y}_a) p(f_\star | \mathbf{y}_a, \sigma_a^i) \\ &= \sum_i p(\sigma_a^i | \mathbf{y}_a) \mathcal{N}(f_\star; m(f_\star | \mathbf{y}_a, \sigma_a^i), C(f_\star | \mathbf{y}_a, \sigma_a^i)), \end{aligned}$$

the weighted sum of Gaussian predictions made using the different possible values for  $\sigma_a$ . We now use approximation 2. It is our hope that our predictions for  $f_\star$  are not so sensitive to the noise in our observations that all the Gaussians in this sum become dramatically different. In any case, the quality of this approximation will improve over time—if  $f_\star$  is far removed from our old data  $\mathcal{D}_a$ , then our predictions really will not be very sensitive to  $\sigma_a$ . So, we take

$$\begin{aligned} p(f_\star | \mathbf{y}_a) &\simeq \mathcal{N}(f_\star; K_{\star,a} \tilde{V}_a^{-1} \mathbf{y}_a, \\ &\quad K_{\star,\star} - K_{\star,a}(\tilde{V}_a^{-1} - \tilde{W}_a^{-1})K_{a,\star} - (K_{\star,a} \tilde{V}_a^{-1} \mathbf{y}_a)^2), \end{aligned}$$

where<sup>1</sup>

$$\begin{aligned} \tilde{V}_a^{-1} &= \sum_i p(\sigma_a^i | \mathbf{y}_a) (V_a^i)^{-1}, \\ \tilde{W}_a^{-1} &= \sum_i p(\sigma_a^i | \mathbf{y}_a) (V_a^i)^{-1} \mathbf{y}_a \mathbf{y}_a^\top (V_a^i)^{-1}. \end{aligned} \quad (5)$$

With these calculations performed, imagine receiving further data  $\mathcal{D}_b$ . To progress, we make approximation 3; we assume that faults will not persist longer than  $|\mathcal{D}_b|$ . To be precise, we assume

$$p(\sigma_{a,b}^{i,j} | \mathbf{y}_{a,b}) \simeq p(\mathbf{y}_a) p(\sigma_a^i | \mathbf{y}_a) p(\sigma_b^j | \mathbf{y}_b) p(\mathbf{y}_b | \sigma_{a,b}^{i,j}, \mathbf{y}_a) \quad (6)$$

Our predictions are now

$$\begin{aligned} p(f_\star | \mathbf{y}_{a,b}) &\simeq \sum_j p(\sigma_b^j | \mathbf{y}_{a,b}) \sum_i p(\sigma_a^i | \mathbf{y}_a) \\ &\quad \mathcal{N}(f_\star; m(f_\star | \mathbf{y}_{a,b}, \sigma_{a,b}^{i,j}), C(f_\star | \mathbf{y}_{a,b}, \sigma_{a,b}^{i,j})). \end{aligned} \quad (7)$$

Before trying to manage these sums, we will determine  $p(\sigma_b | \mathbf{y}_{a,b})$ . As before, this distribution gives

<sup>1</sup> Note that for  $\tilde{W}_a$ , explicitly computing (unstable) matrix inverses can be avoided by solving the appropriate linear equations using Cholesky factors. For  $\tilde{V}_a$ , we can rewrite  $(A^{-1} + B^{-1})^{-1} = A(A + B)^{-1}B$ . If  $i \in \{0, 1\}$  (as it would be if  $a$  identified a single observation which could be either faulty or not),

$$\tilde{V}_a = V_a^0 (p(\sigma_a^1 | \mathbf{y}_a) V_a^0 + p(\sigma_a^0 | \mathbf{y}_a) V_a^1)^{-1} V_a^1. \quad (4)$$

If  $i$  takes more than two values, we can simply iterate using the same technique. We can then use the Cholesky factor of  $\tilde{V}_a$  to compute our required equations.

us the probability of the observations  $\mathcal{D}_b$  being faulty. For example, if we have only a single observation  $\mathcal{D}_b$ ,  $p(\text{fault}(\mathcal{D}_b) | \mathbf{y}_{a,b}) = p(\sigma_b = \sigma_f | \mathbf{y}_{a,b})$ . We define

$$\begin{aligned} \tilde{m}(\mathbf{y}_b | \mathbf{y}_a) &= K_{b,a} \tilde{V}_a^{-1} \mathbf{y}_a \\ \tilde{C}(\mathbf{y}_b | \mathbf{y}_a, \sigma_b) &= V_b - K_{b,a}(\tilde{V}_a^{-1} - \tilde{W}_a^{-1})K_{a,b} \\ &\quad - \tilde{m}(\mathbf{y}_b | \mathbf{y}_{a,b})^2, \end{aligned}$$

where both  $\tilde{V}_a$  (or its Cholesky factor) and  $\tilde{W}_a^{-1}$  were computed previously. By using approximations 2 and 3,

$$\begin{aligned} p(\sigma_b | \mathbf{y}_{a,b}) &= \frac{\sum_i p(\mathbf{y}_b | \mathbf{y}_a, \sigma_{a,b}^i) p(\mathbf{y}_a, \sigma_{a,b}^i)}{p(\mathbf{y}_{a,b})} \\ &\simeq \frac{\mathcal{N}(\mathbf{y}_b; \tilde{m}(\mathbf{y}_b | \mathbf{y}_a), \tilde{C}(\mathbf{y}_b | \mathbf{y}_a, \sigma_b)) p(\sigma_b)}{p(\mathbf{y}_b | \mathbf{y}_a)}, \end{aligned}$$

where we have

$$\begin{aligned} p(\mathbf{y}_b | \mathbf{y}_a) &= \sum_i \sum_j p(\mathbf{y}_b | \mathbf{y}_a, \sigma_{a,b}^i) p(\sigma_{a,b}^{i,j} | \mathbf{y}_a) \\ &\simeq \sum_j \mathcal{N}(\mathbf{y}_b; \tilde{m}(\mathbf{y}_b | \mathbf{y}_a), \tilde{C}(\mathbf{y}_b | \mathbf{y}_a, \sigma_b^j)) p(\sigma_b^j). \end{aligned}$$

Note that the product of  $p(\mathbf{y}_b | \mathbf{y}_a)$  and  $p(\mathbf{y}_a)$  (previously computed in (3)) gives the likelihood of our hyperparameters. Now, returning to (7), we will once again use approximation 2. We aim to reuse our previously evaluated sums over  $i$  to resolve future sums over  $i$ . As we gain more data, the faultiness of old data becomes less important. We arrive at

$$\begin{aligned} p(f_\star | \mathbf{y}_{a,b}) &\simeq \mathcal{N}(f_\star; K_{\star,\{a,b\}} \tilde{V}_{a,b}^{-1} \mathbf{y}_{a,b}, \\ &\quad K_{\star,\star} - K_{\star,a}(\tilde{V}_{a,b}^{-1} - \tilde{W}_{a,b}^{-1})K_{a,\star} - (K_{\star,\{a,b\}} \tilde{V}_{a,b}^{-1} \mathbf{y}_{a,b})^2), \end{aligned}$$

where we have

$$\begin{aligned} \tilde{V}_{a,b}^{-1} &= \sum_j p(\sigma_b^j | \mathbf{y}_{a,b}) \sum_i p(\sigma_a^i | \mathbf{y}_a) (V_{a,b}^{i,j})^{-1} \\ \tilde{W}_{a,b}^{-1} &= \sum_j p(\sigma_b^j | \mathbf{y}_{a,b}) \sum_i p(\sigma_a^i | \mathbf{y}_a) \\ &\quad (V_{a,b}^{i,j})^{-1} \mathbf{y}_{a,b} \mathbf{y}_{a,b}^\top (V_{a,b}^{i,j})^{-1}. \end{aligned}$$

Now, using the inversion by partitioning formula (PRESS, *et al.*, 1992, Section 2.7),

$$\begin{aligned} (V_{a,b}^{i,j})^{-1} &= \\ &\begin{bmatrix} S_{a,b}^{i,j} & -S_{a,b}^{i,j} K_{a,b} (V_b^j)^{-1} \\ -(V_b^j)^{-1} K_{b,a} S_{a,b}^{i,j} & (V_b^j)^{-1} + (V_b^j)^{-1} K_{b,a} S_{a,b}^{i,j} K_{a,b} (V_b^j)^{-1} \end{bmatrix} \end{aligned}$$

where  $S_{a,b}^{i,j} = (V_a^i - K_{a,b} (V_b^j)^{-1} K_{b,a})^{-1}$ . Note that  $(V_{a,b}^{i,j})^{-1}$  is affine in  $S_{a,b}^{i,j}$ , so that when  $V_a \gg K_{a,b} V_b^{-1} K_{b,a}$ ,  $(V_{a,b}^{i,j})^{-1}$  is effectively affine in  $(V_a^i)^{-1}$ . This is true if given  $\mathcal{D}_b$ , it is impossible to accurately

predict  $\mathcal{D}_a$ . This might be the case if  $\mathcal{D}_a$  represents a lot of information relative to  $\mathcal{D}_b$  (if, for example,  $\mathcal{D}_a$  is our entire history of observations where  $\mathcal{D}_b$  is simply the most recent observation), or if  $\mathcal{D}_b$  and  $\mathcal{D}_a$  are simply not particularly well correlated. On this basis, we make approximation 4. Additionally noting that  $\sum_i p(\sigma_a^i | \mathbf{y}_a) = 1$ , we have <sup>2</sup>

$$\begin{aligned}\tilde{V}_{a,b}^{-1} &\simeq \sum_j p(\sigma_b^j | \mathbf{y}_{a,b}) \begin{bmatrix} \tilde{V}_a & K_{a,b} \\ K_{b,a} & V_b^j \end{bmatrix}^{-1}, \\ \tilde{V}_{a,b} &\simeq \begin{bmatrix} \tilde{V}_a & K_{a,b} \\ K_{b,a} & \tilde{V}_{b|a} + K_{b,a} \tilde{V}_a^{-1} K_{a,b} \end{bmatrix} \\ \tilde{V}_{b|a}^{-1} &= \sum_j p(\sigma_b^j | \mathbf{y}_{a,b}) (V_b^j - K_{b,a} \tilde{V}_a^{-1} K_{a,b})^{-1}.\end{aligned}$$

Note that the lower right hand element of  $\tilde{V}_{a,b}$  defines the noise variance to be associated with observations  $\mathcal{D}_b$ . In effect, we represent each observation as having a known variance lying between  $\sigma_n^2$  and  $\sigma_f^2$ . The more likely an observation's faultiness, the closer its assigned variance will be to the (large) fault variance and the less relevant it will become for inference about the latent process. This approximate observation is then used for future predictions; we need never consider the full sum over all observations.

We now turn to  $\tilde{W}_{a,b}^{-1}$ . Unfortunately, even if  $V_a \gg K_{a,b} V_b^{-1} K_{b,a}$ ,  $\tilde{W}_{a,b}^{-1}$  is quadratic in  $(V_a^i)^{-1}$ . We will nonetheless again make approximation 4 and assume that  $\tilde{W}_{a,b}^{-1}$  is affine in  $(V_a^i)^{-1}$ . The quality of our approximation for  $\tilde{W}_{a,b}^{-1}$  is much less critical than for  $\tilde{V}_{a,b}^{-1}$ , because the former only influences the variance of our predictions for the current predictant; any flaws in that approximation will not be propagated forward. Further, of course, if one probability dominates,  $p(\sigma_a^i | \mathbf{y}_a) \gg p(\sigma_a^{i'} | \mathbf{y}_a), \forall i' \neq i$ , then the approximation is valid. With this,

$$\begin{aligned}\tilde{W}_{a,b}^{-1} &\simeq \sum_j p(\sigma_b^j | \mathbf{y}_{a,b}) \\ &\quad \begin{bmatrix} \tilde{V}_a & K_{a,b} \\ K_{b,a} & V_b^j \end{bmatrix}^{-1} \mathbf{y}_{a,b} \mathbf{y}_{a,b}^\top \begin{bmatrix} \tilde{V}_a & K_{a,b} \\ K_{b,a} & V_b^j \end{bmatrix}^{-1}.\end{aligned}$$

If we now receive further data  $\mathcal{D}_c$ , our existing data is simply treated as old data ( $a \leftarrow \{a, b\}$ ,  $b \leftarrow c$ ), and another iteration of our algorithm performed.

An outline of our approach is depicted in Algorithm ??.

<sup>2</sup>  $\tilde{V}_{b|a}^{-1}$  can be computed using the same trick as in (4) if  $b$  identifies a single observation and  $j \in \{0, 1\}$ . The Cholesky factor of  $\tilde{V}_{a,b}$  required to solve the linear equations for our predictions can be efficiently determined (OSBORNE, 2010, Appendix B) using the previously evaluated Cholesky factor of  $\tilde{V}_a$ .

### 3.1 Discussion

We return to the management of our hyperparameters  $\theta$ . Unfortunately, analytically marginalising  $\theta$  is impossible. Most of the hyperparameters of our model can be set by optimising their likelihood on a large training set, giving a likelihood close to a delta function. This is not true of the hyperparameters  $\sigma_n^2$  and  $\sigma_f^2$ , due to exactly the same problematic sums discussed earlier. Instead, we marginalise these hyperparameters online using Bayesian Monte Carlo (OSBORNE, 2010, Chapter 7), taking a fixed set of samples in their values and using the hyperparameter likelihoods  $p(\mathbf{y}_{a,b})$  to construct weights over them. Essentially, we proceed as described above independently in parallel for each sample, and combine the predictions from each in a final weighted mixture for prediction. Note that we can use a similar procedure (GARNETT, *et al.*, 2010) to determine the full posterior distributions of  $\sigma_n^2$  and  $\sigma_f^2$ , if desired. It would be desirable to use non-fixed samples, but, unfortunately, this would require reconstructing our full covariance matrix from scratch each time a sample is moved.

The proposal above can be extended in several ways. First, we may wish to sum over more than one fault variance, which could be useful if observations are prone to faultiness in more than one mode. If, instead of summing over a small number of known variances, we wished to marginalise with respect to a density over noise variance, we can simply replace the sums over  $i$  and  $j$  with appropriate integrals. Obviously this will only be analytically possible if the posteriors for  $\sigma_a^2$  take appropriate, simple forms. In such a way our algorithm might tackle the general problem of heteroskedasticity.

Our proposed algorithm steps through our data one at a time, so that  $\mathcal{D}_b$  always contains only a single observation. With approximation 3, this means that our algorithm is not expecting faults to last more than a single observation. The results that follow, however, will show that we can nonetheless manage sustained faults. It would also be possible to step in larger chunks, evaluating larger sums. Although more computationally demanding, this might be expected to improve results. It would also allow us to consider non-diagonal noise contributions.

We have so far not specified our prior for faultiness (as expressed by  $p(\sigma_a^2)$  and  $p(\sigma_b^2)$ ). Within this paper, we consider exclusively a time-independent probability of faultiness, but our framework does not necessarily require this to be so. In some contexts it might be useful to perform inference about the fault contribution, rather than the signal of interest. To do so, we merely switch the roles of the fault and non-fault contributions. Note that, using our full posteriors for faultiness, we

can also trivially use Bayesian decision theory to make hard decisions as required.

## 4 Results

We test the effectiveness of the fault bucket algorithm on several time-series that are indicative of problems found in environmental monitoring. In particular, we test on water-level readings; such data are often characterised by complex dynamics and will therefore provide a good indicator of our algorithm’s performance on real-world tasks. We aim to improve upon the simple, human-supervised approaches to fault detection used in this field (WAGNER and US GEOLOGICAL SURVEY, 2006). For a quantitative assessment, we used two semi-synthetic datasets where a typical fault has been injected into clean sensor data. We then analyzed qualitative performance on two real data sets with actual faults. All measurements (other than for pH) are given in meters, with samples spaced in increments of approximately 30 minutes.

Our first synthetic example, a bias fault, concerns a simple sensor error where measurements are temporarily adjusted by a constant offset, but otherwise remain accurate. This could happen if the sensor undergoes physical trauma which results in a loss of calibration. The next dataset contains a synthetic anomaly where the water level rises quickly, but smoothly, before returning back to normal. This would be indicative of a genuine environmental event such as a flash flood. Both synthetic datasets represented sustained faults, of length 335 and 161 observations respectively. Our first real dataset deals with pH measurements from the United States Geological Survey, a common indicator of water quality. In this series, a clear sensor fault can be seen in which the observations undergo a sudden, sustained decrease in value. The next (real) dataset contains a fault type called “painting,” which is an error that occurs when ice builds on a sensor, obscuring some of the readings. It is characterised by frequent sensor spikes interlaced with the original, and still accurate, signal.

We implemented the algorithm described in Section 3 in MATLAB to address the task of 1-step-lookahead time-series prediction. A sliding window of size 100 was used to predict the value of the next observation. Each dataset was recentered so that a zero prior mean function was appropriate, and the functions were all modeled using a Matérn covariance with parameter  $\nu = 5/2$  (RASMUSSEN and WILLIAMS, 2006). The hyperparameters for this covariance, including the normal observation noise variance  $\sigma_n^2$ , were learned using training data similar to but disjoint from the test datasets. The unknown fault noise variance  $\sigma_f^2$  was marginalised

using Bayesian Monte Carlo, with a parsimonious 7 samples used. The prior probability of an observation being faulty was set to a constant value of 1% throughout.

We tested against a number of different methods in order to establish the efficacy of the fault bucket algorithm. All GP-based approaches used the same hyperparameters employed by our algorithm. The training set used to learn those hyperparameters was also supplied to other methods for their respective model learning phases. Several methods identify a new observation  $y$  as a fault if

$$|y - m(y|\mathbf{y})| > 3\sigma_T, \quad (8)$$

where  $m(y|\mathbf{y})$  is the method’s *a priori* prediction for  $y$ , and  $\sigma_T$  is the noise standard deviation on the faultless training set. Of course, methods using (8) or similar can not provide the posterior probability of a point’s faultiness, as our algorithm can. Methods tested include:

**XGP:** A GP in which we exhaustively search over the faultiness of the last 10 points, and approximate the noise variance of all previous points in the window as having the value  $\sigma_f^2 p(\text{fault}|\mathbf{y}) + \sigma_n^2 p(\neg \text{fault}|\mathbf{y})$ , fixed at the time the point was observed (when data  $\mathcal{D}$  was available). Clearly, this method is very much more computationally expensive than the fault bucket algorithm (roughly  $2^9$  times more), but offers a useful way to quantify the influence of approximations 2–4.

**EPGP:** A GP with our observation likelihood (2) and the posterior determined by expectation propagation MINKA (2001). Note that expectation propagation will not give us a posterior probability of faultiness, for the purposes of explicitly identifying faults, we use (8).

**TGP:** A GP in which a point was flagged as a fault using (8); if faulty, a point was treated as having noise variance  $\sigma_f^2$ .

**STGP:** A GP with student-t likelihood; the posterior was determined using a Laplace approximation as per VANHATALO, *et al.* (2009). To identify faults, (8) was used.

**MLH:** The most likely heteroscedastic GP (KERSTING, *et al.*, 2007).

**EKF:** An autoregressive neural net trained with the extended Kalman filter to capture nonstationarity. Again, (8) was used to identify and discard faulty data.

Note that for EPGP, STGP and MLH, we perform retrospective prediction (so that all data is available to make predictions about even the first predictant), as these methods are usually used. Clearly this allows these approaches an unfair predictive advantage relative to

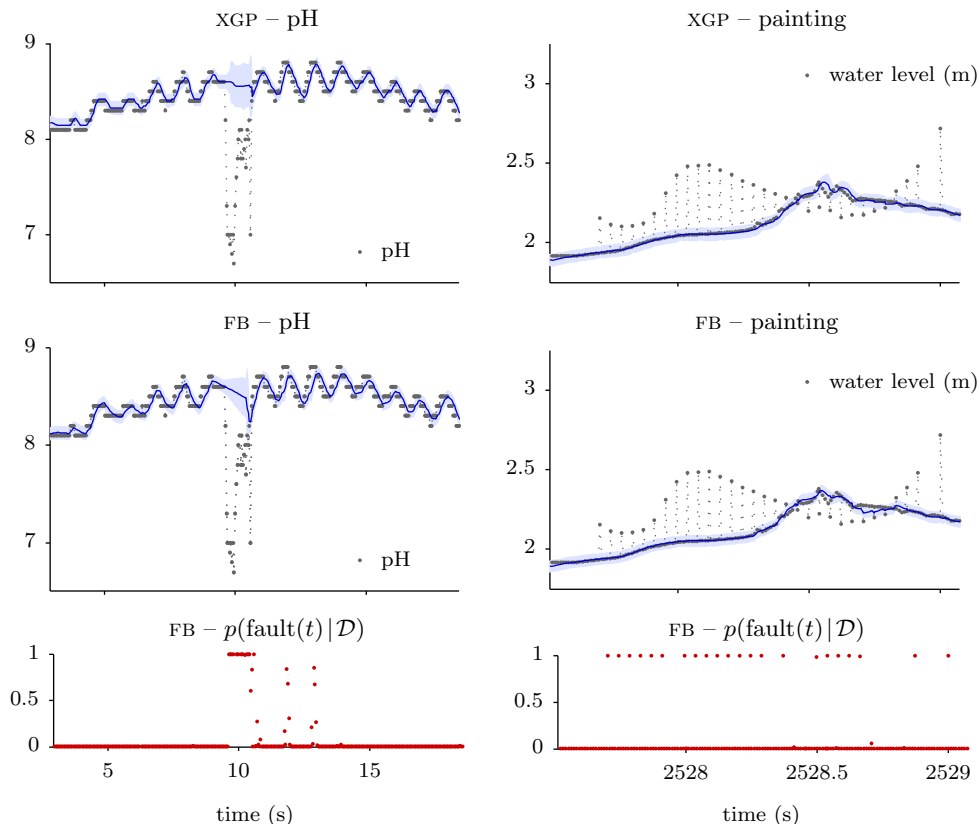


Figure 2: Mean and  $\pm 3\sigma$  standard-deviation bounds for the predictions of the exhaustive (XGP) and fault-bucket (FB) algorithms on the pH and painting datasets. For the fault-bucket algorithm, the posterior faultiness of each observation is also shown underneath the predictions. Note that each column shares the same  $x$ -axis.

sequential methods. Note also that the multiple passes over the data effected by approximation schemes such as expectation propagation cannot be readily applied to the sequential problem without requiring a great deal of expensive computation. For  $N$  observations, the computational cost of expectation propagation is  $\mathcal{O}(N^3)$ , with a large constant of proportionality, rendering it impractical for real-time problems. With efficient Cholesky factor updates, our scaling is  $\mathcal{O}(N^2)$ . For this reason, we did not consider even more demanding, non-sequential, expectation propagation approaches such as the twinned GP (NAISH-GUZMAN and HOLDEN, 2008).

Figures 2 show the performance of the fault-bucket algorithm and the exhaustive alternative on the two real datasets. The fault-bucket algorithm did an excellent job of identifying faults when they occur, and made excellent predictions, even for the sustained fault in the pH dataset.

Table 1 displays quantitative measures of performance for the various algorithms on the synthetic datasets. In addition to superior predictive performance, our detection rates for the faulty points are generally ex-

cellent. The results reveal that approximations 2–4 do not result in significant loss of performance relative to exhaustive search. Our naïve approach to faults may, of course, suffer relative to better-informed models, but its probabilistic estimates provide a human operator with an indication as to whether more sophisticated analysis is necessary.

## 5 Conclusion

We have proposed a novel algorithm, the “fault bucket,” for managing time-series data corrupted by faults of type unknown ahead of time. Our chief contribution is a sequential algorithm for marginalising the faultiness of observations in a GP framework, allowing for fast, effective prediction in the presence of unknown faults. Unlike most robust regression approaches (such as those using student-t likelihoods), we can also compute the posterior probability of faultiness. This capacity is crucial to its utility for the domain, serving as a means of alarming a human operator to the possible need for corrective action.

As to future work, addressing multivariate signals is of

Table 1: Quantitative comparison of different algorithms on the synthetic datasets. For each dataset, we show the mean squared error (MSE), the log likelihood of the true data ( $\log p(\mathbf{y}|\mathbf{x})$ ), and the true-positive and false-positive rates of detection for faulty points (TPR and FPR), respectively, with all methods permitted a ‘burn-in’ period of 50 points. The best value for each set of results is highlighted in bold.

Method	Bias dataset				“Flash-flood” dataset			
	MSE	$\log p(\mathbf{y} \mathbf{x})$	TPR	FPR	MSE	$\log p(\mathbf{y} \mathbf{x})$	TPR	FPR
FB	<b>0.024</b>	334	<b>0.997</b>	0.031	0.069	$-5.77 \times 10^3$	<b>0.829</b>	0.016
XGP	0.037	<b>439</b>	0.982	<b>0.022</b>	<b>0.042</b>	<b><math>-1.52 \times 10^3</math></b>	0.805	<b>0.012</b>
EPGP	0.880	$-3.39 \times 10^3$	0.000	0.000	2.179	$-1.43 \times 10^4$	0.000	0.000
TGP	0.033	278	<b>0.997</b>	0.031	0.075	$-8.29 \times 10^3$	<b>0.829</b>	0.083
STGP	0.939	$-2.49 \times 10^5$	0.027	0.025	0.375	$-1.08 \times 10^5$	0.689	0.005
MLH	0.940	$-5.43 \times 10^7$	0.065	0.031	2.369	$-2.27 \times 10^7$	0.045	0.262
EKF	0.060	$-1.26 \times 10^4$	0.551	0.258	0.613	$-1.81 \times 10^4$	0.169	0.768

great interest. Unfortunately, this extension is not trivial and would itself require additional approximations.

## References

- CHANDOLA, V., BANERJEE, A. and KUMAR, V. (2009 July). Anomaly detection: A survey. *ACM Comput. Surv.*, vol. 41, pp. 15:1–15:58.
- DE FREITAS, N., MACLEOD, I.M. and MALTZ, J.S. (1996). Neural networks for pneumatic actuator fault detection. *Transactions of the South African Institute of Electrical Engineers*, vol. 90, pp. 28–34.
- DERESZYNSKI, E. and DIETTERICH, T.G. (2011). Spatiotemporal models for anomaly detection in dynamic environmental monitoring campaigns. *ACM Transactions on Sensor Networks*.
- DING, S.X. (2008). *Model-based Fault Diagnosis Techniques: Design Schemes, Algorithms, and Tools*. 1st edn. Springer.
- ECIOLAZA, L., ALKAROURI, M., LAWRENCE, N.D., KADIRKAMANATHAN, V. and FLEMING, P.J. (2007). Gaussian Process Latent Variable Models for Fault Detection. In: *IEEE Symposium on Computational Intelligence and Data Mining*, pp. 287–292.
- GARNETT, R., OSBORNE, M.A., REECE, S., ROGERS, A. and ROBERTS, S.J. (2010). Sequential Bayesian Prediction in the Presence of Changepoints and Faults. *The Computer Journal*, vol. 53.
- ISERMANN, R. (2005). Model-based fault-detection and diagnosis – status and applications. *Annual Reviews in Control*, vol. 29, no. 1, pp. 71–85.
- KERSTING, K., PLAGEMANN, C., PFAFF, P. and BURGARD, W. (2007). Most likely heteroscedastic Gaussian process regression. In: *Proceedings of the 24th international conference on Machine learning*, pp. 393–400. ACM.
- KHAN, S. and MADDEN, M. (2010). A survey of recent trends in one class classification. In: COYLE, L. and FREYNE, J. (eds.), *Artificial Intelligence and Cognitive Science*, vol. 6206 of *Lecture Notes in Computer Science*, pp. 188–197. Springer Berlin / Heidelberg.
- MARKOU, M. and SINGH, S. (2003). Novelty detection: a review – Part 1: Statistical approaches. *Signal Processing*, vol. 83, no. 12, pp. 2481–2497.
- MINKA, T. (2001). Expectation propagation for approximate bayesian inference. In: *Uncertainty in Artificial Intelligence*, vol. 17, pp. 362–369. Citeseer.
- NAISH-GUZMAN, A. and HOLDEN, S. (2008). Robust regression with twinned Gaussian processes. *Advances in Neural Information Processing Systems*, vol. 20, pp. 1065–1072.
- OSBORNE, M.A. (2010). *Bayesian Gaussian Processes for Sequential Prediction, Optimisation and Quadrature*. Ph.D. thesis, University of Oxford. Available at [www.robots.ox.ac.uk/~mosb/full\\_thesis.pdf](http://www.robots.ox.ac.uk/~mosb/full_thesis.pdf).
- PRESS, W.H., TEUKOLSKY, S.A., VETTERLING, W.T. and FLANNERY, B.P. (1992). *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge University Press, New York, NY, USA.
- RASMUSSEN, C.E. and WILLIAMS, C.K.I. (2006). *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA.
- VANHATALO, J., JYLÄNKI, P. and VEHTARI, A. (2009). Gaussian process regression with student-t likelihood. In: *Advances in Neural Information Processing Systems*, vol. 22, pp. 1910–1918.
- WAGNER, R.J. and US GEOLOGICAL SURVEY (2006). *Guidelines and standard procedures for continuous water-quality monitors: Station operation, record computation, and data reporting*. US Department of the Interior, US Geological Survey.