

Anomaly Detection in Star Light Curves using Hierarchical Gaussian Processes

Haoyan Chen¹ and Tom Diethe^{1,2} * and Niall Twomey¹ and Peter Flach¹

1- University of Bristol - Intelligent Systems Group, Bristol - UK

2- Amazon Research Cambridge, Cambridge - UK

Abstract. Here we examine astronomical time-series called *light-curve* data, which represent the brightness of celestial objects over a period of time. We focus specifically on the task of finding anomalies in three sets of light-curves of periodic variable stars. We employ a hierarchical Gaussian process to create a general and stable model of time series for anomaly detection, and apply this approach to the light curve problem. Hierarchical Gaussian processes require only a few additional parameters than Gaussian processes and incur negligible additional computational complexity. Additionally, the additional parameters are objectively optimised in a principled probabilistic framework. Experimentally, our approach outperforms several baselines and highlights several anomalous light curves in the datasets investigated.

1 Introduction

Anomaly detection refers to the problem of finding patterns in data that do not conform to expected behaviour and we often call these non-conforming patterns anomalies, outliers, aberrations, exceptions, etc. in different applications [1]. One of the major tasks in astronomy is to detect when aberrant phenomena are encountered from historical observations [2], and it is almost impossible to find the anomalous objects through manual inspection due to the scale of the data. In this work we apply an unsupervised anomaly detection method to an astronomical time-series data. Fig. 1 shows a typical light-curve from the Optical Gravitational Lensing Experiment (OGLE) [3] for periodic variable stars (CEPH, EB and RRL) after data pre-processing [4]. Cepheid (CEPH), Eclipsing Binaries (EB) and RR Lyrae (RRL) are common types of periodic variable stars, and the details of these three stars can be found in [5]. The three light-curves in fig. 1 are typical, but the dataset may contain an unknown number of additional outliers, and this work introduces a robust approach to detecting them.

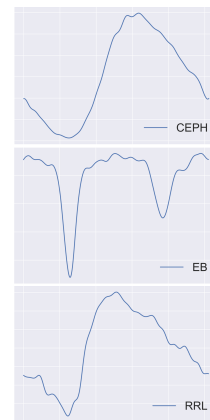


Fig. 1: Example light-curves

Probabilistic models making use of Gaussian Processes (GPs) have become a standard approach to solving a variety of machine learning problems, due to their flexibility, quantification of uncertainty, and calibrated probabilistic outputs. Hierarchical GPs (HGPs) [6] are a recent method that provides a statistical

*Work done prior to joining Amazon.

model for tasks involving multiple related time-series. HGPs were originally proposed for the analysis of gene expression data, where the multiple time series are seen as noisy realisations of a common underlying driver function. Our central hypothesis is that this intuition is also a fundamental property of light curves, where the underlying function represents the periodic physical variation of the celestial bodies.

2 Related Work

Statistical models for anomaly detection assume that anomalies usually occur in low probability regions [1]. In this work we focus on *non-parametric* methods that do not assume a fixed parametric form for the underlying stochastic model.

Histogram-based models count the frequency of values in the training data, and declare a test instance as an anomaly if it cannot be assigned to any bin of that histogram [1]. This approach is particularly sensitive to the bin size which can produce unstable models. Nearest Neighbour-based techniques can also be applied to anomaly detection, where the key assumption is that anomalies are far away from their closest neighbour [1], for example by computing a *local outlier factor* [7]. Periodic Curve Anomaly Detection (PCAD) [2] is a modified k-means algorithm for time-series data that uses cross-correlation as distance metric and updates phase information at each iteration. PCAD will form a baseline for this work as it was developed in the context of light curves.

Kernel-based methods have been used to estimate probability density functions for “normal” instances [8]. GPs harness the power of kernel methods, with the advantage of capturing the uncertainty in the data. GPs have been combined with Active Learning [9] and Extreme Value Theory [10] to detect abnormal behaviours in a maritime dataset. However this method is limited to finding anomalies *within* a time-series, rather than anomalous time-series within a collection, as is the case here.

3 Methods

A GP is a distribution over functions, and is specified by a mean function $m(t)$ and its covariance function $k(t, t')$, compactly $f(t) \sim GP(m(t), k(t, t'))$. When the function f is evaluated at points t , the marginal distribution is a multivariate normal [11]. The covariance (kernel) function k takes two inputs and reveals how similar they are. When performing Bayesian inference, we have that the set of known function values of the training examples \mathbf{f} , and the set of known function values corresponding to a set of testing examples \mathbf{y} , are jointly Gaussian:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_f \\ \mathbf{m}_y \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{fy} \\ \mathbf{K}_{fy}^\top & \mathbf{K}_{yy} \end{bmatrix} \right), \quad (1)$$

where $\mathbf{m}_f, \mathbf{m}_y$ are the means for training and testing examples respectively, and similarly $\mathbf{K}_{ff}, \mathbf{K}_{yy}, \mathbf{K}_{fy}$ denote the train, test, train-test covariances. The

conditional distribution of \mathbf{f} given by \mathbf{y} is then:

$$\mathbf{f}|\mathbf{y} \sim \mathcal{N}\left(\mathbf{K}_{fy}\mathbf{K}_{yy}^{-1}(\mathbf{y} - \mathbf{m}_y) + \mathbf{m}_f, \mathbf{K}_{ff} - \mathbf{K}_{fy}\mathbf{K}_{yy}^{-1}\mathbf{K}_{fy}^T\right). \quad (2)$$

This is the posterior distribution for a specific set of testing examples. For a training set \mathcal{D} , the posterior process can be described by [11]:

$$\begin{aligned} y|\mathcal{D} &\sim GP(m_{\mathcal{D}}, k_{\mathcal{D}}), \\ m_{\mathcal{D}}(t) &= m(t) + \mathbf{k}_{Tt}^T \mathbf{K}^{-1}(y - m) \\ k_{\mathcal{D}}(t, t') &= k(t, t') - \mathbf{k}_{Tt}^T \mathbf{K}^{-1} \mathbf{k}_{Tt'}, \end{aligned} \quad (3)$$

where $\mathbf{k}_{Tt} = k(T, t)$ is the covariance vector between every training example and t . Notice that the posterior variance $k_{\mathcal{D}}(t, t')$ is always smaller than the prior variance $k(t, t')$ because $\mathbf{k}_{Tt}^T \mathbf{K}^{-1} \mathbf{k}_{Tt'}$ is positive. This process means that we are reducing the degree of uncertainty by using training examples.

The choice of covariance function is based on an understanding of the domain. In the case of light-curves, so we need a kernel that can express both smooth variation and small fluctuations (see fig. 1). An appropriate choice is the Matern 3/2 kernel, which is twice differentiable, and is given by:

$$k_{\text{Matern32}}(x, x') = \sigma_v^2 \left(1 + \frac{\sqrt{3} \|x - x'\|_2}{\ell}\right) \exp\left(-\frac{\sqrt{3} \|x - x'\|_2}{\ell}\right), \quad (4)$$

where $\|\cdot\|_2$ denotes the ℓ_2 -norm, and ℓ, σ are the kernel hyperparameters. We will use type-II maximum likelihood to optimise the hyperparameters [11].

3.1 A hierarchy across time series

In this work, we adopt the following notation: \mathbf{l}_{no} denotes the vector of light-curve measurements of the \mathbf{n}^{th} star in the \mathbf{o}^{th} observation, and \mathbf{t}_{no} is the sequence of time of that measurement. We can use the following expression to indicate the data for the \mathbf{n}^{th} star: $\mathbf{L}_n = \{\mathbf{l}_{no}\}_{o=1}^{N_n}$, $\mathbf{T}_n = \{\mathbf{t}_{no}\}_{o=1}^{N_n}$.

$$\begin{aligned} f_n(t) &\sim GP(0, k_f(t, t')), \\ s_{no}(t) &\sim GP(f_n(t), k_s(t, t')). \end{aligned} \quad (5)$$

Thus, the probability of a set of N_n observations $\mathbf{L}_n = \{\mathbf{l}_{no}\}_{o=1}^{N_n}$, measured at time $\mathbf{T}_n = \{\mathbf{t}_{no}\}_{o=1}^{N_n}$, a likelihood expression is denoted by:

$$p(\mathbf{L}_n|\mathbf{T}_n, \theta) = \mathcal{N}(\bar{\mathbf{I}}_n|\mathbf{0}, \mathbf{K}_n), \quad (6)$$

where $\bar{\mathbf{I}}_n$ is used to indicate the concatenation of \mathbf{L}_n , $\bar{\mathbf{I}}_n = [\mathbf{l}_{n,1}, \mathbf{l}_{n,2} \cdots, \mathbf{l}_{n,N_n}]$, and \mathbf{K}_n denotes the augmented covariance matrix, which is defined as follows:

$$\mathbf{K}_n(o, o') = \begin{cases} \mathbf{K}_f(t_{no}, t_{no'}) + \mathbf{K}_s(t_{no}, t_{no'}) + \sigma_n^2 I & \text{if } o = o' \\ \mathbf{K}_f(t_{no}, t_{no'}) & \text{otherwise.} \end{cases} \quad (7)$$

Owing to space limitations, we refer the reader to [6] for intuition on this.

Since the light curves exhibit both low frequency oscillations and high frequency noise, we have selected the Matern kernel for k_f and k_s since they are capable of capturing such dynamics. Fig. 2 illustrates how the hierarchical GP model we build performs on the three types of astronomical stars. The left-most column defines the underlying latent functions inferred by the HGPAD model, and the remaining subplots show example light curves from the OGLE dataset.

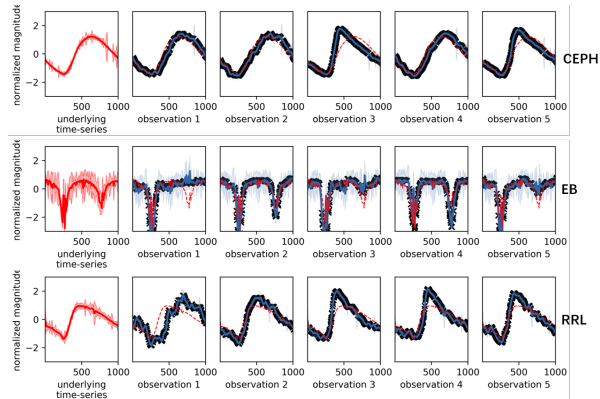


Fig. 2: HGP on the three star types.

The remaining subplots show example light curves from the OGLE dataset.

3.2 Anomaly scores

The process of evaluating the anomaly score requires us to compute the degree to which all points in the test data conform to our understanding of typical light curves. The log-marginal likelihood can be used to achieve this:

$$\log p(y_t) = \log \int p(y_t|f)p(f)df = -\frac{1}{2}y_t^T \mathbf{K}^{-1}y_t - \frac{1}{2}\log |\mathbf{K}| - \frac{n}{2}\log 2\pi. \quad (8)$$

where f is the latent distribution over hierarchical functions. A natural anomaly score is then the expectation of the over the time series of the negative log-marginal likelihood $\mathbb{E}_t [-\log p(y_t)]$, for which the empirical analogue is:

$$S(\mathbf{y}) = -\sum_{t=1}^T \log p(y_t) = \sum_{t=1}^T \frac{1}{2}y_t^T \mathbf{K}^{-1}y_t + \frac{t}{2}\log |\mathbf{K}| + \frac{nt}{2}\log 2\pi. \quad (9)$$

We can now infer the HGP model, and use the anomaly score in order to evaluate the degree to which new instances conform to the hierarchical distribution.

4 Results

We first validate our method on a benchmark dataset for which anomalous time-series are known, followed by analysis on

the light-curve data. Both datasets are available from the UCR Time Series Classification Archive [12]. PCAD and RAND-C [2] are used as baselines throughout.

MALLAT is a synthetic dataset initially created for the research of wavelets in signal processing [13]. This dataset consists of 8 classes, with 300 examples

Table 1: Precision for MALLAT & OGLE.

Data size	MALLAT			OGLE		
	0.4	0.7	1.0	0.01	0.08	0.1
RAND-C	0.01	0.03	0.01	0.57	0.59	0.77
SPCAD	0.01	0.01	0.01	0.40	0.16	0.12
HGPAD	1.00	1.00	1.00	0.98	0.99	0.99

per class. Each example has 1024 time points. Classes 3 and 6 are visually rather similar to one another and hence we consider these to be the ‘inlier’ classes.

4.1 Detection of known anomalies

We can quantify the ability of our model to detect anomalous data by performing inference with HGPAD on one class, and testing the anomaly score on all other classes.

In this experiment we apply HGPAD to infer an underlying function for MALLAT dataset by using its different proportions of the inliers as training instances, then construct a mixed test set including 10% outliers (class 7) and 90% normal instances (classes 3 and 6) to test the precision of our model. We treat the OGLE dataset in a similar manner, where we use CEPH and RRL as the inlier classes and EB as the outlier class. Since we know the number of anomalies in the testing dataset (n), the precision can be calculated by looking at the top n instances in the output of HGPAD model.

We investigate the effect of dataset size to the robustness of anomaly detection in our HGPAD model. We show the results in Table 1 for the MALLAT and OGLE datasets. We can see that the proposed model consistently out-performs the baseline models in terms of precision, even when using small fractions of the total dataset. We believe that the implicit uncertainty quantification of the HGPAD model contributes to this since it reduces the likelihood of overfitting.

4.2 Detection of unknown anomalies

For the analysis of unknown anomalies within each class of light curves, we first infer a HGPAD model on a particular light curve class and compute anomaly scores. We then rank the held out dataset from least anomalous to most anomalous. We illustrate this idea in Figs. 3a to 3c for CEPH, EB, and RRL light curves respectively. In each figure, the leftmost column depicts the underlying

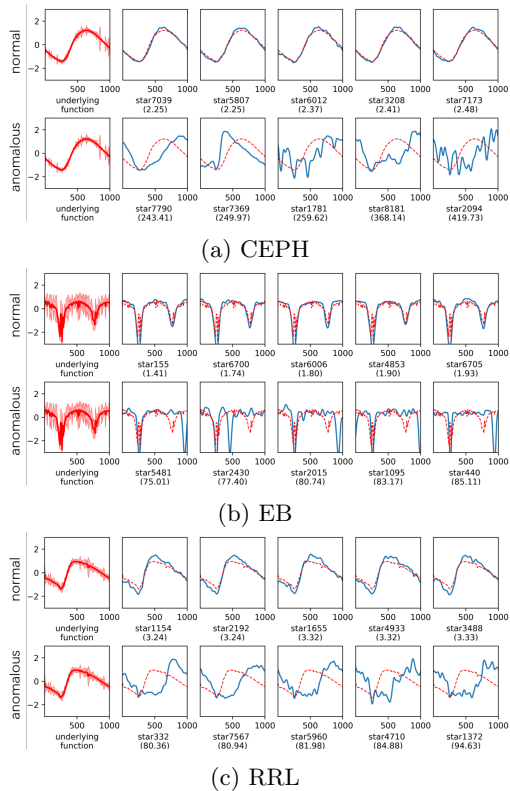


Fig. 3: Normal (odd rows) and anomalous (even rows) star light curves.

latent light curve function that was inferred by HGPAD. The top row depicts the 5 light curves that received the smallest anomaly score. We can see that there is a close match between the latent function and these instances in this figure. The bottom row of each sub-figure shows the instances that received the highest anomaly score. We can also see that the light curves do not closely resemble any of the prototypical curves, and are likely to be outlier instances. These could then be checked by an expert and either assigned to one of the other existing classes or classified as an entirely new class of star.

5 Conclusions

This paper introduces hierarchical Gaussian process anomaly detection for time-series and starlight curves. The approach used is a simple generalisation of Gaussian processes. While we incur several additional parameters, the probabilistic programming framework will optimise these during inference. We show experimentally that HGPAD is able to adapt to different situations and consistently outperforms baseline methods. We also show that our model is able to generalise well after using only 1% of the available training cases.

References

- [1] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):15, 2009.
- [2] Umaa Rebbapragada, Pavlos Protopapas, Carla E. Brodley, and Charles Alcock. Finding anomalous periodic time series. *Machine Learning*, 74(3):281–313, 12 2008.
- [3] OGLE. <http://ogle.astrouw.edu.pl/>.
- [4] Dragomir Yankov et al. Disk aware discord discovery: Finding unusual time series in terabyte sized datasets. *Knowledge and Information Systems*, 17(2):241–262, 2008.
- [5] Christiaan Sterken and Carlos Jaschek. *Light curves of variable stars: a pictorial atlas*. Cambridge University Press, 2005.
- [6] James Hensman et al. Hierarchical Bayesian modelling of gene expression time series across irregularly sampled replicates and clusters. *BMC bioinformatics*, 14(1):252, 2013.
- [7] Markus M Breunig et al. LOF: identifying density-based local outliers. In *ACM sigmod record*, volume 29, pages 93–104. ACM, 2000.
- [8] MJ Desforges et al. Applications of probability density estimation to the detection of abnormal conditions in engineering. *Proceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering Science*, 212(8):687–703, 1998.
- [9] Kira Kowalska and Leto Peel. Maritime anomaly detection using Gaussian Process active learning. In *Information Fusion (FUSION), 2012 15th International Conference on*, pages 1164–1171. IEEE, 2012.
- [10] Mark Smith et al. Maritime abnormality detection using Gaussian processes. *Knowledge and Information Systems*, 38(3):717–741, 08 2013.
- [11] Carl Edward Rasmussen and Christopher KI Williams. Gaussian processes in machine learning. *Lecture notes in computer science*, 3176:63–71, 2004.
- [12] Yanping Chen, Eamonn Keogh, Bing Hu, Nurjahan Begum, Anthony Bagnall, Abdullah Mueen, and Gustavo Batista. The UCR time series classification archive, July 2015. www.cs.ucr.edu/~eamonn/time_series_data/.
- [13] Stéphane Mallat. *A wavelet tour of signal processing*. Academic press, 1999.