

pgmuvi: Quick and easy Gaussian Process Regression for multi-wavelength astronomical timeseries

Peter Scicluna^{1,2}, Stefan Waterval^{3,4}, Diego A. Vasquez-Torres⁵, Sundar Srinivasan⁵, and Sara Jamal⁶

¹ European Southern Observatory, Alonso de Córdova 3107, Vitacura, Santiago, Chile ² Space Science Institute, 4750 Walnut Street, Suite 205, Boulder, CO 80301, USA ³ New York University Abu Dhabi, PO Box 129188, Abu Dhabi, United Arab Emirates ⁴ Center for Astro, Particle and Planetary Physics (CAP³), New York University Abu Dhabi, PO Box 129188, Abu Dhabi, United Arab Emirates ⁵ IRyA, Universidad Nacional Autónoma de México, Morelia, Michoacán, México ⁶ Max Planck Institute for Astronomy, Königstuhl 17, 69117 Heidelberg, Germany ¶ Corresponding author

DOI: 10.xxxxxx/draft

Software

- Review
- Repository
- Archive

Editor: Axel Donath

Reviewers:

- @joshspeagle
- @baptklein

Submitted: 31 July 2023

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License (CC BY 4.0).

Summary

Time-domain observations are increasingly important in astronomy, and are often the only way to study certain objects. The volume of time-series data is increasing dramatically as new surveys come online - for example, the Vera Rubin Observatory will produce 15 terabytes of data per night, and its Legacy Survey of Space and Time (LSST) is expected to produce five-year lightcurves for $> 10^7$ sources, each consisting of 5 photometric bands. Historically, astronomers have worked with Fourier-based techniques such as the Lomb-Scargle periodogram or information-theoretic approaches; however, in recent years Bayesian and data-driven approaches such as Gaussian Process Regression (GPR) have gained traction. However, the computational complexity and steep learning curve of GPR has limited its adoption. pgmuvi makes GPR of multi-band timeseries accessible to astronomers by building on cutting-edge open-source machine-learning libraries, and hence pgmuvi retains the speed and flexibility of GPR while being easy to use. It provides easy access to GPU acceleration and Bayesian inference of the hyperparameters (e.g. the periods), and is able to scale to large datasets.

Statement of need

Astronomical objects are in general not static, but vary in brightness over time. This is especially true for objects that are variable by nature, such as pulsating stars, or objects that are variable due to their orbital motion, such as eclipsing binaries. The study of these objects is called time-domain astronomy, and is a rapidly growing field. A wide range of approaches to time-series analysis have been developed, ranging from simple period-finding algorithms to more complex machine learning techniques (e.g. Donoso-Oliva et al., 2023; Friedman, 1984; Huijse et al., 2018; Palmer, 2009 and many more). Perhaps the most popular in astronomy is the Lomb-Scargle periodogram (Lomb, 1976; Scargle, 1982), which is a Fourier-based technique to find periodic signals in unevenly sampled data. However, the handling of unevenly sampled data is not the only challenge in time-series analysis.

A particular challenge in astronomy is handling heterogeneous, multiwavelength data (VanderPlas & Ivezić, 2015). Data must often be combined from a wide variety of instruments, telescopes or surveys, and so the systematics or noise properties of different datasets vary widely. In addition, by combining multiple wavelengths, we gain a better understanding of the physical processes driving the variability of the object. For example, some variability mechanisms differ as a function of wavelength only in amplitude (e.g. eclipsing binaries), while others may vary in

42 phase (e.g. pulsating stars) or even period (e.g. multiperiodic systems). Thus, it is important
43 to combine data from multiple wavelengths in a way that accounts for these differences.

44 Gaussian processes (GPs) have recently become a popular tool to handle these challenges.
45 GPs are a flexible way to forward-model arbitrary signals, by assuming that the signal is drawn
46 from a multivariate Gaussian distribution. By constructing a covariance function describing the
47 covariance between any two points in the signal, we can model the signal as a Gaussian process.
48 By doing so, we are freed from any assumptions about sampling, and can model the signal
49 as a continuous function. We are also able to model the noise in the data, and thus account
50 for heteroscedastic noise. We also gain the ability to directly handle multiple wavelengths, by
51 constructing a multi-dimensional covariance function. Hence, Gaussian Process Regression
52 (GPR) is a machine learning technique that is able to model non-periodic signals in unevenly
53 sampled data, and is thus well suited for the analysis of astronomical time-series data.

54 However, GPR is not without its challenges. The most popular covariance functions are often
55 not tailored to modelling complex signals, which implies that users must construct their own
56 covariance functions. GPs are also computationally expensive, and thus approximations must
57 be used to scale to large datasets. Finally, many GP packages either have very steep learning
58 curves or only provide limited features, and thus are not suitable for the average astronomer.

59 In this paper we present a new Python package, `pgmuvi`, to perform GPR on multi-wavelength
60 astronomical time-series data. The package is designed to be easy to use, and provide a quick
61 way to perform GPR on multi-wavelength data by exploiting uncommon but powerful kernels
62 and approximations. The package is designed to be flexible and allow the user to customize
63 the GPR model to their needs. `pgmuvi` exploits multiple strategies to scale regression to large
64 datasets, making it suitable for the current era of large-scale astronomical surveys.

65 A number of other software packages exist to perform GPR, some of which, such as `celerite`
66 (D. Foreman-Mackey et al., 2017; D. Foreman-Mackey, 2018), `tinygp` (Daniel Foreman-
67 Mackey, 2023) or `george` (Ambikasaran et al., 2015) were developed within the astronomical
68 community with astronomical time-series in mind. However, these each have their own
69 limitations. `celerite` is extremely fast, but is limited to one-dimensional inputs, and thus
70 cannot handle multiwavelength data, except under certain restrictive assumptions. Furthermore,
71 since it achieves its speed through a specific form of kernel decomposition, it is not able to
72 handle arbitrary covariance functions. It is therefore restricted to a small number of kernels
73 with specific forms; while it is able to handle the most common astronomical timeseries by
74 combining these kernels, not all signals can be modelled in this way. `tinygp` is a more general
75 package, able to retain a high degree of flexibility while still being fast thanks to a JAX-based
76 implementation, which makes it feasible to implement models in `tinygp` that are equivalent to
77 those in `pgmuvi`. However, `pgmuvi` is designed to provide an easier learning curve by packaging
78 GPs with data transforms and inference routines. In essence, `tinygp` could in principle be used
79 by `pgmuvi` as a GP backend instead of `GPyTorch`. For a summary of the state of the art of
80 GPR in astronomy, see the recent review by Aigrain & Foreman-Mackey (2023).

81 `pgmuvi` is used in two ongoing projects by our group: one of the authors' (DAVT) masters
82 thesis and the paper resulting from this work deals with the analysis of multiwavelength
83 light curves for targets from the Nearby Evolved Stars Survey (NESS; Scicluna et al. (2022),
84 <https://evolvedstars.space>). This work served as the first test of the code and has analyzed
85 thousands of light curves at optical and infrared wavelengths for over seven hundred dusty stars
86 within 3 kpc of the Solar Neighborhood. The paper will be published in 2023 (Vasquez-Torres
87 et al., in prep.). A different project related to dusty variable stars in M33 has also used `pgmuvi`
88 to estimate the periods of these objects from infrared light curves. This work will be published
89 in 2023 (Srinivasan et al., in prep.).

Method and Features

pgmuvi builds on the popular GPyTorch library. GPyTorch (Gardner et al., 2018) is a Gaussian process library for PyTorch (Paszke et al., 2019), which is a popular machine learning library for Python. By default, pgmuvi exploits the highly-flexible Spectral Mixture kernel (Wilson & Adams, 2013) in GPyTorch, able to model a wide variety of signals. This kernel is particularly interesting for astronomical time-series data, as it can effectively model multi-periodic and quasi-periodic signals. The spectral mixture kernel models the power spectrum of the covariance matrix as Gaussian mixture model (GMM), making it highly flexible, easy to interpret and adaptable to multi-dimensional data. This kernel is known for its ability to extrapolate, and is thus well suited to cases where prediction is important (for example, preparing astronomical observations of variable stars). By modelling the power spectrum in this way, pgmuvi effectively filters out noise in the periodogram, and thus is able to find the dominant periods in noisy data more effectively than, for example, the Lomb-Scargle periodogram.

However, the flexibility of this kernel comes at a cost; for more than one component in the mixture, the solution space becomes highly non-convex, and the optimization of the kernel hyperparameters becomes difficult. pgmuvi addresses this by first exploiting the Lomb-Scargle periodogram to find the dominant periods in the data, and then using these periods as initial guesses for the means of the mixture components.

Multiple options are available to accelerate inference depending on the size of the dataset. For small datasets, the exact GPs can be used to scale to datasets of up to ~ 1000 points. pgmuvi can exploit the Structured Kernel Interpolation (SKI) approximation (Wilson & Nickisch, 2015) to scale to datasets of up to $\sim 10^5$ points. Future work will include implementing approximations for even larger datasets: pgmuvi can in principle exploit the Sparse Variational GP (SVGP) or Variational Nearest Neighbour (VNN) approximations (Hensman et al., 2013; Wu et al., 2022) to scale to datasets of arbitrary size. pgmuvi can employ also GPU computing for both exact and variational GPs.

For exact GPs and SKI, pgmuvi performs maximum a posteriori (MAP) estimation of the hyperparameters, or full Bayesian inference. MAP estimation can exploit any PyTorch optimizer, but by default it uses ADAM (Kingma & Ba, 2014). Bayesian inference uses the pyro (Bingham et al., 2018) implementation of the No-U-Turn Sampler (NUTS) (Hoffman et al., 2014), which is a Hamiltonian Monte Carlo (HMC) sampler.

Finally, by allowing arbitrary GPyTorch likelihoods to be used, pgmuvi can be extended to a wide range of problems. For example, an instance of `gpytorch.likelihoods.StudentTLikelihood` can be dropped in to turn pgmuvi into a T-Process regressor, or missing data can be handled using `gpytorch.likelihoods.GaussianLikelihoodWithMissingObs`.

To summarise, the key features of pgmuvi are:

- Highly flexible kernel, able to model a wide range of multiwavelength signals
- Able to exploit multiple strategies to scale to large datasets
- GPU acceleration for both exact and variational GPs
- Fully Bayesian inference using HMC
- Initialisation of kernel hyperparameters using Lomb-Scargle periodogram
- Automated creation of diagnostic and summary plots
- Automated reporting of kernel hyperparameters and their uncertainties, and summary of MCMC chains.

Acknowledgements

This project was developed in part at the 2022 Astro Hack Week, hosted by the Max Planck Institute for Astronomy and Haus der Astronomie in Heidelberg, Germany. This work was

137 partially supported by the Max Planck Institute for Astronomy, the European Space Agency,
 138 the Gordon and Betty Moore Foundation, the Alfred P. Sloan foundation.
 139 SS and DAVT acknowledge support from UNAM-PAPIIT Program IA104822.

140 References

- 141 Aigrain, S., & Foreman-Mackey, D. (2023). Gaussian process regression for astronomical time
 142 series. *Annual Review of Astronomy and Astrophysics*, 61(1), null. <https://doi.org/10.1146/annurev-astro-052920-103508>
 143
- 144 Ambikasaran, S., Foreman-Mackey, D., Greengard, L., Hogg, D. W., & O'Neil, M. (2015).
 145 Fast Direct Methods for Gaussian Processes. *IEEE Transactions on Pattern Analysis and*
 146 *Machine Intelligence*, 38, 252. <https://doi.org/10.1109/TPAMI.2015.2448083>
- 147 Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., Singh,
 148 R., Szerlip, P., Horsfall, P., & Goodman, N. D. (2018). Pyro: Deep Universal Probabilistic
 149 Programming. *Journal of Machine Learning Research*.
- 150 Donoso-Oliva, C., Becker, I., Protopapas, P., Cabrera-Vives, G., Vishnu, M., & Vardhan, H.
 151 (2023). ASTROMER. A transformer-based embedding for the representation of light curves.
 152 670, A54. <https://doi.org/10.1051/0004-6361/202243928>
- 153 Foreman-Mackey, D. (2018). Scalable Backpropagation for Gaussian Processes using Celerite.
 154 *Research Notes of the American Astronomical Society*, 2(1), 31. <https://doi.org/10.3847/2515-5172/aaaf6c>
 155
- 156 Foreman-Mackey, Daniel. (2023). *dfm/tinygp: The tiniest of Gaussian Process libraries*
 157 (Version v0.2.4rc1). Zenodo. <https://doi.org/10.5281/zenodo.7646759>
- 158 Foreman-Mackey, D., Agol, E., Ambikasaran, S., & Angus, R. (2017). Fast and Scalable
 159 Gaussian Process Modeling with Applications to Astronomical Time Series. 154, 220.
 160 <https://doi.org/10.3847/1538-3881/aa9332>
- 161 Friedman, J. H. (1984). A variable span scatterplot smoother. Laboratory for Computational
 162 Statistics, Stanford University Technical Report, 5.
- 163 Gardner, J., Pleiss, G., Weinberger, K. Q., Bindel, D., & Wilson, A. G. (2018). Gpytorch:
 164 Blackbox matrix-matrix gaussian process inference with gpu acceleration. *Advances in*
 165 *Neural Information Processing Systems*, 31.
- 166 Hensman, J., Fusi, N., & Lawrence, N. D. (2013). Gaussian processes for big data. *arXiv*
 167 *Preprint arXiv:1309.6835*.
- 168 Hoffman, M. D., Gelman, A., & others. (2014). The no-u-turn sampler: Adaptively setting
 169 path lengths in hamiltonian monte carlo. *J. Mach. Learn. Res.*, 15(1), 1593–1623.
- 170 Huijse, P., Estévez, P. A., Förster, F., Daniel, S. F., Connolly, A. J., Protopapas, P., Carrasco,
 171 R., & Príncipe, J. C. (2018). Robust period estimation using mutual information for
 172 multiband light curves in the synoptic survey era. *The Astrophysical Journal Supplement*
 173 *Series*, 236(1), 12. <https://doi.org/10.3847/1538-4365/aab77c>
- 174 Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv Preprint*
 175 *arXiv:1412.6980*.
- 176 Lomb, N. R. (1976). Least-Squares Frequency Analysis of Unequally Spaced Data. 39(2),
 177 447–462. <https://doi.org/10.1007/BF00648343>
- 178 Palmer, D. M. (2009). A FAST CHI-SQUARED TECHNIQUE FOR PERIOD SEARCH
 179 OF IRREGULARLY SAMPLED DATA. *The Astrophysical Journal*, 695(1), 496. <https://doi.org/10.1088/0004-637X/695/1/496>
 180

- 181 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,
182 Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M.,
183 Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., ... Chintala, S. (2019). PyTorch:
184 An imperative style, high-performance deep learning library. In *Proceedings of the 33rd*
185 *international conference on neural information processing systems*. Curran Associates Inc.
- 186 Scargle, J. D. (1982). Studies in astronomical time series analysis. II. Statistical aspects of
187 spectral analysis of unevenly spaced data. 263, 835–853. <https://doi.org/10.1086/160554>
- 188 Scicluna, P., Kemper, F., McDonald, I., Srinivasan, S., Trejo, A., Wallström, S. H. J.,
189 Wouterloot, J. G. A., Cami, J., Greaves, J., He, J., Hoai, D. T., Kim, H., Jones, O. C.,
190 Shinnaga, H., Clark, C. J. R., Dharmawardena, T., Holland, W., Imai, H., van Loon, J. T.,
191 ... Zijlstra, A. A. (2022). The Nearby Evolved Stars Survey II: Constructing a volume-limited
192 sample and first results from the James Clerk Maxwell Telescope. 512(1), 1091–1110.
193 <https://doi.org/10.1093/mnras/stab2860>
- 194 VanderPlas, J. T., & Ivezić, Ž. (2015). Periodograms for Multiband Astronomical Time Series.
195 812(1), 18. <https://doi.org/10.1088/0004-637X/812/1/18>
- 196 Wilson, A., & Adams, R. (2013). Gaussian process kernels for pattern discovery and ex-
197 trapolation. In S. Dasgupta & D. McAllester (Eds.), *Proceedings of the 30th inter-*
198 *national conference on machine learning* (Vol. 28, pp. 1067–1075). PMLR. <https://proceedings.mlr.press/v28/wilson13.html>
199
- 200 Wilson, A., & Nickisch, H. (2015). Kernel interpolation for scalable structured gaussian
201 processes (KISS-GP). *International Conference on Machine Learning*, 1775–1784.
- 202 Wu, L., Pleiss, G., & Cunningham, J. P. (2022). Variational nearest neighbor gaussian process.
203 *International Conference on Machine Learning*, 24114–24130.