

Positive and Scale Invariant Gaussian Process Latent Variable Model for Astronomical Spectra

Nikolaos Gianniotis, Iliana Isabel Cortés Pérez and Kai Lars Polsterer*

Astroinformatics - HITS gGmbH
Schloss-Wolfsbrunnenweg 35, D-69118 Heidelberg, Germany

Abstract. We propose a probabilistic model that reduces the dimensionality of positive-valued data in a scale-invariant way, treating data items that differ only in scaling as identical. Extending the Gaussian Process Latent Variable Model, we ensure positive function values by applying a non-linear transformation to latent function values. To address the intractable marginal log-likelihood, we utilize a variational lower bound and amortized inference to reduce the number of variational parameters. We apply our model to reconstructing partially observed spectra and show how its scale-invariant property leads to better reconstructions.

1 Introduction

Gaussian processes (GPs) [1] are a versatile tool in machine learning, but imposing on them constraints like positivity, monotonicity, or physical constraints is challenging [2]. Past works have considered constraining GPs as solutions to differential equations [3], temporal and spectral reconstruction problems [4], or injecting domain-specific constraints via linear operators [5]. Other works constrain GP outputs with non-linear functions [6, 7], bound outputs to positive values by constraining the marginal likelihood [8], or cast linear constraints as conditional expectations of the truncated multivariate Gaussian distribution [9].

In this work, we aim to discover a latent space for positive-valued astronomical spectra. Amongst past works on dimensionality reduction for spectra [10, 11, 12], [13] uniquely incorporates a non-negativity constraint. We extend the Gaussian process latent variable model (GPLVM) [14] by bounding its outputs to positive values. The amplitude of astronomical spectra is not an intrinsic physical property and should not be reflected in the latent space. We introduce scale-invariance and show that it leads to better reconstructions.

1.1 Gaussian Process latent variable model

We consider the inference of a function $f(\cdot)$ from observed input-output pairs $\{\mathbf{x}_n \in \mathbb{R}^Q, y_n \in \mathbb{R}\}_{n=1}^N$. We take a probabilistic view and model $f(\cdot)$ with a GP:

$$f(\mathbf{x}) \sim \mathcal{GP}(m, k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})),$$

where $m \in \mathbb{R}$ specifies the unknown mean value of $f(\cdot)$ and $k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ specifies the covariance function [1]. A simple regression model for multidimensional

*The authors gratefully acknowledge the support of the Klaus Tschira Foundation. Software available under <https://github.com/HITS-AIN/GPLVMplus.jl>

outputs $\mathbf{y}_n \in \mathbb{R}^D$, models each d -th output dimension with an independent GP:

$$\log \int \prod_{d=1}^D \prod_{n=1}^N \mathcal{N}(y_{nd}|f_{nd}, \beta^{-1}) \mathcal{N}(\mathbf{f}_d|\mathbf{0}, \mathbf{K}) d\mathbf{f} = \sum_{d=1}^D \log \mathcal{N}(\mathbf{y}_d|\mathbf{0}, \mathbf{K} + \beta^{-1} \mathbf{I}_N) ,$$

where f_{nd} are the function values, and all D dimensions share the same $N \times N$ covariance \mathbf{K} (for brevity we suppress the dependence of \mathbf{K} on \mathbf{x} in the notation). This model can be also employed for dimensionality reduction and is known as the GPLVM [14]: we view \mathbf{y} as the images of unknown, latent low-dimensional input coordinates \mathbf{x} under a probabilistic map.

1.2 Basic statements

We use the following in the formulation of our model:

$$\langle \ln \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{C}) \rangle_{\mathcal{N}(\mathbf{x}|\mu, \Sigma)} = \ln \mathcal{N}(\mathbf{y}|\mathbf{A}\mu + \mathbf{b}, \mathbf{C}) - \frac{1}{2} \text{tr}(\Sigma \mathbf{A}^T \mathbf{C}^{-1} \mathbf{A}) .$$

We use the expectations wrt $\mathcal{N}(x|\mu, \sigma^2)$, $E(a, \mu, \sigma, b) \equiv \langle \exp(ax+b) \rangle = \exp(a\mu + \frac{(a\sigma)^2}{2} + b)$ and $B(a, \mu, \sigma, b) \equiv \langle \exp(ax+b)^2 \rangle = E(2a\mu, 2a\sigma, 2b)$. The variance of $\exp(ax+b)$ is $V(a, \mu, \sigma, b) \equiv B(a, \mu, \sigma, b) - E(a, \mu, \sigma, b)^2$.

2 Proposed model

We extend the GPLVM to enforce positivity and scale invariance. We call the extension GPLVM₊. The map between low-dimensional coordinates \mathbf{x} to high-dimensional data \mathbf{y} is modelled with D independent GPs. We impose positivity, by applying the $\exp(a \cdot f + b)$ function on the latent function values f of the map, where $a \geq 0$, b respectively control the amplitude and mean value of f . Additionally, we ensure scale invariance by assigning a scaling coefficient $s_n > 0$ to each data item, so that two data items differing only in scale are projected to the same location \mathbf{x} . The marginal log-likelihood of GPLVM₊ reads:

$$\ln \int \prod_{d=1}^D \prod_{n=1}^N \mathcal{N}(y_{nd}|s_n \exp(a f_{nd} + b), \beta^{-1}) \mathcal{N}(\mathbf{f}_d|\mathbf{0}, \mathbf{K}) d\mathbf{f}_d . \quad (1)$$

2.1 Derivation of lower bound

The introduced non-linearity prohibits analytically integrating out the latent values f_{nd} . We resort to variational inference [15] to approximate the marginal log-likelihood with a lower bound. Denoting the collection of all observations y_{nd} as \mathbf{Y} , latent values f_{nd} as \mathbf{F} and x_{nq} as \mathbf{X} , the lower bound reads:

$$\ln \int p(\mathbf{Y}, \mathbf{F}|\mathbf{X}) d\mathbf{F} \geq \int q(\mathbf{F}) \ln p(\mathbf{Y}, \mathbf{F}|\mathbf{X}) d\mathbf{F} - \int q(\mathbf{F}) \ln q(\mathbf{F}) d\mathbf{F} ,$$

where $q(\cdot)$ is a density we need to postulate, which if equal to the true posterior, makes the lower bound equal to the marginal log-likelihood. We postulate a

factorised posterior $q(\mathbf{F}) = \prod_{d=1}^D q(\mathbf{f}_d)$, where $q(\mathbf{f}_d) = \mathcal{N}(\mathbf{f}_d | \boldsymbol{\mu}_d, \boldsymbol{\Sigma})$, with $\boldsymbol{\mu}_d \in \mathbb{R}^N$ and all D dimensions share the same $N \times N$ covariance $\boldsymbol{\Sigma}$. In order to reduce the number of free parameters in $\boldsymbol{\Sigma}$, we parametrise, according to [16], as $\boldsymbol{\Sigma} = (\mathbf{K}^{-1} + \boldsymbol{\Lambda})^{-1}$, where $\boldsymbol{\Lambda}$ is a $N \times N$ diagonal matrix of positive elements. For this choice, the lower bound reads:

$$\left\langle \sum_{d=1}^D \sum_{n=1}^N \ln \mathcal{N}(y_{nd} | s_n \exp(a f_{nd} + b), \beta^{-1}) \right\rangle_q + \left\langle \sum_{d=1}^D \ln \mathcal{N}(\mathbf{f}_d | \mathbf{0}, \mathbf{K}) \right\rangle_q + \mathcal{H}[q], \quad (2)$$

where the entropy term reads $\mathcal{H}[q] = \frac{D}{2} \ln((2\pi e)^N |\boldsymbol{\Sigma}|)$. Given the basic results in Section 1.2, the first expectation in Eq. (2) reads:

$$\begin{aligned} & \sum_{d=1}^D \sum_{n=1}^N \frac{1}{2} \ln \left(\frac{\beta}{2\pi} \right) - \frac{\beta}{2} (y_{nd} - s_n E(a, \mu_{nd}, \sqrt{\Sigma_{nn}}, b))^2 - \frac{\beta s_n^2}{2} V(a, \mu_{nd}, \sqrt{\Sigma_{nn}}, b) = \\ & \sum_{d=1}^D \sum_{n=1}^N \ln \mathcal{N}(y_{nd} | s_n E(a, \mu_{nd}, \sqrt{\Sigma_{nn}}, b), \beta^{-1}) - \frac{\beta s_n^2}{2} V(a, \mu_{nd}, \sqrt{\Sigma_{nn}}, b). \end{aligned}$$

The second expectation in Eq. (2) reads $\sum_{d=1}^D \ln \mathcal{N}(\boldsymbol{\mu}_d | \mathbf{0}, \mathbf{K}) - \frac{D}{2} \text{tr}(\mathbf{K}^{-1} \boldsymbol{\Sigma})$. To learn the latent space, we optimise the lower bound wrt the free variational parameters $\boldsymbol{\mu}_d$, $\boldsymbol{\Lambda}$, kernel hyperparameters $\boldsymbol{\theta}$, latent coordinates \mathbf{x}_n and parameters s_n, a, b, β . We also note that only \mathbf{K} depends on \mathbf{x}_n .

2.2 Amortised inference

The number of free variational parameters in posterior $q(\mathbf{F})$ grows with the number of observations: $N \cdot D$ for the means $\boldsymbol{\mu}_d \in \mathbb{R}^N$ and N for $\boldsymbol{\Lambda}$. This leads to a computational burden on the optimizer. We also note that in the lower bound in Eq. (2), the latent coordinates \mathbf{x} in matrix \mathbf{K} are decoupled from the latent function values f_{nd} , i.e. changing \mathbf{x} does not influence directly f_{nd} , which exacerbates the optimisation problem.

We address both issues via amortised inference [17]: we parametrise the variational parameters $\boldsymbol{\mu}_d \in \mathbb{R}^N$ as the images of the latent coordinates \mathbf{x}_n under a neural network $g(\mathbf{x}_n; \mathbf{w}) : \mathbb{R}^Q \rightarrow \mathbb{R}^D$, where \mathbf{w} are the weights of the network. If g_d is the d -th output of the network, then $g_d(\mathbf{x}_n; \mathbf{w})$ returns μ_{nd} . Evaluating g_d on all \mathbf{x}_n , we obtain the vector $\boldsymbol{\mu}_d \in \mathbb{R}^N$. Thus, we optimise the lower bound of Eq. (2) wrt the weights \mathbf{w} of network g , $\boldsymbol{\Lambda}$, the latent coordinates \mathbf{x}_n and parameters $s_n, a, b, \beta, \boldsymbol{\theta}$ using scaled conjugate gradients[18].

2.3 Inference for test data

Given N_* number of yet unseen, test data items \mathbf{y}_n^* , we want to infer the corresponding latent coordinates \mathbf{x}_n^* . The corresponding latent function values are f_{nd}^* . The likelihood for the test data reads:

$$\prod_{n=1}^{N_*} \prod_{d=1}^D \mathcal{N}(y_{nd}^* | s_n^* \exp(a f_{nd}^* + b), \beta^{-1}). \quad (3)$$

The latent function values \mathbf{f}_d^* of the test data depend on the latent function values \mathbf{f}_d of the training data as follows [1]:

$$p(\mathbf{f}_d^*|\mathbf{f}_d) = \mathcal{N}(\mathbf{f}_d^*|\mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{f}_d, \mathbf{K}_{*,*} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*), \quad (4)$$

where $\mathbf{K}_* \in \mathbb{R}^{N_* \times N}$ is the cross-covariance between test and training data; $\mathbf{K}_{*,*} \in \mathbb{R}^{N_* \times N_*}$ is the covariance for the test data items¹. Integrating out \mathbf{f}_d in Eq. (4) with respect to posterior $q(\mathbf{f}_d)$ reads:

$$p(\mathbf{f}_d^*) = \mathcal{N}(\mathbf{f}_d^*|\mathbf{K}_*^T \mathbf{K}^{-1} \boldsymbol{\mu}_d, \mathbf{K}_{*,*} - \mathbf{K}_*^T (\mathbf{K} + \boldsymbol{\Lambda}^{-1})^{-1} \mathbf{K}_*).$$

Density $p(\mathbf{f}_d^*)$ acts as the prior for the latent function values for the testing data:

$$\ln \int \prod_{n=1}^{N_*} \prod_{d=1}^D \mathcal{N}(y_{nd}^*|s_n^* \exp(a f_{nd}^* + b), \beta^{-1}) p(\mathbf{f}_d^*) d\mathbf{f}_d. \quad (5)$$

We note the similarity of Eq. (5) to the marginal log-likelihood in Eq. (1). Again, we face an intractable integral and we resort to a variational lower bound: we postulate an approximating posterior $r(\mathbf{F}^*) = \prod_{d=1}^D r(\mathbf{f}_d^*)$, $r(\mathbf{f}_d^*) = \mathcal{N}(\mathbf{f}_d^*|\boldsymbol{\nu}_d, \mathbf{A})$. We parametrise the covariance as $\mathbf{A} = (\mathbf{C}^{-1} + \mathbf{L})^{-1}$, where \mathbf{L} is a $N_* \times N_*$ diagonal matrix of positive elements [16]. The lower bound reads:

$$\left\langle \ln \prod_{n=1}^{N_*} \prod_{d=1}^D \mathcal{N}(y_{nd}^*|s_n^* \exp(a f_{nd}^* + b), \beta^{-1}) + \ln \prod_{d=1}^D p(\mathbf{f}_d^*) \right\rangle_r + \mathcal{H}[r], \quad (6)$$

where $\mathcal{H}[r]$ is the entropy of $r(\mathbf{F}^*)$. The expectations in Eq. (6) are calculated in the same way as the ones in Section 2.1. To infer latent coordinates \mathbf{x}_n^* for the test data \mathbf{y}_n^* , we optimise the lower bound in Eq. (6) wrt $\boldsymbol{\nu}_d$, \mathbf{L} , \mathbf{x}_n^* and s_n^* .

3 Numerical Experiments

We use a two-hidden layer network $g(\cdot; \mathbf{w})$, each layer has 50 units, with tanh activations. We employ the kernel $k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = \theta_1^2 \exp(-\theta_2^2 \|\mathbf{x} - \mathbf{x}'\|^2)$.

Rubber duck. We use 64×64 images of a rubber duck photographed from 72 angles of a full rotation [19]. The images form an intrinsically two-dimensional closed curve. We scale each image with a factor sampled from the uniform density $\mathcal{U}(0.5, 1.5)$. We train the GPLVM and GPLVM₊ on the images for $Q = 2$. Fig. 1 demonstrates the scale invariance of the proposed model.

Spectra. We work with astronomical spectra from the SDSS survey² observed on a grid of 500 wavelengths. Each spectrum is observed only in some wavelengths, the rest are treated as missing values. We test the ability of the model to reconstruct partially observed test spectra: in each test spectrum we randomly remove 10% of the values and reserve them for checking model predictions. We

¹These matrices are actually functions of the latent coordinates i.e. $\mathbf{K}_*(\mathbf{X}^*)$ and $\mathbf{K}_{*,*}(\mathbf{X}^*)$. We refrain from denoting them as such in order to lighten notation.

²<https://www.sdss.org/>. We shift spectra to the restframe.

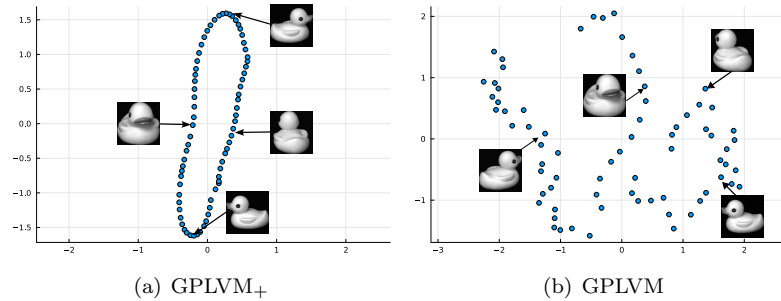


Fig. 1: GPLVM₊ finds the intrinsic structure of the data, while GPLVM cannot.

fix $Q = 3$. In a first experiment (Fig. 2a, 2b), we randomly take 1000 spectra from the BOSS catalogue for training, and 500 for testing. Spectra contain negative values that are artefacts due to noise or cosmic rays. The GPLVM and GPLVM₊ achieve a normalised mean squared error of 0.61 and 0.42 respectively. A 95% bootstrap confidence interval shows that the difference of ~ 0.19 is statistically significant and lies in (0.098, 0.362). In a second experiment, we take 500 spectra from the Stripe-82 catalogue that exhibit *great variation in amplitude*. The scale invariant GPLVM₊ produces good reconstructions, while GPLVM produces unrealistic spectra (Fig. 2c), unlike any spectra in the dataset.

4 Conclusion

We introduced the GPLVM₊ for dimensionality reduction of positive-valued data with scale-invariance. Numerical experiments validate its effectiveness. Currently, GPLVM₊ incurs high inference costs for latent coordinates of new data. We plan to tackle this in future research.

References

- [1] C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. MIT Press, 2006.
- [2] L. P. Swiler, M. Gulian, A. L. Frankel, C. Safta, and J. D. Jakeman. A survey of constrained gaussian process regression: Approaches and implementation challenges. *Journal of Machine Learning for Modeling and Computing*, 1(2), 2020.
- [3] M. Raissi, P. Perdikaris, and G. Em Karniadakis. Machine learning of linear differential equations using gaussian processes. *Journal of Computational Physics*, 348:683–693, 2017.
- [4] Felipe Tobar, Lerko Araya-Hernandez, Pablo Huijse, and Petar M. Djuric. Bayesian reconstruction of fourier pairs. *IEEE Transactions on Signal Processing*, 69:73–87, 2021.
- [5] Simo Särkkä. Linear operators and stochastic partial differential equations in gaussian process regression. In *ICANN 2011*, pages 151–158, 2011.
- [6] B. S. Jensen, J. B. Nielsen, and J. Larsen. Bounded gaussian process regression. In *IEEE International Workshop on Machine Learning for Signal Processing*, 2013.
- [7] Edward Snelson, Zoubin Ghahramani, and Carl Rasmussen. Warped gaussian processes. *Advances in neural information processing systems*, 16, 2003.

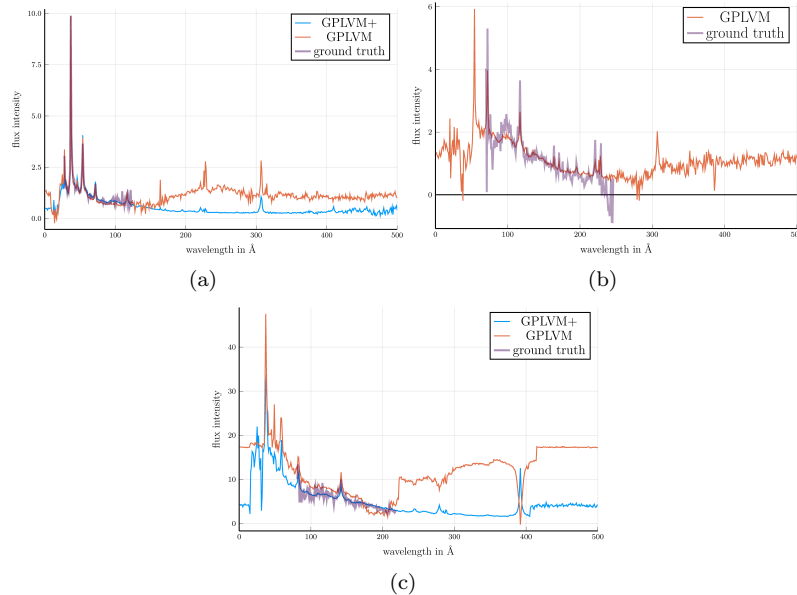


Fig. 2: (a) Both GPLVM and GPLVM₊ produce good reconstruction for BOSS data, but (b) GPLVM produces negative values for some cases. (c) GPLVM₊ produces good reconstructions for the Stripe-82 data, but the GPLVM does not.

- [8] Andrew Pensoneault, Xiu Yang, and Xueyu Zhu. Nonnegativity-enforced gaussian process regression. *Theoretical and Applied Mechanics Letters*, 10(3):182–187, 2020.
- [9] S. Da Veiga and A. Marrel. Gaussian process modeling with inequality constraints. *Annales de la Faculté des sciences de Toulouse: Mathématiques*, Ser. 6, 21(3):529–555, 2012.
- [10] S. D. Kügler, N. Gianniotis, and K. L. Polsterer. A spectral model for multimodal redshift estimation. In *2016 IEEE SSCI*, pages 1–8, 2016.
- [11] S. KN Portillo, J. K Parejko, J. R Vergara, and A. J Connolly. Dimensionality reduction of sdss spectra with variational autoencoders. *The Astronomical Journal*, 160(1):45, 2020.
- [12] AC. Eilers, DW. Hogg, B. Schölkopf, D. Foreman-Mackey, FB. Davies, and JT. Schindler. A generative model for quasar spectra. *The Astrophysical Journal*, 938(1):17, 2022.
- [13] Markus Harva and Ata Kabán. Variational learning for rectified factor analysis. *Signal Processing*, 87(3):509–527, 2007.
- [14] N. Lawrence. Gaussian process latent variable models for visualisation of high dimensional data. *Advances in neural information processing systems*, 16, 2003.
- [15] CM. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 2007.
- [16] M. Opper and C. Archambeau. The variational gaussian approximation revisited. *Neural Computation*, 21(3), 2009.
- [17] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations (ICLR)*, 2014.
- [18] Patrick Kofod Mogensen and Asbjørn Nilsen Riseth. Optim: A mathematical optimization package for Julia. *Journal of Open Source Software*, 3(24):615, 2018.
- [19] S. A. Nene, S. K. Nayar, and H. Murase. Columbia object image library (coil-20). Technical report, Technical Report CUCS-005-96, February 1996.