

Efficient Gaussian Neural Processes for Regression

Stratis Markou^{*1} James R. Requeima^{*12} Wessel Bruinsma¹² Richard E. Turner¹

Abstract

Conditional Neural Processes (CNP; Garnelo et al., 2018a) are an attractive family of meta-learning models which produce well-calibrated predictions, enable fast inference at test time, and are trainable via a simple maximum likelihood procedure. A limitation of CNPs is their inability to model dependencies in the outputs. This significantly hurts predictive performance and renders it impossible to draw coherent function samples, which limits the applicability of CNPs in downstream applications and decision making. Neural Processes (NPs; Garnelo et al., 2018b) attempt to alleviate this issue by using latent variables, relying on these to model output dependencies, but introduces difficulties stemming from approximate inference. One recent alternative (Bruinsma et al., 2021), which we refer to as the FullConvGNP, models dependencies in the predictions while still being trainable via exact maximum-likelihood. Unfortunately, the FullConvGNP relies on expensive $2D$ -dimensional convolutions, which limit its applicability to only one-dimensional data. In this work, we present an alternative way to model output dependencies which also lends itself maximum likelihood training but, unlike the FullConvGNP, can be scaled to two- and three-dimensional data. The proposed models exhibit good performance in synthetic experiments.

1. Introduction and Motivation

Conditional Neural Processes (CNP; Garnelo et al., 2018a) are a recently proposed class of meta-learning models which promises to combine the modelling flexibility, robustness, and fast inference of neural networks with the calibrated uncertainties of Gaussian processes (GPs; Rasmussen, 2003). CNPs are trained using a simple maximum-likelihood pro-

cedure and make predictions with complexity linear in the number of data points. Recent work has extended CNPs by incorporating attentive mechanisms (Kim et al., 2019) or accounting for symmetries in the prediction problem (Gordon et al., 2020; Kawano et al., 2021), achieving impressive performance on a variety of tasks.

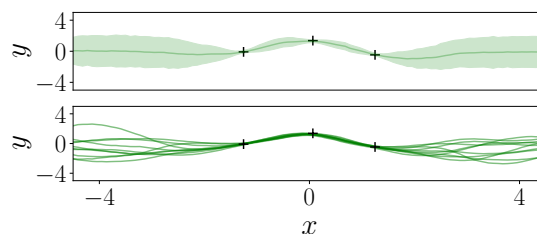


Figure 1. Unlike the ConvCNP (top) which produces only a marginal predictive, the ConvGNP (bottom) provides a correlated predictive, which can be used to draw coherent function samples.²

Despite these favourable qualities, CNPs are limited to predictions which do not model output dependencies, treating different input locations as independent (fig. 1). In this paper, we call such predictions *mean field*. The inability to model dependencies hurts performance and renders CNPs unable to produce coherent function samples, limiting their applicability in downstream applications. For example, in precipitation modelling, we might wish to evaluate the probability of the event that the amount of rainfall per day within some region remains above some specified threshold throughout a sustained length of time, which could help assess the likelihood of a flood. Mean-field predictions, which model every location independently, would assign unreasonably low probabilities to such events. If we were able to draw coherent samples from the predictive, however, then the probabilities of these events and numerous other useful quantities could be more reasonably estimated.

To address the inability of CNPs to model dependencies in the predictions, follow-up work (Garnelo et al., 2018b) introduced latent variables, introducing difficulties stemming from approximate inference (Le et al., 2018; Foong et al., 2020). More recently Bruinsma et al. (2021) introduced a variant of the CNP called the Gaussian Neural Process, hereafter referred to as the FullConvGNP, which directly parametrises the predictive covariance of the outputs. However, for D -dimensional data, the architecture of the FullConvGNP involves $2D$ -dimensional convolutions,

^{*}Equal contribution ¹Department of Engineering, University of Cambridge, Cambridge, UK ²Invenia Labs, Cambridge, UK. Correspondence to: Stratis Markou <em626@cam.ac.uk>, James R. Requeima <jrr41@cam.ac.uk>.

which are costly, and, for $D > 1$, poorly supported by most Deep Learning libraries. This work introduces an alternative method to directly parametrise output dependencies, which circumvents the costly convolutions of the FullConvGNP and can be applied to higher dimensional data.

2. Conditional and Gaussian Neural processes

Following the work of Foong et al. (2020), we present CNPs from the viewpoint of *prediction maps*. A prediction map π is a function which maps (1) a *context set* $(\mathbf{x}_c, \mathbf{y}_c)$ where $\mathbf{x}_c = (x_{c,1}, \dots, x_{c,N})$ are the inputs and $\mathbf{y}_c = (y_{c,1}, \dots, y_{c,N})$ the outputs and (2) a set of *target inputs* $\mathbf{x}_t = (x_{t,1}, \dots, x_{t,M})$ to a distribution over the corresponding *target outputs* $\mathbf{y}_t = (y_{t,1}, \dots, y_{t,M})$:

$$\pi(\mathbf{y}_t; \mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t) = p(\mathbf{y}_t | \mathbf{r}), \quad (1)$$

where $\mathbf{r} = r(\mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t)$ is a vector which parameterises the distribution over \mathbf{y}_t . For a fixed context set $(\mathbf{x}_c, \mathbf{y}_c)$, using Kolmogorov’s extension theorem (Oksendal, 2013), the collection of finite-dimensional distributions (f.d.d.s) $\pi(\mathbf{y}_t; \mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t)$ for all $\mathbf{x}_t \in \mathbb{R}^M$, $M \in \mathbb{N}$ defines a stochastic process if these f.d.d.s are consistent under (i) permutations of any entries of $(\mathbf{x}_t, \mathbf{y}_t)$ and (ii) marginalisations of any entries of \mathbf{y}_t — see appendix A. Prediction maps include, but are not limited to, Bayesian posteriors. One familiar example of such a map is the Bayesian GP posterior

$$\pi(\mathbf{y}_t; \mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{m}, \mathbf{K}), \quad (2)$$

where $\mathbf{m} = m(\mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t)$ and $\mathbf{K} = k(\mathbf{x}_c, \mathbf{x}_t)$ are given by the usual GP posterior expressions (Rasmussen, 2003). Another prediction map is the CNP (Garnelo et al., 2018a):

$$\pi(\mathbf{y}_t; \mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t) = \prod_{m=1}^M p(y_{t,m} | \mathbf{r}_m), \quad (3)$$

where each $p(y_{t,m} | \mathbf{r}_m)$ is an independent Gaussian and $\mathbf{r}_m = r(\mathbf{x}_c, \mathbf{y}_c, x_{t,m})$ is parameterised by a DeepSet³ (Zaheer et al., 2017). CNPs are permutation and marginalisation consistent and thus correspond to valid stochastic processes. However, CNPs do not respect the product rule in general — see appendix A and Foong et al. (2020). Nevertheless, CNPs and their variants (Gordon et al., 2020) have been demonstrated to give competitive performance and robust predictions in a variety of tasks and are a promising class of meta-learning models.

A central problem with the predictive in eq. (3) is that is mean field: eq. (3) does not model correlations between $y_{t,m}$ and $y_{t,m'}$ for $m \neq m'$. Mean field predictives severely hurt the predictive log-likelihood. In addition, one cannot

³The DeepSet ensures the posterior map is invariant to permutations of the context set. This is a desirable property in general, which should not be conflated with Kolmogorov consistency.

use a mean field predictive to draw coherent function samples. To remedy these issues, we consider parameterising a correlated multivariate Gaussian

$$\pi(\mathbf{y}_t; \mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{m}, \mathbf{K}) \quad (4)$$

where, instead of the expressions for the Bayesian GP posterior, we use neural networks to parameterise the mean $\mathbf{m} = m(\mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t)$ and covariance $\mathbf{K} = K(\mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t)$. We refer to this class of models as Gaussian Neural Processes (GNPs). The first such model, the FullConvGNP, was introduced by Bruinsma et al. (2021) with promising results. Unfortunately, the FullConvGNP relies on $2D$ -dimensional convolutions for parameterising \mathbf{K} , which are challenging to scale to higher dimensions. To overcome this difficulty we propose parameterising \mathbf{m} and \mathbf{K} by

$$\mathbf{m}_i = f(x_{t,i}, \mathbf{r}), \quad (5)$$

$$\mathbf{K}_{ij} = k(g(x_{t,i}, \mathbf{r}), g(x_{t,j}, \mathbf{r})) \quad (6)$$

where $\mathbf{r} = r(\mathbf{x}_c, \mathbf{y}_c)$, f and g are neural networks with outputs in \mathbb{R} and \mathbb{R}^{D_g} respectively, and k is an appropriately chosen positive-definite function. Note that, since k models a posterior covariance, it cannot be stationary. Equations (5) and (6) define a class of GNPs which, unlike the FullConvGNP, do not require costly convolutions. GNPs can be readily trained via the log-likelihood

$$\theta^* = \arg \max_{\theta} \log \pi(\mathbf{y}_t; \mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t), \quad (7)$$

also used in (Garnelo et al., 2018a), where θ collects all the parameters of the neural networks f , g , and r . In this work, we consider two methods to parameterise \mathbf{K} . The first method is the `linear` covariance

$$\mathbf{K}_{ij} = g(x_{t,i}, \mathbf{r})^\top g(x_{t,j}, \mathbf{r}) \quad (8)$$

which can be interpreted as a linear-in-the-parameters model with D_g basis functions and a unit Gaussian distribution on their weights. This model meta-learns D_g context dependent basis functions, which attempt to best approximate the true distribution of the target given the context. By Mercer’s theorem (Rasmussen, 2003), up to regularity conditions, every positive-definite function k can be decomposed as

$$k(z, z') = \sum_{d=1}^{\infty} \phi_d(z) \phi_d(z') \quad (9)$$

where $(\phi_d)_{d=1}^{\infty}$ is a set of orthogonal basis functions. We therefore expect eq. (8) to be able to recover arbitrary (sufficiently regular) GP predictives as D_g grows large. Further, the `linear` covariance has the attractive feature that sampling from it scales linearly with the number of query locations. A drawback is that the finite number of basis functions may limit its expressivity. An alternative method which sidesteps this issue, is to parametrise \mathbf{K} using the `kvv` covariance:

$$\mathbf{K}_{ij} = k(g(x_{t,i}, \mathbf{r}), g(x_{t,j}, \mathbf{r})) v(x_{t,i}, \mathbf{r}) v(x_{t,j}, \mathbf{r}) \quad (10)$$

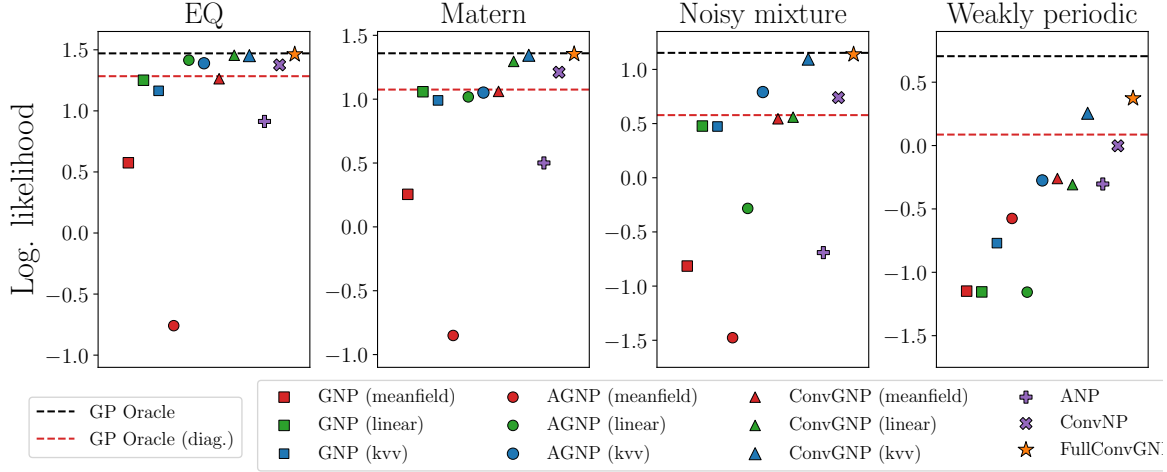


Figure 2. Predictive log-likelihood performance of the models, across datasets. The oracle GP performance is shown in dashed black. The dashed red line marks the performance of a modified version of the oracle, in which the off-diagonal terms of the covariance are set to 0.

where k is the Exponentiated Quadratic (EQ) covariance and v is a neural network with its output in \mathbb{R} . The v factors modulate the magnitude of the covariance, which would otherwise not be able to shrink near the context points. Unlike `linear`, `kvv` is not limited by a finite number of basis functions. A drawback of `kvv` is that the cost of drawing samples from it, scales cubically in the number of query locations, which may impose important practical limitations.

Both `linear` and `kvv` leave room for choosing f , g , and r according to the task at hand, giving rise to a collection of different models of the GNP family. For example, we may choose these to be feedforward DeepSets, giving rise to Gaussian Neural Processes (GNPs); attentive DeepSets, giving rise to Attentive Gaussian Neural Processes (AGNPs); or convolutional architectures, giving rise to Convolutional Gaussian Neural Processes (ConvGNPs). In this work, we explore these three alternatives, proposing the ConvGNP as a scalable alternative to the FullConvGNP. This approach can be extended to multiple outputs, which we will address in future work.

3. Experiments

We apply the proposed models to synthetic datasets generated from GPs with various covariance functions and known hyperparameters. We sub-sample these datasets into context and target sets, and train via the log-likelihood (eq. (7)).

We also train the ANP and ConvNP models as discussed in Foong et al. (2020). These latent variable models place a distribution q over \mathbf{r} and rely on q for modelling output dependencies. Following Foong et al. we train the ANP and ConvNP via a biased Monte Carlo estimate of the objective

$$\theta^* = \arg \max_{\theta} \log \left[\mathbb{E}_{\mathbf{r} \sim q(\mathbf{r})} [p(\mathbf{y}_t; \mathbf{x}_c, \mathbf{y}_c, \mathbf{x}_t, \mathbf{r})] \right]. \quad (11)$$

Figure 2 compares the predictive log-likelihood of the models, evaluated on in-distribution data, from which we observe the following trends.

Dependencies improve performance: We expected that modelling output dependencies would allow the models to achieve better log-likelihoods. Indeed, for a fixed architecture, we see that the correlated GNPs (■, ■, ●, ●, ▲, ▲) typically outperform their mean-field counterparts (■, ●, ▲). This result is encouraging and suggests that GNPs can learn meaningful dependencies in practice, in some cases recovering oracle performance.

Comparison with the FullConvGNP: The correlated ConvGNPs (▲, ▲) are often competitive with the FullConvGNP (★). The `kvv` ConvGNP (▲) is the only model, from those examined here, which competes with the FullConvGNP in all tasks. Unlike the latter, however, the former is scalable to $D = 2, 3$ dimensions.

Comparison with the ANP and ConvNP: Correlated GNPs typically outperform the latent-variable ANP (⊕) and ConvNP (⊗) models, which could be explained by the fact that the GNPs have a Gaussian predictive while ANP and ConvNP do not, and all tasks are Gaussian. Despite experimenting with different architectures, and even allowing for many more parameters in the ANP and ConvNP compared to the AGNP (●, ●) and ConvGNP (▲, ▲), we found it difficult to make the latent variable models competitive with the GNPs. We typically found the GNP family easier to train than these latent variable models.

Kvv outperformed linear: We generally observed that the `kvv` models (■, ●, ▲) performed as well, and occasionally better than, their `linear` counterparts (■, ●, ▲). To test whether the `linear` models were limited by the number of basis functions D_g , we experimented with various

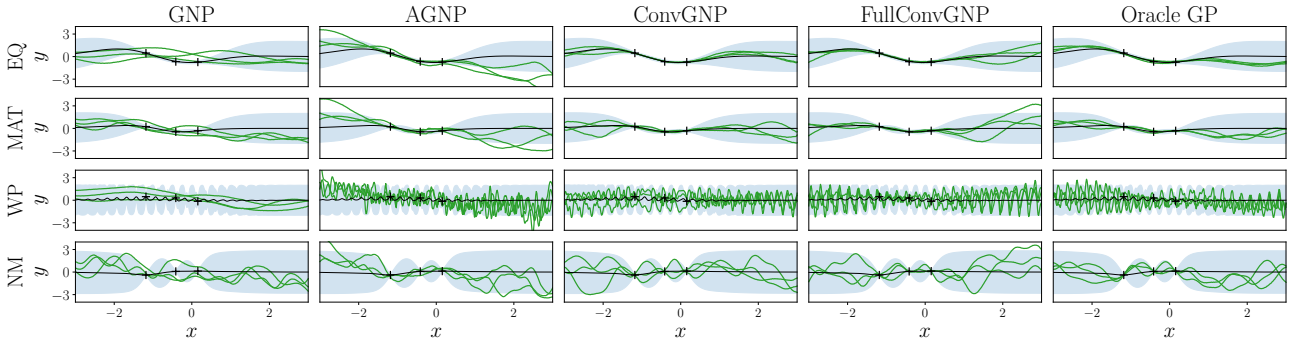


Figure 3. Samples drawn from the models’ predictive posteriors (green) compared to the ground truth (blue), using the `kvv` covariance. Note that the models are trained with context and target points uniformly distributed in the $[-2, 2]$ range.

settings $D_g \in \{16, 128, 512, 2048\}$. We did not observe a performance improvement for large D_g , suggesting that the models are not limited by this factor. This is surprising because, as $D_g \rightarrow \infty$ and assuming flexible enough f , g , and r , the linear models should, by Mercer’s theorem, be able to recover any (sufficiently regular) GP posterior. From preliminary investigations, we leave open the possibility that the linear models might be more difficult to optimise and thus struggle to compete with `kvv`. We hope to conduct a more careful study on our training protocol in the future, to determine whether the training method can account for this performance gap, or whether the `kvv` model is fundamentally more powerful than the linear model.

Figure 3 shows samples drawn from the GNP models, from which we qualitatively observe that, like the FullConvGNP, the ConvGNP produces good quality function samples. These samples are consistent with the observed data, whilst maintaining uncertainty and capturing the behaviour of the underlying process. The ConvGNP is the only conditional model (other than the FullConvGNP) which produces high-quality posterior samples. Figure 4 shows plots of the models’ covariances. Observe that, like the FullConvGNP, the ConvGNP is able to recover intricate covariance structure.

4. Conclusion and further work

This work introduced an alternative method for parametrising a correlated Gaussian predictive in CNPs. This approach can be combined with existing CNP architectures such as the feedforward (GNP), attentive (AGNP), or convolutional networks (ConvGNP). The resulting models are computationally cheaper and easier to scale to higher dimensions than the existing FullConvGNP of Bruinsma et al. (2021), whilst still being trainable via exact maximum-likelihood. The ConvGNP outperforms the other conditional and latent-variable models which we consider in this work, with the exception of the FullConvGNP. We found that modelling dependencies in the output improves the predictive log-likelihood over mean-field models. It also allows us to draw coherent function samples, which means that GNPs

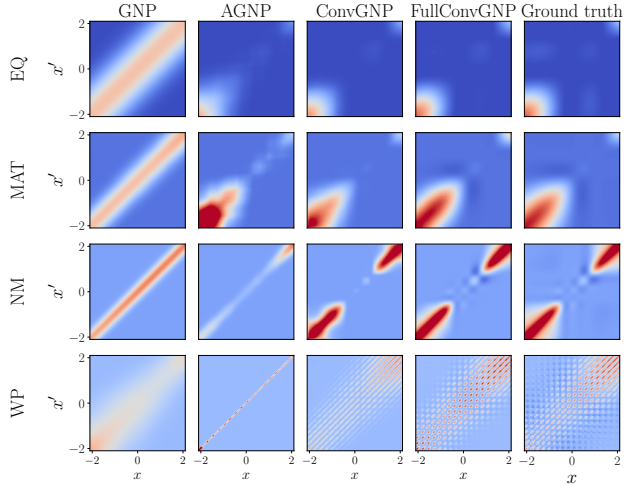


Figure 4. Plots of the predictive covariance of the GNP class of models, using the `kvv` covariance. The FullConvGNP and ground truth covariances are shown for comparison.

can be chained with more elaborate downstream estimators. Unlike the ANP and ConvNP models whose predictive is non-analytic, we expect the evaluation of, e.g., Active Learning acquisition functions to be significantly easier and more tractable in GNPs, a use-case we hope to explore in future work. We also note that, although ConvGNPs exhibit favourable scaling over FullConvGNPs, they still require 2- or 3-dimensional convolutions when applied to higher dimensions, which are also very costly. We wish to explore ways to reduce this cost, as well as to how other kinds of equivariances such as rotational and reflective equivariance (Kawano et al., 2021; Holderrieth et al., 2020) can be scaled to higher dimensions, in a computationally cheaper manner. For this, an approach similar to the work of Satorras et al. (2021) in the context of equivariant GNNs is a promising direction. We believe that cheap and scalable conditional neural processes for higher dimensional data could be highly valuable in a wide range of applications, including weather and environmental modelling, simulations, graphics and vision.

5. Acknowledgements

Richard E. Turner is supported by Google, Amazon, ARM, Improbable and EPSRC grant EP/T005386/1.

References

- Bruinsma, W. P., Requeima, J., Foong, A. Y. K., Gordon, J., and Turner, R. E. The gaussian neural process, 2021.
- Foong, A. Y. K., Bruinsma, W. P., Gordon, J., Dubois, Y., Requeima, J., and Turner, R. E. Meta-learning stationary stochastic process prediction with convolutional neural processes, 2020.
- Garnelo, M., Rosenbaum, D., Maddison, C. J., Ramalho, T., Saxton, D., Shanahan, M., Teh, Y. W., Rezende, D. J., and Eslami, S. M. A. Conditional neural processes. *CoRR*, abs/1807.01613, 2018a.
- Garnelo, M., Schwarz, J., Rosenbaum, D., Viola, F., Rezende, D. J., Eslami, S. M. A., and Teh, Y. W. Neural processes. *CoRR*, abs/1807.01622, 2018b.
- Gordon, J., Bruinsma, W. P., Foong, A. Y. K., Requeima, J., Dubois, Y., and Turner, R. E. Convolutional conditional neural processes, 2020.
- Holderrieth, P., Hutchinson, M., and Teh, Y. W. Equivariant conditional neural processes. *CoRR*, abs/2011.12916, 2020.
- Kawano, M., Kumagai, W., Sannai, A., Iwasawa, Y., and Matsuo, Y. Group equivariant conditional neural processes. *CoRR*, abs/2102.08759, 2021.
- Kim, H., Mnih, A., Schwarz, J., Garnelo, M., Eslami, S. M. A., Rosenbaum, D., Vinyals, O., and Teh, Y. W. Attentive neural processes. *CoRR*, abs/1901.05761, 2019.
- Le, T. A., Kim, H., Garnelo, M., Rosenbaum, D., Schwarz, J., and Teh, Y. W. Empirical evaluation of neural process objectives. In *NeurIPS workshop on Bayesian Deep Learning*, 2018.
- Oksendal, B. *Stochastic differential equations: an introduction with applications*. Springer Science & Business Media, 2013.
- Rasmussen, C. E. Gaussian processes in machine learning. In *Summer school on machine learning*, pp. 63–71. Springer, 2003.
- Satorras, V. G., Hoogeboom, E., and Welling, M. E (n) equivariant graph neural networks. *arXiv preprint arXiv:2102.09844*, 2021.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Póczos, B., Salakhutdinov, R., and Smola, A. J. Deep sets. *CoRR*, abs/1703.06114, 2017.

A. Consistency

Here we briefly discuss the consistency of CNPs and GNPs in the Kolmogorov and Bayesian sense. Informally, Kolmogorov’s extension theorem (KET) (Oksendal, 2013) states the following. Suppose that for every list of inputs $\mathbf{x} = (x_1, \dots, x_M)$ and corresponding outputs $\mathbf{y}_t = (y_1, \dots, y_M)$, there is a probability measure $\mu_{\mathbf{x}}$ over \mathbf{y}_t . Suppose that these laws $\{\mu_{\mathbf{x}} : \mathbf{x} \in \mathbb{R}^M, M \in \mathbb{R}\}$ are consistent under permutations

$$\begin{aligned} \mu_{x_1, \dots, x_M}(A_1 \times \dots \times A_M) &= \\ &= \mu_{x_{\sigma(1)}, \dots, x_{\sigma(M)}}(A_{\sigma(1)} \times \dots \times A_{\sigma(M)}), \end{aligned}$$

and marginalisations

$$\begin{aligned} \mu_{x_1, \dots, x_K}(A_1 \times \dots \times A_K) &= \\ &= \mu_{x_1, \dots, x_M}(A_1 \times \dots \times A_K \times \mathbb{R} \times \dots \times \mathbb{R}) \end{aligned}$$

where $1 \leq K \leq M$, $(A_i)_{i=1}^M$ are Borel measurable sets, and σ is any permutation, then there exists a stochastic process such that the finite-dimensional distributions of this process coincide with $\{\mu_{\mathbf{x}} : \mathbf{x} \in \mathbb{R}^M, M \in \mathbb{R}\}$. We can see directly from their definition in eq. (3) that CNPs are both permutation as well as marginalisation consistent. GNPs satisfy permutation consistency because by eqs. (4) to (6), a permutation σ of the target set indices, i.e. $x_{t,i}, y_{t,i} \rightarrow x_{t,\sigma(i)}, y_{t,\sigma(i)}$, leaves the RHS of eq. (4) invariant. They are also marginalisation consistent because marginalising out any of the entries of \mathbf{y}_t in eq. (4), gives the same result as querying the GNP at the same set of target points except for the marginalised ones.

However, we stress that consistency of the CNP and GNP in the sense that they satisfy KET is not the same as Bayesian consistency. In particular, the C/GNP predictive posteriors are not expected to satisfy Bayes’ rule in general

$$\begin{aligned} \pi(y^*; (\mathbf{x}_c, x^\dagger), (\mathbf{y}_c, y^\dagger), x^*) \pi(y^\dagger; \mathbf{x}_c, \mathbf{y}_c, x^\dagger) &\neq \\ \pi(y^\dagger; (\mathbf{x}_c, x^*), (\mathbf{y}_c, y^*), x^\dagger) \pi(y^*; \mathbf{x}_c, \mathbf{y}_c, x^*), \end{aligned}$$

where we have used parentheses to denote the inclusion of an additional data to the input/output context sets. Therefore, as Foong et al. (2020) have pointed out, such models do not correspond to a single consistent Bayesian model.

B. Experimental details

Each synthetic task consists of a collection of datasets sampled from the same distribution. To generate each of these datasets, we first determine the number of context and target points. We use a random number between 3 and 50 of context points and a fixed number of 50 target points. For each dataset we sample the inputs of both the context and target points, that is $\mathbf{x}_c, \mathbf{x}_t$ uniformly at random in the

region $[-2, 2]$. We then sample the corresponding outputs $\mathbf{y}_c, \mathbf{y}_t$ as follows.

Exponentiated Quadratic (EQ): We sample $\mathbf{y}_c, \mathbf{y}_t$ from a GP with an EQ covariance

$$k_{EQ}(x, x') = \sigma_v^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right),$$

with parameters $(\sigma_v^2, \ell) = (1.00, 1.00)$.

Matern 5/2: We sample $\mathbf{y}_c, \mathbf{y}_t$ from a GP with a covariance

$$k_M(x, x') = \sigma_v^2 \left(1 + \frac{r}{\ell} + \frac{r^2}{3\ell^2}\right) \exp\left(-\frac{r}{\ell}\right), \quad (12)$$

where $r = |x - x'|$, with parameters $(\sigma_v^2, \ell) = (1.00, 1.00)$.

Noisy mixture: We sample $\mathbf{y}_c, \mathbf{y}_t$ from a GP which is a sum of two EQ kernels

$$k_{NM}(x, x') = k_{EQ,1}(x, x') + k_{EQ,2}(x, x'),$$

with the following parameters $(\sigma_{v,1}^2, \ell_1) = (1.00, 1.00)$ and $(\sigma_{v,2}^2, \ell_2) = (1.00, 0.25)$.

Weakly periodic: We sample $\mathbf{y}_c, \mathbf{y}_t$ from a GP which is the product of an EQ and a periodic covariance

$$k_{WP}(x, x') = k_{EQ}(x, x') \exp\left(-\frac{2 \sin^2(\pi|x - x'|/p)}{\ell_p^2}\right),$$

with EQ parameters $(\sigma_{v,EQ}^2, \ell_p) = (1.00, 1.00)$ and periodic parameters $(p, \ell_{EQ}) = (0.25, 1.00)$.

Lastly, for all tasks we add iid Gaussian noise with zero mean and variance $\sigma_n^2 = 0.05^2$. This noise level was not given to the models, which in every case learned a noise level from the data.

The models were trained for 100 epochs, each consisting of 1024 iterations, at each of which 16 datasets were presented as a minibatch to the models. All models were trained with the Adam optimiser, with a learning rate of 5×10^{-4} .