
NP-PROV: Neural Processes with Position-Relevant-Only Variances

Xuesong Wang

University of New South Wales
xuesong.wang1@student.unsw.edu.au

Lina Yao

University of New South Wales
lina.yao@unsw.edu.au

Xianzhi Wang

University of Technology Sydney
xianzhi.wang@uts.edu.au

Feiping Nie

Northwestern Polytechnical University
feipingnie@gmail.com

Abstract

Neural Processes (NPs) families encode distributions over functions to a latent representation, given context data, and decode posterior mean and variance at unknown locations. Since mean and variance are derived from the same latent space, they may fail on out-of-domain tasks where fluctuations in function values amplify the model uncertainty. We present a new member named Neural Processes with Position-Relevant-Only Variances (NP-PROV). NP-PROV hypothesizes that a target point close to a context point has small uncertainty, regardless of the function value at that position. The resulting approach derives mean and variance from a function-value-related space and a position-related-only latent space separately. Our evaluation on synthetic and real-world datasets reveals that NP-PROV can achieve state-of-the-art likelihood while retaining a bounded variance when drifts exist in the function value.

1 Introduction

Neural networks (NNs) are proven effective in various machine learning tasks for making deterministic predictions. They have the flexibility of describing and fitting any type of data distributions. Nevertheless, most neural networks are incapable of evaluating model stochasticity. They cannot handle tasks where prediction uncertainty is equally crucial, e.g., autonomous driving Wei et al. (2011), disease diagnosis Lorenzi et al. (2019), and stock market forecasting Christou et al. (2017). In contrast, Gaussian processes (GPs) model data with an infinite sequence of correlated normal distributions and are intrinsically suitable for such problems. Conditioned on some prior knowledge, e.g., context points with positions \mathbf{X} and function values \mathbf{Y} , they are able to infer the likelihood of target values \mathbf{Y}_* at unknown locations \mathbf{X}_* . Despite the advantages, GPs require designing data-specific kernel functions; the cubic computational cost of matrix inversion impedes GPs from handling large-scale data.

The recent progress in models based on Neural Process (NP) Garnelo et al. (2018b) have brought GPs the advantages of fast forward propagation and powerful feature representations of NNs. Basic NPs represent a stochastic process with an encoder-decoder network structure. It encodes context data (\mathbf{X}, \mathbf{Y}) to a latent representation $\mathbf{Z} \sim \mathbf{P}(\mathbf{Z}|\mathbf{X}, \mathbf{Y})$ and decodes it to a posterior probability of the target data based on its relationship with the context points $\mathbf{Y} \sim \mathbf{P}(\mathbf{Y}_*|\mathbf{Z}; \mathbf{X}_*)$. Recent years have witnessed a family of NP-variants, such as Attentive NP Kim et al. (2019), Convolutional Conditional NP Gordon et al. (2019), and Sequential NP Singh et al. (2019). They improve NP via aggregating context knowledge non-linearly and considering spatial or temporal relationships among target points but still decode mean and variance from the same latent variable—meaning that their variances are

correlated to \mathbf{Y} . When future testing sets have shifts compared to training sets (e.g., stock market price with incremental trends), fluctuations in \mathbf{Y} can severely amplify model uncertainty. A model can be stable on out-of-domain tasks only when the variance inference is relevant to locations \mathbf{X} yet irrelevant to function values \mathbf{Y} .

In this paper, we introduce a new member named Neural Processes with Position-Relevant-Only Variances (NP-PROV), which derives mean and variance functions from two coupled latent spaces. Mean values are related to function values \mathbf{Y} , self-correlation within context locations \mathbf{X} , and cross-correlations between context positions \mathbf{X} and target positions \mathbf{X}_* . Variance values exclude function values yet are also associated with the self-correlations within the target positions. Our main contributions are as follows:

- We designed an auto-encoder module associated with self-correlations between data points. It reconstructs context data through the module and reduces model uncertainty in scenarios where context points have higher self-correlations.
- The approach derives mean and variance values from two coupled latent spaces. The position-relevant-only variances prevent the model from estimating oscillating uncertainty when function values are out of the training range.
- The proposed model achieves state-of-the-art performance over on-the-grid and off-the-grid datasets. Specifically, three GP-based kernels and a real-world time series are adopted to sample off-the-grid tasks. Four image inpainting tasks are evaluated: MNIST, SVHN, celebA and miniImageNet.

2 Background

Gaussian Processes. Let $\mathbf{Y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$ be a set of N observations at input locations $\mathbf{X} = [x_1, \dots, x_N]^\top \in \mathbb{R}^{N \times D}$, a function $f: \mathbb{R}^D \mapsto \mathbb{R}$, and a likelihood $p(\mathbf{Y} | \mathbf{F}; \mathbf{X})$, where $\mathbf{F} = f(\mathbf{X})$ denotes the function values at the input locations. A Gaussian process (GP) prior is placed on the function f ; it models all function values as a jointly Gaussian distribution to infer f . The prior has a mean function $m(\cdot): \mathbb{R}^D \mapsto \mathbb{R}$ and a covariance function $\mathcal{K}(\cdot, \cdot): \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}$. The generative model of the corresponding GP is as follows:

$$p(\mathbf{Y} | \mathbf{F}; \mathbf{X}) = \mathcal{N}(m(\mathbf{X}), \mathcal{K}(\mathbf{X}, \mathbf{X}^\top)) \quad (1)$$

GP hypothesizes similar function values of two inputs that are highly correlated. Given some context samples (\mathbf{X}, \mathbf{Y}) , we can perform probabilistic inference and assign posterior distributions over unknown locations from the same function. For new input locations $\mathbf{X}_* = [x_{*1}, \dots, x_{*M}]^\top$, the posterior distribution also follows a joint Gaussian distribution $p(\mathbf{Y}_* | \mathbf{F}; \mathbf{X}_*, \mathbf{X}, \mathbf{Y}) = \mathcal{N}(\mu_*, \Sigma_*)$:

$$\begin{aligned} \mu_* &= m(\mathbf{X}_*) + \mathbf{K}_*^\top \mathbf{K}^{-1} (\mathbf{Y} - m(\mathbf{X})) \\ \Sigma_* &= \mathbf{K}_{**} - \mathbf{K}_*^\top \mathbf{K}^{-1} \mathbf{K}_* \end{aligned} \quad (2)$$

where $\mathbf{K} = \mathcal{K}(\mathbf{X}, \mathbf{X}^\top)$, $\mathbf{K}_* = \mathcal{K}(\mathbf{X}, \mathbf{X}_*^\top)$, and $\mathbf{K}_{**} = \mathcal{K}(\mathbf{X}_*, \mathbf{X}_*^\top)$. Intuitively, when there is no trend in observations (i.e., $m(\mathbf{X}_*) = m(\mathbf{X}) = 0$), the mean of \mathbf{Y}_* is a linear combination of elements in \mathbf{Y} . The weights of those elements are associated with the self-correlation of \mathbf{X} and cross-correlation between \mathbf{X} and \mathbf{X}_* . The variance equals the total uncertainty of \mathbf{X}_* minus the certainty induced from \mathbf{X} and is irrelevant to function values \mathbf{Y} .

Convolutional Conditional Neural Processes. As a type of latent variable model, NPs inherit distribution consistency from GPs, where context data (\mathbf{X}, \mathbf{Y}) and target data $(\mathbf{X}_*, \mathbf{Y}_*)$ are sampled from the same function and share latent variables.

Convolutional Conditional Neural Processes (ConvCNP) differ from other NP families in handling testing data outside the range of training data. They introduce a discretization of a continuous range of $\mathbf{X} \cup \mathbf{X}_*$ named \mathbf{X}_t . Convolution operations enable NPs to focus on local receptive fields

of gridded inputs and to generalize to out-of-domain predictions. The posterior distribution of $p(\mathbf{Y}_*|\mathbf{X}_*, \mathbf{X}_t, \mathbf{X}, \mathbf{Y}) = \mathcal{N}(\mu_*, \Sigma_*)$ is also in the form of a Gaussian distribution:

$$\begin{aligned}\mu_* &= \psi_{*\mu}\{\mathbf{K}_*^\top \text{CNN}[\psi_t(\mathbf{K}_t^\top \mathbf{Y})]\} \\ \Sigma_* &= \psi_{*\Sigma}\{\mathbf{K}_*^\top \text{CNN}[\psi_t(\mathbf{K}_t^\top \mathbf{Y})]\}\end{aligned}\quad (3)$$

where $\mathbf{K}_* = \mathcal{K}_*(\mathbf{X}_*, \mathbf{X}_t^\top)$, $\mathbf{K}_t = \mathcal{K}_t(\mathbf{X}_t, \mathbf{X}^\top)$ are covariance functions with learnable length scales; ψ_t , $\psi_{*\mu}$ and $\psi_{*\Sigma}$ are positive-definite transformations associated with a Reproducing Kernel Hilbert Space; the graphical model of ConvCNP is described in Fig 1 (a); $\psi_t(\mathbf{K}_t^\top \mathbf{Y})$ encodes prior knowledge about the function distribution on the entire grid \mathbf{X}_t ; $\text{CNN}(\cdot)$ enables model to exhibit translation invariance; $\psi_{*\mu}(\mathbf{K}_*^\top \cdot)$ and $\psi_{*\Sigma}(\mathbf{K}_*^\top \cdot)$ translate the knowledge to predict target mean and standard deviation for $p(\mathbf{Y}_*)$.

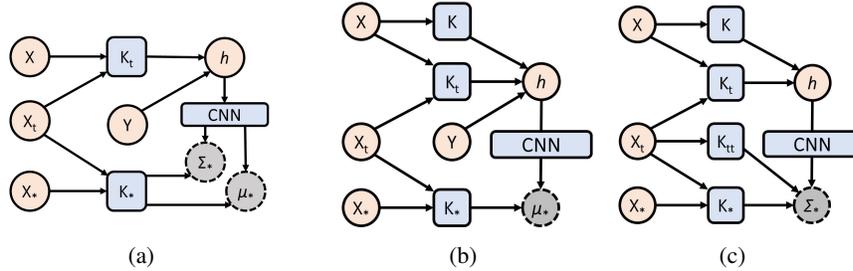


Figure 1: Graphical Models for (a) ConvCNP (b) NP-PROV mean and (c) NP-PROV variance

3 Neural Processes with Position-Relevant-Only Variances

\mathbf{K}_t evaluates correlation between \mathbf{X}_t and \mathbf{X} , and \mathbf{K}_* evaluates cross correlation between \mathbf{X}_* and \mathbf{X}_t . Therefore, $\psi_t(\mathbf{K}_t^\top \mathbf{Y})$ and $\psi_{*}(\mathbf{K}_*^\top \cdot)$ in ConvCNP act as two cascaded cross-correlation modules in GPs. We propose a new member Neural Processes with Position-Relevant-Only Variances (NP-PROV) to improve ConvCNP in two ways:

1. Considering self-correlations within context data. ConvCNP usually considers cross-correlations between target data and context data only. It is possible to further reduce the model uncertainty near the region where context data have high self-correlations.
2. Changing variance inference to position-relevant-only. The variance derivation is related to function values in ConvCNP. It can effectively prevent the model from amplifying uncertainty caused by value fluctuations if it is made to be only associated with relative distances among locations (as in GP).

In the following, we detail in an off-the-grid scenario with continuous inputs and an on-the-grid scenario with intrinsically discretized inputs.

3.1 Off-the-grid Scenario.

Mean function We follow convCNP to construct the model with two cascaded GP-alike layers. The first layer maps \mathbf{X} to a uniformly discretized grid space $\mathbf{X}_t = [x_1, \dots, x_t]^\top$ built on the lower and upper range of $\mathbf{X} \cup \mathbf{X}_*$. The second layer maps the space back to \mathbf{X}_* . As illustrated in Fig 1(b), an auto-encoder structure is computed to match $\mathbf{K}^{-1}\mathbf{Y}$ first:

$$\begin{aligned}h_{ni}^{self} &= \psi_E(\mathbf{K}_{ni}y_n) \\ \tilde{y}_i &= \psi_D\left(\sum_{n=1}^N h_{ni}^{self} \mathbf{K}_{ni}\right), \quad \mathcal{L}_2 = \frac{1}{N} \sum_{i=1}^N (y_i - \tilde{y}_i)^2\end{aligned}\quad (4)$$

where h_{ni}^{self} is self-correlation mapping from the n -th context point to the i -th context point; \mathbf{K}_{ni} represents corresponding self-covariance elements in $\mathbf{K} = \mathcal{K}(\mathbf{X}, \mathbf{X}^\top) \in \mathbb{R}^{n \times n}$; ψ_E and ψ_D are

positive-definite linear transformations. Since $\mathbf{K}\mathbf{K}^{-1}\mathbf{Y} = \mathbf{Y}$, we use auto-encoder structure to map h_{ni}^{self} back to y_i and minimize reconstruction loss \mathcal{L}_2 to mimic matrix inversion. Then, we can compute cross-correlation weights associated with $\mathbf{K}^t = \mathcal{K}(\mathbf{X}_t, \mathbf{X}^\top) \in \mathbb{R}^{t \times n}$:

$$h_j^{cross} = \sum_{n=1}^N [1, y_n]^T \mathbf{K}_{nj}^t, \quad h_j^{cross(1)} = h_j^{cross(1)} / h_j^{cross(0)} \quad (5)$$

$$h_j = \psi_t \left(\left[\frac{1}{N^2} \sum_{i=1}^N \sum_{n=1}^N h_{ni}^{self}, \quad h_j^{cross} \right] \right)$$

where h_j^{cross} maps the n -th context point to the j -th data in \mathbf{X}_t . $\mathbf{I} = \mathbf{Y}^0$ is an additional identity density channel appended to $\mathbf{Y} \leftarrow [\mathbf{I} \mathbf{Y}]$ and is the powers of \mathbf{Y} up to order-1 for normalizing large variance. Dividing $h_i^{(0)}$ helps to normalize scales when there are large variations in \mathbf{X} . A new linear transformation ψ_t fuses averaged self-correlation and cross-correlation information to a latent variable h_j on the j -th data of \mathbf{X}_t .

A CNN takes the obtained condition h_j to suffice translation invariance. It uses a UNet structure with skipped concatenation of 6 convolution layers and 6 transposed convolution layers to capture global and local patterns. The CNN enables NP-PROV to predict out-of-domain tasks by handling locations outside the training range through filter sliding. When projecting the output back from \mathbf{X}_t to \mathbf{X}_* , ψ_* transforms the cross-correlation weights associated with $\mathbf{K}^* = \mathcal{K}(\mathbf{X}_*, \mathbf{X}_t^\top) \in \mathbb{R}^{m \times t}$. The overall mean function of the m -th target point in \mathbf{X}_* is as follows:

$$\mu_m^* = \psi_\mu^* \left(\sum_{j=1}^t \text{CNN}(h_j) \mathbf{K}_{jm}^* \right) \quad (6)$$

Covariance function Based on Eq 2, we replace original values \mathbf{Y} with \mathbf{K}^t in covariance functions and add a module associated with self-correlations within the grided data \mathbf{X}_t . As in Fig 1 (c), we define two new mappings: ψ and ψ_{tt} :

$$h_i = \psi \left(\sum_{j=1}^t \mathbf{K}_{ji}^{t^\top} \right), \quad h_j^{self} = \psi_{tt} \left(\sum_{k=1}^t \mathbf{K}_{kj}^{tt} \right) \quad (7)$$

where $\mathbf{K}^{t^\top} = \mathcal{K}(\mathbf{X}, \mathbf{X}_t^\top) \in \mathbb{R}^{n \times t}$; $\mathbf{K}^{tt} = \mathcal{K}(\mathbf{X}_t, \mathbf{X}_t^\top) \in \mathbb{R}^{t \times t}$. h_i maps \mathbf{X}_t to \mathbf{X} and replaces values \mathbf{Y} in the mean function. This value-irrelevant input will generate a new h_j . After the convolution on the new h_j , the result is concatenated with h_j^{self} . The overall covariance function is as follows:

$$\Sigma_m^* = \psi_\Sigma^* \left(\sum_{j=1}^t [\text{CNN}(h_j), \quad h_j^{self}] \mathbf{K}_{jm}^* \right) \quad (8)$$

Again, as in Fig1(c), the added elements in NP-PROV can retain reasonable variance regardless of drifts in function values when compared with ConvCNP.

Inference and Learning With μ_* and Σ_* , the posterior distribution of \mathbf{Y}_* follows a multivariate normal distribution: $\mathcal{N}(\mu_*, \Sigma_*)$. The training objective is to maximize the log-likelihood of target outputs meanwhile minimizing the reconstruction loss, given all parameters defined as Θ :

$$\theta^* = \arg \max_{\theta \in \Theta} \sum_{m=1}^M \log p(y_m | \mathcal{N}(\mu_*, \Sigma_*)) - \frac{1}{N} \sum_{i=1}^N (y_i - \tilde{y}_i)^2 \quad (9)$$

3.2 On-the-grid Scenario.

CNN can be applied to On-the-grid data (e.g., images), which are discretized and represent better spatial correlations between context and target data with little efforts. We describe the context inputs \mathbf{X} using a context mask \mathbf{M} , where an element value is 1 if a pixel is revealed and 0 otherwise. The target inputs \mathbf{X}_* formulate the whole image.

Mean function Similar to off-the grid data, self-correlations and cross-correlations weights are computed for the latent variable extraction. Instead of adopting kernel functions, convolutions are implemented on the context masks \mathbf{M} and the context values $\mathbf{M} \odot \mathbf{Y}$:

$$\begin{aligned} h_{self} &= \psi_E(\mathbf{M}), \quad \tilde{M} = \psi_D(h_{self}), \quad \mathcal{L}_2 = \|\mathbf{M} - \tilde{M}\|_2 \\ h_{cross} &= \psi([\mathbf{M}, \mathbf{M} \odot \mathbf{Y}]) \end{aligned} \quad (10)$$

where \mathbf{Y} represents a full image, ψ_E, ψ_D, ψ are 1-layer convolution operations. ψ_E aims to extract self-correlations within unmasked points. ψ shows cross correlations within unmasked pixel values. Then, a UNet CNN structure maps the concatenated latent variable to a translation invariant representation. Finally, a new convolution $\psi_{*\mu}$ maps the representation to the posterior mean:

$$\mu_* = \psi_{*\mu}(\text{CNN}([h_{self} \quad h_{cross}])) \quad (11)$$

Covariance function We substitute the masked image $\mathbf{M} \odot \mathbf{Y}$ with the mask \mathbf{M} to extract a output-irrelevant latent variable h for variance functions:

$$h_{cross} = \psi([\mathbf{M}, \mathbf{M}]) \quad h = [h_{self} \quad h_{cross}] \quad (12)$$

Also, we add a self-correlation layer on the target masks, i.e., an identity matrix \mathbf{I} :

$$h_{**} = \psi_{**}(\mathbf{I}) \quad (13)$$

The overall function learns a new transformation $\psi_{*\Sigma}$ that maps h and h_{**} to the posterior covariance. The objective is to maximize a log-likelihood to recover the whole image and to minimize the reconstruction error of the mask simultaneously.

$$\Sigma_* = \psi_{*\Sigma}([\text{CNN}(h) \quad h_{**}]) \quad (14)$$

4 Experiments

We conduct few-shot regressions tasks on off-the-grid 1d datasets and on-the-grid images to evaluate the effectiveness of NP-PROV.

4.1 Off-the-grid datasets

We aim to maximize the likelihood of the outputs \mathbf{Y}_* at unknown locations \mathbf{X}_* , given observations \mathbf{Y} at input locations \mathbf{X} . We are interested in the following questions: (a) Are the prediction and uncertainty estimation reasonable? (b) How will the model perform when the testing range of \mathbf{X} and \mathbf{Y} goes beyond the training data, i.e., out-of-domain tasks? (c) How will the self-correlation affect model prediction? We compare the results of four state-of-the-art neural process members: Neural Process (NP) Garnelo et al. (2018b), Conditional Neural Process (CNP) Garnelo et al. (2018a), Attentive Neural Process (ANP) Kim et al. (2019), and Convolutional Conditional Neural Process(ConvCNP)Gordon et al. (2019). We use three challenging kernels to generate synthesised Gaussian processes functions, according to Gordon et al. (2019):

- EQ : $\mathcal{K}(x, x') = e^{-\frac{1}{2}(\frac{x-x'}{0.25})^2}$
- Matern $-\frac{5}{2}$: $\mathcal{K}(x, x') = (1 + 4\sqrt{5}d + \frac{5}{3}d^2)e^{-\sqrt{5}d}$ with $d = 4|x - x'|$
- Weakly periodic: $\mathcal{K}(x, x') = e^{-\frac{1}{2}(f_1(x) - f_1(x'))^2 - \frac{1}{2}(f_2(x) - f_2(x'))^2} \cdot e^{-\frac{1}{8}(x-x')^2}$, with $f_1(x) = \cos(8\pi x)$ and $f_2(x) = \sin(8\pi x)$

The training data range within $x \in [-2, 2]$ for GP-sampled datasets. Also, a real-world time series dataset Smart Meter 3springs (2019) is added. It contains energy consumption readings from a sample of 5,567 London Households between November 2011 and February 2014. We select timestamp in days as the input \mathbf{X} and consumption in kWh/ half-hour as the output \mathbf{Y} . The x range is set to $[0, 2]$, representing a relative 2-day window from a random clip on the time axis. The number of context points and target points in each task are uniformly distributed in $\mathcal{U}(3, 50)$. A batch of 16 tasks is

generated and the total number of tasks in one epoch is 256. We train each model for 200 epochs and then test them using 2,048 new generated tasks for 6 times.

Table 1 shows the log-likelihood (on the probability density function) of five methods. Model performance is presented in Fig 2, which supplement mean and variance results from the original GP as a reference for synthesized datasets. Table 1 shows both ConvCNP and NP-PROV outperform others significantly. The result of NP is quite unstable, due to stochastic encoding of the latent variable. The first column of Fig 2 shows the results with the testing range $[-5, 5]$ for the GP-sampled datasets. NP, CNP, and ANP predict oscillated mean values; therefore, we omit their variance values for display clarity. Both NP-PROV and ConvCNP match well with the original GP. This advantage sources from convolution, where filters can slide along the axis. NP-PROV usually predicts a narrower variance than ConvCNP when the target points are adjacent to context points. This might be attributed to more uncertainty reduction from context self-correlation. When Smart Meter is extended to the interval $[-1, 5]$, the predicted variance becomes much higher than NP-PROV (as it only adopts cross-correlation), and the target points become far from the context points. NP-PROV mitigates this issue by leveraging the self-correlation on target data.

Table 1: Log-likelihood of off-the-grid datasets (mean \pm standard deviation)

Model	EQ	Matern	Weakly Periodic	Smart Meter
NP	-1.20 ± 0.43	-0.82 ± 1.95	-1.75 ± 1.36	0.93 ± 0.23
CNP	-1.10 ± 0.02	-1.36 ± 0.03	-2.04 ± 0.02	1.81 ± 0.06
ANP	-0.35 ± 0.01	-0.75 ± 0.02	-2.09 ± 0.03	1.66 ± 0.05
ConvCNP	2.15 ± 0.05	0.83 ± 0.06	-1.15 ± 0.01	2.65 ± 0.06
NP-PROV	2.20 ± 0.02	0.90 ± 0.03	-1.00 ± 0.02	2.32 ± 0.05
ConvCNP (self)	77.60 ± 2.33	40.26 ± 0.81	-47.63 ± 0.75	0.95 ± 0.00
NP-PROV (self)	77.55 ± 2.61	44.14 ± 1.03	-40.96 ± 0.89	0.99 ± 0.01

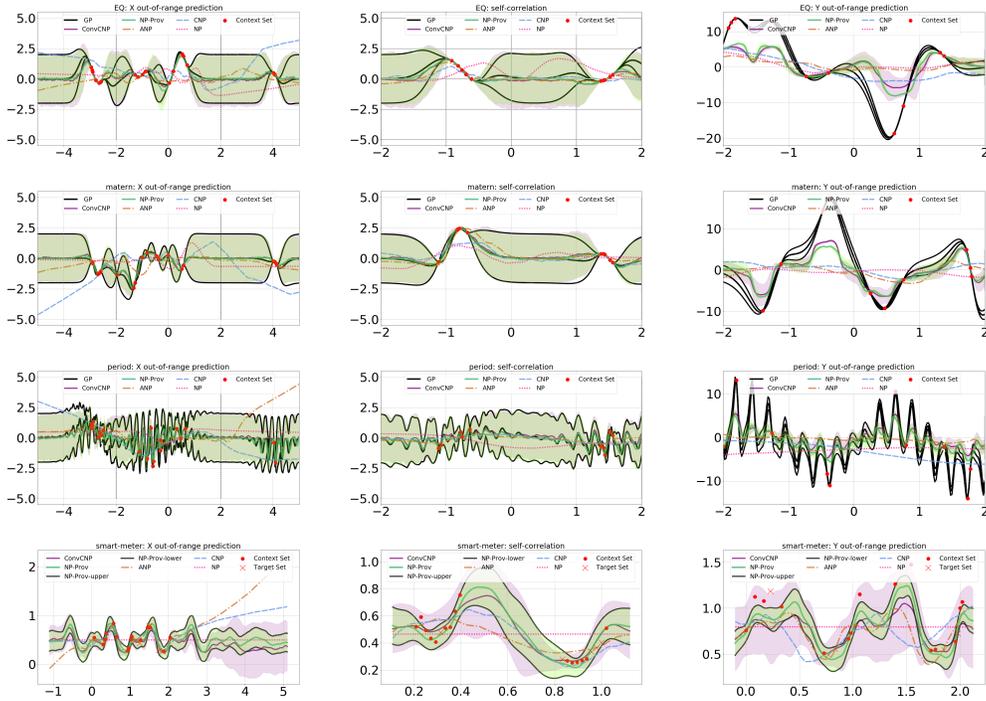


Figure 2: Model predictions on 4 datasets: EQ (the first row), Matern (the second row), Weakly period (the third row) and Smart meter (the fourth row). The first column is out-of-x-range prediction, the second column is self-correlation prediction, the third column is out-of-y-range prediction.

We analyze the impact of self-correlation on the model. In the second column of Fig2, we give context data in smaller intervals of $[-1, -0.5] \cup [1.2, 1.7]$ for GP-sampled datasets and $[0.2, 0.4] \cup [0.8, 1.1]$ for Smart Meter to get higher self-correlations between context points. Also, in Table 1, we present the log-likelihood of ConvCNP (self) and NP-PROV (self) where target points are adjacent to these compacted regions of context data. NP-PROV predicted lower variances on target points near compacted regions, indicating the ability to capture the self-correlation between context points.

Finally, we modify the testing GP generator to $y = 10 \times \mathcal{GP}(x)$ and 1.5 times of the original Smart Meter to evaluate model performance when testing \mathbf{Y} goes beyond the training range. The third column of Fig 2 reveals that all methods are too conservative in adaptation due to the normalization effect of activation layers. However, the variance of NP-PROV on the context data is close to zero, and it is not affected by the explosion of target values, which suffice the hypothesis of a stochastic process. Such an advantage is crucial when there are shifts in time series data.

4.2 On-the-grid datasets

Given different proportions of pixel values, the objective is to complete the whole image and estimate uncertainty at unknown pixels. Four image datasets are introduced: MNIST, SVHN, celebA 32×32 and miniImageNet Deleu et al. (2019). Except MNIST, all the other datasets are RBG images. MiniImageNet is used to verify the model capability of recovering images of unseen classes. The context points are sampled from $\mathcal{U}(\frac{n_{total}}{100}, \frac{n_{total}}{2})$ and the target points are n_{total} , i.e., the whole image.

Table 2 demonstrates the log-likelihood of the comparing methods on testing sets. The results of miniImageNet are only displayed for ConvCNP and NP-PROV since other baselines are overwhelmed by the image size 84×84 . Fig 3 shows the prediction results on the datasets. NP-PROV and miniImageNet massively outperform other baselines. NP-PROV achieves better results on SVHN and miniImageNet. From the variance results in Fig 3, most of the methods don't exhibit explainable variances in these two methods, because they are similar to "unsupervised" scenarios where one or a few images represent a new class, hence the methods are unable to gain enough variance information from the pixel values. Whereas there are only limited patterns for a digit or an outline of a face in MNIST and celebA, other baselines adopting pixel values can learn a stabilized mean and display variance on the edges.

Table 2: Log-likelihood of on-the-grid datasets (mean \pm standard deviation)

Model	MNIST	SVHN	celebA 32×32	miniImageNet
NP	$0.65 \pm 4e-4$	$3.21 \pm 6e-4$	$2.78 \pm 1e-3$	–
CNP	$1.94 \pm 4e-2$	$4.48 \pm 5e-3$	$3.09 \pm 4e-2$	–
ANP	$0.95 \pm 3e-3$	$3.50 \pm 5e-3$	$2.32 \pm 4e-2$	–
ConvCNP	$2.98 \pm 4e-2$	$6.03 \pm 2e-1$	$6.35 \pm 2e-1$	$3.65 \pm 4e-2$
NP-PROV	$2.66 \pm 3e-2$	$8.24 \pm 5e-2$	$5.11 \pm 1e-2$	$4.39 \pm 2e-1$

5 Related Work

Since Gaussian processes suffer high computational cost on matrix inverse, recent work focuses on adopting fast forward-feeding neural networks to substitute original GP. A group of Neural Processes-based models is proposed based on variational inference and ELBO (Evidence Lower Bound) Rudner et al. (2018). Neural process families abstract a latent variable or a function from context data and decode the function to the target data based on the relationships between target data and context data. The variations of NPs mostly depend on how the relationships are presented. The original NP Garnelo et al. (2018b) and conditional NPGarnelo et al. (2018a) adopt mean function to extract stochastic and deterministic latent variables that suffice the exchangeability of a stochastic process. Attentive NP Kim et al. (2019) considers the self-attention between context points and cross-attention between context points and target points so that attentive aggregation towards latent variables can be used. Regarding nonstructural relationships between data points, Sequential NP Singh et al. (2019) and recurrent NPWilli et al. (2019) address the temporal correlations in time series. Convolutional conditional NP Gordon et al. (2019) utilizes the translation equivariance of convolution to predict on

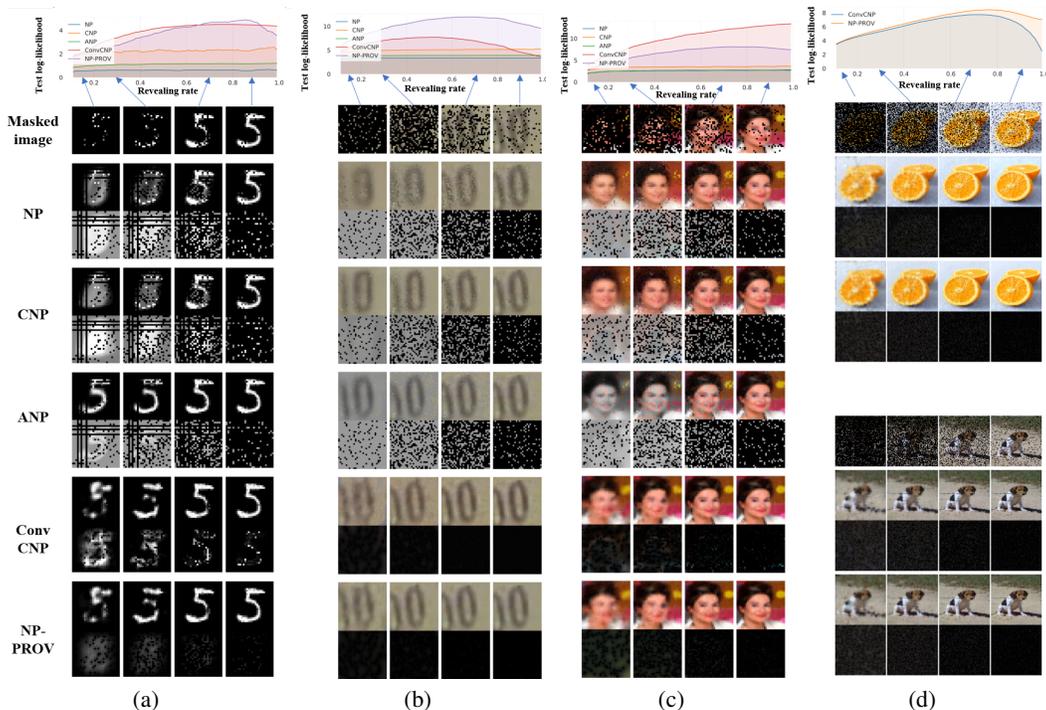


Figure 3: Test log-likelihood w.r.t revealing rates (the upper row). The lower row presents some samples from (a) MNIST, (b) SVHN, and (c) celeBA $\times 32$ and miniImageNet (d) with different revealing rates: 10%, 30%, 70%, 90%. The first row in every method predicts mean and the second row predicts variance.

out-of-training range target locations. Functional NP Louizos et al. (2019) and Graph NP Carr and Wingate (2019) exploit topological relationships between context and target nodes. Similar to Meta-Agnostic-Meta-Learning(MAML), residual NPLee et al. (2020) assumes that few shot tasks share a unified latent variable with minor variations on each task. Thus, the latent variable of the context data is adapted based on the prediction residues. Besides, numerous work focuses on enhancing Gaussian Processes with non-linear feature extraction using neural networks. NNGP Lee et al. (2017) and ConvNet Garriga-Alonso et al. (2018) explore the relationships between Gaussian processes and one layer neural networks from the perspective of Bayesian inference and approximation. Deep Gaussian processes Damianou and Lawrence (2013) and deep kernel learning Wilson et al. (2016) substitute manual designed kernels with neural networks to extract higher-dimensional features.

6 Conclusions and Discussions

We introduced NP-PROV, a new member of Neural Processes that derive variances from a position-relevant-only latent space. We verify that NP-PROV can estimate bounded uncertainty when context data has high self-correlations or function values are out-of-the-training range. We mitigate the problem of predicting stabilized variance under shifts in function values. We believe that NP-PROV opens a door of rethinking the relationships between the mean and variance in Neural Processes. This work leaves multiple avenues for future improvements. It would be interesting to see the mean derivation be more adaptive to out-of-the training range. Also, unifying on-and-off-the-grid version of NP-PROV to fit in higher-dimensional space, e.g., by designing a hyper grid for convolution.

Broader Impact

The proposed neural processes have implications for both researchers and practitioners who rely on dynamic data or rapidly shifting contexts to make reliable predictions. This work enables a learning

model to adapt to new samples in testing environments and, therefore, overcomes the limitations or biases in pre-prepared training data. This work fundamentally supports better modeling that improves all relevant applications that are based on neural processes. It has the potential to mitigate the challenge of model uncertainty posed by the lack of training data, overfitting, and imbalanced training and testing data distribution in traditional machine learning research. Potential applications benefiting from this technology include super-resolution image reconstruction, fine-grained map interpolation in order to save costs from densely distributed sensors, and medical data generation using only a few shots. One may also argue that the related research can be used for military computer vision tasks, high precision GPS search, and drone state estimation under dynamic environments.

References

- 3springs (2019). Neural processes for sequential data. Available at: <https://github.com/3springs/attentive-neural-processes>.
- Carr, A. N. and Wingate, D. (2019). Graph neural processes: Towards bayesian graph neural networks. *arXiv preprint arXiv:1902.10042*.
- Christou, C., Cunado, J., Gupta, R., and Hassapis, C. (2017). Economic policy uncertainty and stock market returns in pacificrim countries: Evidence based on a bayesian panel var model. *Journal of Multinational Financial Management*, 40:92–102.
- Damianou, A. and Lawrence, N. (2013). Deep gaussian processes. In *Artificial Intelligence and Statistics*, pages 207–215.
- Deleu, T., Würfl, T., Samiei, M., Cohen, J. P., and Bengio, Y. (2019). Torchmeta: A Meta-Learning library for PyTorch. Available at: <https://github.com/tristandeleu/pytorch-meta>.
- Garnelo, M., Rosenbaum, D., Maddison, C. J., Ramalho, T., Saxton, D., Shanahan, M., Teh, Y. W., Rezende, D. J., and Eslami, S. (2018a). Conditional neural processes. *arXiv preprint arXiv:1807.01613*.
- Garnelo, M., Schwarz, J., Rosenbaum, D., Viola, F., Rezende, D. J., Eslami, S., and Teh, Y. W. (2018b). Neural processes. *arXiv preprint arXiv:1807.01622*.
- Garriga-Alonso, A., Rasmussen, C. E., and Aitchison, L. (2018). Deep convolutional networks as shallow gaussian processes. *arXiv preprint arXiv:1808.05587*.
- Gordon, J., Bruinsma, W. P., Foong, A. Y., Requeima, J., Dubois, Y., and Turner, R. E. (2019). Convolutional conditional neural processes. *arXiv preprint arXiv:1910.13556*.
- Kim, H., Mnih, A., Schwarz, J., Garnelo, M., Eslami, A., Rosenbaum, D., Vinyals, O., and Teh, Y. W. (2019). Attentive neural processes. *arXiv preprint arXiv:1901.05761*.
- Lee, B.-J., Hong, S., and Kim, K. (2020). Residual neural processes. In *AAAI 2020*.
- Lee, J., Bahri, Y., Novak, R., Schoenholz, S. S., Pennington, J., and Sohl-Dickstein, J. (2017). Deep neural networks as gaussian processes. *arXiv preprint arXiv:1711.00165*.
- Lorenzi, M., Filippone, M., Frisoni, G. B., Alexander, D. C., Ourselin, S., Initiative, A. D. N., et al. (2019). Probabilistic disease progression modeling to characterize diagnostic uncertainty: application to staging and prediction in alzheimer’s disease. *NeuroImage*, 190:56–68.
- Louizos, C., Shi, X., Schutte, K., and Welling, M. (2019). The functional neural process. In *Advances in Neural Information Processing Systems*, pages 8743–8754.
- Rudner, T. G., Fortuin, V., Teh, Y. W., and Gal, Y. (2018). On the connection between neural processes and gaussian processes with deep kernels. In *Workshop on Bayesian Deep Learning, NeurIPS*.
- Singh, G., Yoon, J., Son, Y., and Ahn, S. (2019). Sequential neural processes. In *Advances in Neural Information Processing Systems*, pages 10254–10264.

- Wei, J., Dolan, J. M., Snider, J. M., and Litkouhi, B. (2011). A point-based mdp for robust single-lane autonomous driving behavior under uncertainties. In *2011 IEEE International Conference on Robotics and Automation*, pages 2586–2592.
- Willi, T., Masci, J., Schmidhuber, J., and Osendorfer, C. (2019). Recurrent neural processes. *arXiv preprint arXiv:1906.05915*.
- Wilson, A. G., Hu, Z., Salakhutdinov, R., and Xing, E. P. (2016). Deep kernel learning. In *Artificial Intelligence and Statistics*, pages 370–378.