

---

# Neural Processes with Stability

---

**Huafeng Liu, Liping Jing, Jian Yu**

Beijing Key Lab of Traffic Data Analysis and Mining, Beijing Jiaotong University  
The School of Computer and Information Technology, Beijing Jiaotong University  
{hfliu1, lpjing, jianyu}@bjtu.edu.cn

## Abstract

Unlike traditional statistical models depending on hand-specified priors, neural processes (NPs) have recently emerged as a class of powerful neural statistical models that combine the strengths of neural networks and stochastic processes. NPs can define a flexible class of stochastic processes well suited for highly non-trivial functions by encoding contextual knowledge into the function space. However, noisy context points introduce challenges to the algorithmic stability that small changes in training data may significantly change the models and yield lower generalization performance. In this paper, we provide theoretical guidelines for deriving stable solutions with high generalization by introducing the notion of algorithmic stability into NPs, which can be flexible to work with various NPs and achieves less biased approximation with theoretical guarantees. To illustrate the superiority of the proposed model, we perform experiments on both synthetic and real-world data, and the results demonstrate that our approach not only helps to achieve more accurate performance but also improves model robustness.

## 1 Introduction

Neural processes (NPs) [9, 10] constitute a family of variational approximation models for stochastic processes with promising properties in computational efficiency and uncertainty quantification. Different from traditional statistical modeling for which a user typically hand-specifies a prior (e.g., smoothness of functions quantified by a Gaussian distribution in Gaussian process [25]), NPs implicitly define a broad class of stochastic processes with neural networks in a data-driven manner. When appropriately trained, NPs can define a flexible class of stochastic processes well suited for highly non-trivial functions that are not easily represented by existing stochastic processes.

NPs meta-learn a distribution over predictors and provide a way to select an inductive bias from data to adapt quickly to a new task. Incorporating the data prior into the model as an inductive bias, NPs can reduce the model complexity and improve model generalization. Usually, an NP predictor is described as predicting a set of data (target set) given a set of labeled data (context set). However, the number of noise in data introduces challenges to the algorithmic stability. In NPs, models are biased to the meta-datasets (a dataset of datasets), so small changes in the dataset (noisy or missing) may significantly change the models. As demonstrated in previous work [9, 10, 14, 11], existing NPs cannot provide stable predictions under noisy conditions, which may introduce high training error variance, and minimizing the training error may not guarantee consistent error reduction on the test set, i.e. low generalization performance [2]. In this case, algorithmic stability and generalization performance have strong connections and an unstable NP model has low generalization performance.

A stable model is one for which the learned solution does not change much with small changes in training set [2]. In general, heuristic techniques, such as cross-validation and ensemble learning, can be adopted to improve the generalization performance. Cross-validation needs to sacrifice the limited training data, while ensemble learning is computationally expensive on training sub-models. Recently, there are several improved NPs focused on considering model stability and

improving generalization performance empirically, such as hierarchical prior [27], stochastic attention mechanism [15], bootstrap [13], and Mixture of Expert [28]. However, most of them are unable to investigate the theoretical bound of the generalization performance of NPs. It is desirable to develop robust algorithms with low generalization error and high efficiency.

In this paper, we investigate NP-related models and explore more expressive stability toward general stochastic processes by proposing a stable solution. Specifically, by introducing the notion of stability into NPs, we focus on developing theoretical guidelines for deriving a stable NPs solution. We propose a method to find out subsets that are harder to predict than average, which is a key step for constructing this optimization problem. Based on it, a new extension of NPs with stable guarantees is formulated, which can be flexible to work with various NPs and achieves less biased approximation with theoretical guarantees. Considering the model adaptivity, an adaptive weighting strategy is proposed. To illustrate the superiority of the proposed stable solution, we perform experiments on synthetic 1D regression, system identification of physics engines, and real-world image completion tasks, and the results demonstrate that NPs with our stable solution are much more robust than original NPs.

## 2 Related Work

In this section, we briefly review two different areas which are highly relevant to the proposed method, neural processes, and algorithmic stability.

**Neural Processes** Neural processes are a well-known member of the stochastic process family by directly capturing uncertainties with deep neural networks, which are not only computationally efficient but also retain a probabilistic interpretation of the model [9, 10, 14, 13]. Starting with conditional neural processes (CNP) [9], there have been several follow-up works to improve NPs in various aspects [6]. Vanilla CNP combines neural networks with the Gaussian process to extract prior knowledge from training data. NP [10] introduces a global latent variable to model uncertainty in a variational manner. Considering the problem of underfitting in the vanilla NP, Attentive NP [14] introduces the attention mechanism to improve the model’s reconstruction quality. [11] introduced convolutional conditional neural process (CONVCNP) models translation equivariance in the data. Wang and Van Hoof [27] presented a doubly stochastic variational process (DSVNP), which combines both global and local latent variables. Lee et al. [18] extended NP using Bootstrap and proposed the bootstrapping neural processes (BNP). Kawano et al. [13] presented a group equivariant conditional neural process by incorporating group equivariant into CNPs in a meta-learning manner. Wang and van Hoof [28] proposed to combine the Mixture of Expert models with NPs to develop more expressive exchangeable stochastic processes. Kim et al. [15] proposed a stochastic attention mechanism for NPs to capture appropriate context information. Although there are many NP variants to improve the model performance, those do not consider stability to yield high generalization performance.

**Algorithmic Stability** Stability, as known as algorithmic stability, is a computational learning theory of how a machine learning algorithm is perturbed by small changes to its inputs [2]. Many efforts have been made to analyze various notions of algorithmic stability and prove that a broad spectrum of learning algorithms are stable in some sense [2, 3, 29, 12]. [3] proved that  $l_2$  regularized learning algorithms are uniformly stable and able to obtain new bounds on generalization performance. [29] generalized [3]’s results and proved that regularized learning algorithms with strongly convex penalty functions on bounded domains. Hardt et al. [12] showed that parametric models trained by stochastic gradient descent algorithms are uniformly stable. Li et al. [19] introduced the stability notation to low-rank matrix approximation. Liu et al. [22] proved that tasks in multi-task learning can act as regularizers and that multi-task learning in a very general setting will therefore be uniformly stable under mild assumptions. This is the first work to investigate the stability of NPs from theoretical guidelines and derive NPs solutions with high stability.

## 3 Preliminary

Let calligraphic letters (e.g.,  $\mathcal{A}$ ) indicate sets, capital letters (e.g.,  $A$ ) indicate scalars, lower-case bold letters (e.g.,  $\mathbf{a}$ ) denote vectors, and capital bold letters (e.g.,  $\mathbf{A}$ ) indicate matrices. Suppose there is a dataset  $\mathcal{D} = (\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  with  $N$  data points  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ , and corresponding labels  $\mathbf{y} = [y_1, y_2, \dots, y_N] \in \mathbb{R}^N$ . Considering an arbitrary number of data points  $\mathcal{D}^{\mathcal{C}} = (\mathbf{X}^{\mathcal{C}}, \mathbf{y}^{\mathcal{C}}) = \{(\mathbf{x}_i, y_i)\}_{i \in \mathcal{C}}$ , where  $\mathcal{C} \subseteq \{1, 2, \dots, N\}$  is an index set defining context information, neural processes model the conditional predictive distribution of the target values  $\mathbf{y}^{\mathcal{T}} = \{y_i\}_{i \in \mathcal{T}}$

at some target data points  $\mathbf{X}^T = \{\mathbf{x}_i\}_{i \in \mathcal{T}}$  based on the context  $\mathcal{D}^C$ , i.e.  $P(\mathbf{y}^T | \mathbf{X}^T, \mathcal{D}^C)$ . Usually, target set is defined as  $\mathcal{T} = \{1, 2, \dots, N\}$ . Only in CNP [9],  $\mathcal{T} \subseteq \{1, 2, \dots, N\}$  and  $\mathcal{T} \cap \mathcal{C} = \emptyset$ . In this paper, we define  $\mathcal{T} = \{1, 2, \dots, N\}$  for all NPs, i.e. conditional predictive distribution is  $P(\mathbf{y} | \mathbf{X}, \mathcal{D}^C) = \prod_{i=1}^N P(y_i | \mathbf{x}_i, \mathcal{D}^C)$ .

Fundamentally, there are two NP variants: deterministic and probabilistic. Deterministic NP [9], i.e. CNP, models the conditional distribution as  $P(\mathbf{y} | \mathbf{X}, \mathcal{D}^C) = P(\mathbf{y} | \mathbf{X}, \mathbf{r}^C)$ , where  $\mathbf{r}^C \in \mathbb{R}^d$  is an aggregated feature vector processed by a function that maps  $\mathcal{D}^C$  into a finite-dimensional vector space in a permutation-invariant way. In probabilistic NPs [10], a latent variable  $\mathbf{z} \in \mathbb{R}^d$  is introduced to capture model uncertainty and the NPs infer  $P_\theta(\mathbf{z} | \mathcal{D}^C)$  given context set using the reparameterization trick [16] and models such a conditional distribution as  $P_\theta(y_i | \mathbf{x}_i, \mathcal{D}^C) = \int P_\theta(y_i | \mathbf{x}_i, \mathcal{D}^C, \mathbf{z}) P_\theta(\mathbf{z} | \mathcal{D}^C) d\mathbf{z}$  and it is trained by maximizing an ELBO:  $\mathbb{E}_{\mathbf{z} \sim P_\theta(\mathbf{z} | \mathbf{X}, \mathcal{Y})} [\log P_\theta(\mathbf{y} | \mathbf{X})] - KL[P_\theta(\mathbf{z} | \mathbf{X}, \mathcal{Y}) \| P_\theta(\mathbf{z} | \mathcal{D}^C)]$ .

**Meta Training NP Prediction** To achieve fast prediction on a new context set at test time, NPs meta-learn a distribution over predictors. To perform meta-learning, we require a meta-dataset (dataset of datasets). We consider an unknown distribution  $\mu$  on an instance space  $\mathcal{X} \times \mathcal{Y}$ , and a set of independent sample  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  drawn from  $\mu$ :  $(\mathbf{x}_i, y_i) \sim \mu$  and  $\mathcal{D} \sim \mu^N$ . Suppose meta-dataset contains  $M$  datasets  $\mathcal{D}_{1:M} = \{\mathcal{D}_m\}_{m=1}^M$  with  $\mathcal{D}_m = \{\mathcal{D}_m^C, \mathcal{D}_m^T\}$ , we assume that all  $M$  datasets drawn from a common environment  $\tau$ , which is a probability measure on the set of probability measures on  $\mathcal{X} \times \mathcal{Y}$ . The draw of  $\mu \sim \tau$  indicates the encounter of a specific learning task  $\mu$  in the environment  $\tau$ . For simplicity, we assume that each dataset has the same sample size  $N$ . Following the previous work related to multi-task learning [24] and meta learning [5], The environment  $\tau$  induces a measure  $\mu_{N,\tau}$  on  $(\mathcal{X} \times \mathcal{Y})^N$  such that  $\mu_{N,\tau}(A) = \mathbb{E}_{\mu \sim \tau} [\mu^N(A)]$ ,  $\forall A \subseteq (\mathcal{X} \times \mathcal{Y})^N$ . Thus a dataset  $\mathcal{D}_m$  is independently sampled from a task  $\mu$  encountered in  $\tau$ , which is denoted as  $\mathcal{D}_m \sim \mu_{N,\tau}$  for  $m \in [M]$ .

Suppose there exists a meta parameter  $\theta$  indicating the shared knowledge among different tasks. In this case, a meta learning algorithm  $\mathcal{A}_{meta}$  for NPs takes meta-datasets  $\mathcal{D}_{1:M}$  as input, and then outputs a meta parameter  $\theta = \mathcal{A}_{meta}(\mathcal{D}_{1:M}) \sim P_{\theta | \mathcal{D}_{1:M}}$ . When given a new test dataset  $\mathcal{D}$ , we can evaluate the quality of the meta parameter  $\theta$  by the following true risk:

$$R_\tau(\theta) = \mathbb{E}_{\mathcal{D} \sim \mu_{N,\tau}} \mathbb{E}_{U \sim P_{\theta | \mathcal{D}_{1:M}}} [R_\mu(\theta)] \quad (1)$$

where  $R_\mu(\theta) = -\mathbb{E}_{(\mathbf{x}_i, y_i) \sim \mu} \log P_\theta(y_i | \mathbf{x}_i, \mathcal{D}^C)$ . Usually,  $\tau$  and  $\mu$  are unknown, we can only estimate the meta parameter  $\theta$  from the observed data  $\mathcal{D}_{1:M}$ . In this case, the empirical risk w.r.t  $\theta$  is:

$$R_{\mathcal{D}_{1:M}}(\theta) = 1/M \sum_{m=1}^M \mathbb{E}_{\theta \sim P_{\theta | \mathcal{D}_m^C}} R_{\mathcal{D}_m}(\theta) \quad (2)$$

where  $R_{\mathcal{D}_m}(\theta) = -(1/N) \sum_{i=1}^N \log P_\theta(y_i | \mathbf{x}_i, \mathcal{D}_m^C)$ .

NPs have various strengths: 1) *Efficiency*: meta-learning allows NPs to incorporate information from a new context set and make predictions with a single forward pass. The complexity is linear or quadratic in the context size instead of cubic as with Gaussian process regression; 2) *Flexibility*: NPs can define a conditional distribution of an arbitrary number of target points, conditioning an arbitrary number of observations; 3) *Permutation invariance*: the encoders of NPs use set property [32] to make the target prediction permutation invariant. Thanks to these properties, NPs are widely-used in lots of tasks, e.g., Bayesian optimization [8], recommendation [20, 21], physics engines controlling [27] etc. While there are many NP variants to improve the performance of NPs [9, 10, 14, 13, 15, 28], those do not take model's stability into consider account yet, which is the key to the robustness of the model.

## 4 Problem Formulation

**Stability of NP** A stable learning algorithm has the property that replacing one element in the training set does not result in a significant change to the algorithm's output [2]. Therefore, if we take the training error as a random variable, the training error of a stable learning algorithm should have a small variance. This implies that stable algorithms have the property that the training errors are close to the testing error [2]. Based on the defined risks, the algorithmic stability of approximate  $\{y_i\}_{i \in \mathcal{T}}$  in NPs is defined as follows.

**Definition 4.1.** (Algorithmic Stability of Neural Processes) For any measure  $\mu_{N,\tau}$  on  $(\mathcal{X} \times \mathcal{Y})^N$  such that  $\mu_{N,\tau}(A) = \mathbb{E}_{\mu \sim \tau} [\mu^N(A)]$ ,  $\forall A \subseteq (\mathcal{X} \times \mathcal{Y})^N$ , sample  $M$  datasets  $\mathcal{D}_{1:M}$  from  $\mu_{N,\tau}$  randomly. For a given  $\epsilon > 0$ , we say that  $R_{\mathcal{D}_{1:M}}(\theta)$  is  $\delta$ -stable if the following holds:

$$P(|R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta)| \leq \epsilon) \geq 1 - \delta. \quad (3)$$

The above stability for NPs has the property that the generalization error is bounded, which indicates that minimizing the training error will have a high probability of minimizing the testing error. This new stability notion makes it possible to measure the generalization performance between different NP approximations. For instance, for any two meta-datasets  $\mathcal{D}_{1:M}^1$  and  $\mathcal{D}_{1:M}^2$  from  $\mu_{N,\tau}$ , train NPs on  $\mathcal{D}_{1:M}^1$  and  $\mathcal{D}_{1:M}^2$  are  $\delta_1$ -stable and  $\delta_2$ -stable, respectively. Then  $R_{\mathcal{D}_{1:M}^1}(\theta)$  is more stable than  $R_{\mathcal{D}_{1:M}^2}(\theta)$  if  $\delta_1 < \delta_2$ . This implies that  $R_{\mathcal{D}_{1:M}^1}(\theta)$  is close to  $R_\tau(\theta)$  with higher probability than  $R_{\mathcal{D}_{1:M}^2}(\theta)$ , i.e. minimizing  $R_{\mathcal{D}_{1:M}^1}(\theta)$  will lead to solutions that are of high probabilities with better generalization performance than minimizing  $R_{\mathcal{D}_{1:M}^2}(\theta)$ .

Based on the above analysis, we can see that the reliability of data points is crucial to the success of NPs and frail NPs are susceptible to noise.

**Stability vs. Generalization Error** The sparsity of the data, incomplete and noisy introduces challenges to the algorithm stability. NP models are biased to the quality of context data and target data, so small changes in the training data (noisy) may significantly change the models. In this case, unstable solutions will introduce high training error variance, and minimizing the training error may not guarantee consistent error reduction on the testing dataset, i.e., low generalization performance. In other words, the algorithm stability has a direct impact on generalization performance, and an unstable NP solution has low generalization performance. We take NPs with 1D regression task as an example [9] to investigate the relationship between generalization performance and stability of NPs.

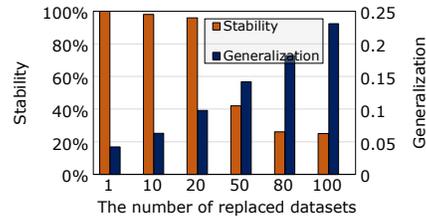


Figure 1: Stability vs. generalization error with different numbers of replaced noisy datasets.

The total number of training and testing datasets is 200 and 100. We trained the NPs model with curves generated from the Gaussian process with RBF kernels by replacing the normal data dataset with a noisy dataset, i.e. the number of replaced datasets is turned in  $\{1, 10, 20, 50, 80, 100\}$ . We quantify stability changes of NPs with the generalization error when the number of replaced datasets increases from 1 to 100. We compute the difference between training error and test error to measure generalization error. We define the difference between test error and training error as  $R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta)$ , and compute  $P(|R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta)| \leq \epsilon)$  with 100 different runs to measure stability. Here we choose  $\epsilon$  in Definition 4.1 as 0.0015 to cover all error differences when the number of replaced datasets is 1. As shown in Figure 1, the generalization error increases when the number of replaced points increases since the testing error becomes lower. On the contrary, the stability of NP decreases with the number of replaced points increases. This indicates that stability decreases with generalization error increases. This study demonstrates that existing NPs suffer from lower generalization performance due to low algorithmic stability. Therefore, it is important to develop a stable solution for NPs that offers good generalization performance.

## 5 Method

In this section, inspired by the previous work [19], we present a stable solution for NPs with stability and high generalization. Algorithmic stability provides an intuitive way to measure the changes in the outputs of a learning algorithm when the input is changed. Various ways have been introduced to measure algorithmic stability. Following the definition of uniform stability [2], given a stable NP, the approximation results remain stable if the change of the datasets. For instance, we can remove a subset of easily predictable data points from  $\mathcal{D}_{1:M}$  to obtain  $\mathcal{D}'_{1:M}$ . It is desirable that the solution of minimizing both  $\mathcal{D}_{1:M}$  and  $\mathcal{D}'_{1:M}$  together will be more stable than the solution of minimizing  $\mathcal{D}_{1:M}$  only. The following Theorem formally proves the statement.

**Theorem 5.1.** *Let  $\mathcal{D}_{1:M}$  ( $M \geq 2$ ) be a sampled meta-dataset of measure  $\mu_{N,\tau}$ . Let  $\mathcal{D}_s \in \mathcal{D}_{1:M}$  be a subset of the meta-dataset, which satisfy that  $\forall (\mathbf{x}_i, y_i) \in \mathcal{D}_s, -\log P_\theta(y_i | \mathbf{x}_i, \mathcal{D}^C) \leq R_{\mathcal{D}_{1:M}}(\theta)$ . Let  $\mathcal{D}'_{1:M} = \mathcal{D}_{1:M} - \mathcal{D}_s$ , then for any  $\epsilon > 0$  and  $1 > w_0 > 0, 1 > w_1 > 0$  ( $w_0 + w_1 = 1$ ),  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta)$  and  $R_{\mathcal{D}_{1:M}}(\theta)$  are  $\delta_1$ -stable and  $\delta_2$ -stable, respectively, then  $\delta_1 \leq \delta_2$ .*

*Proof.* Let's assume that  $R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta) \in [-a_1, a_1]$  and  $R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta)) \in [-a_2, a_2]$  are two random variables with zero mean, where  $a_1 = \sup\{R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta)\}$  and  $a_2 =$

$\sup\{R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta))\}$ . Based on Markov's inequality<sup>1</sup>, for any  $t > 0$ , we have

$$P(R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta) \geq \epsilon) \leq \frac{\mathbb{E}\left[e^{t(R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta))}\right]}{e^{t\epsilon}}. \quad (4)$$

Based on Hoeffding's lemma<sup>2</sup>, we have  $\mathbb{E}[e^{t(R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta))}] \leq e^{\frac{1}{2}t^2 a_1^2}$ , i.e.  $P(R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta) \geq \epsilon) \leq \frac{e^{\frac{1}{2}t^2 a_1^2}}{e^{t\epsilon}}$ . Similarly, we have  $P(R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta) \leq -\epsilon) \leq \frac{e^{\frac{1}{2}t^2 a_1^2}}{e^{t\epsilon}}$ . Combining those two inequalities, we have  $P(|R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta)| \geq \epsilon) \leq \frac{2e^{\frac{1}{2}t^2 a_1^2}}{e^{t\epsilon}}$ , i.e.

$$P(|R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta)| \leq \epsilon) \geq 1 - \frac{2e^{\frac{1}{2}t^2 a_1^2}}{e^{t\epsilon}}. \quad (5)$$

Similarly, we have

$$P\left(|R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta))| \leq \epsilon\right) \geq 1 - \frac{2e^{\frac{1}{2}t^2 a_2^2}}{e^{t\epsilon}}. \quad (6)$$

In this case, the relationship between  $a_1$  and  $a_2$  is

$$\begin{aligned} a_2 &= \sup\left\{R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta) + w_1\left(R_{\mathcal{D}_{1:M}}(\theta) - R_{\mathcal{D}'_{1:M}}(\theta)\right)\right\} \\ &= \sup\{R_\tau(\theta) - R_{\mathcal{D}_{1:M}}(\theta)\} + w_1 \sup\left\{R_{\mathcal{D}_{1:M}}(\theta) - R_{\mathcal{D}'_{1:M}}(\theta)\right\} \\ &= a_1 + \lambda_1 \sup\left\{R_{\mathcal{D}_{1:M}}(\theta) - R_{\mathcal{D}'_{1:M}}(\theta)\right\} \end{aligned} \quad (7)$$

Since  $\forall(\mathbf{x}_i, y_i) \in \mathcal{D}_s$ ,  $-\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}_s^C) \leq R_{\mathcal{D}_{1:M}}(\theta)$ , we have  $-(1/N) \sum_{i=1}^N \log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}_s^C) \leq R_{\mathcal{D}_{1:M}}(\theta)$ . Then, since  $\mathcal{D}_{1:M} = \mathcal{D}_s \cup \mathcal{D}'_{1:M}$ . This means that  $\sup\{R_{\mathcal{D}_{1:M}}(\theta) - R_{\mathcal{D}'_{1:M}}(\theta)\} \leq 0$ . Thus, we have  $a_2 \leq a_1$ . This turns out that  $\frac{2e^{\frac{1}{2}t^2 a_2^2}}{e^{t\epsilon}} \leq \frac{2e^{\frac{1}{2}t^2 a_1^2}}{e^{t\epsilon}}$ , i.e.  $\delta_1 \leq \delta_2$ .  $\square$

Theorem 5.1 indicates that, if we remove a subset that is easier to predict than average from  $\mathcal{D}_{1:M}$  to form  $\mathcal{D}'_{1:M}$ , then  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta)$  has higher probability of being close to  $R_\tau(\theta)$  than  $R_{\mathcal{D}_{1:M}}(\theta)$ . Therefore, minimizing  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta)$  will lead to solutions that have better generalization performance than minimizing  $R_{\mathcal{D}_{1:M}}(\theta)$ .

However, Theorem 5.1 only proves that it is beneficial to remove an easily predictable dataset from  $\mathcal{D}_{1:M}$  to obtain  $\mathcal{D}'_{1:M}$ , but does not show how many datasets we should remove from  $\mathcal{D}_{1:M}$ . Actually, removing more datasets that satisfy  $-\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^C) \leq R_{\mathcal{D}_{1:M}}(\theta)$  can obtain better  $\mathcal{D}'_{1:M}$ , as shown in following Theorem 5.2.

**Theorem 5.2.** Let  $\mathcal{D}_{1:M}$  ( $M \geq 2$ ) be a sampled meta-dataset of measure  $\mu_{N,\tau}$ . Let  $\mathcal{D}_{s1}$  and  $\mathcal{D}_{s2}$  be two subsets of  $\mathcal{D}_{1:M}$ , which satisfy  $\mathcal{D}_{s2} \subset \mathcal{D}_{s1} \subset \mathcal{D}_{1:M}$ .  $\mathcal{D}_{s2}$  and  $\mathcal{D}_{s1}$  satisfy that  $\forall(\mathbf{x}_i, y_i) \in \mathcal{D}_{s1}$ ,  $-\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^C) \leq R_{\mathcal{D}_{1:M}}(\theta)$ . Let  $\mathcal{D}'_{1:M} = \mathcal{D}_{1:M} - \mathcal{D}_{s1}$  and  $\mathcal{D}^2_{1:M} = \mathcal{D}_{1:M} - \mathcal{D}_{s2}$ , then for any  $\epsilon > 0$  and  $1 > w_0 > 0$ ,  $1 > w_1 > 0$  ( $w_0 + w_1 = 1$ ),  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta)$  and  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}^2_{1:M}}(\theta)$  are  $\delta_1$ -stable and  $\delta_2$ -stable, respectively, then  $\delta_1 \leq \delta_2$ .

*Proof.* Let's assume that  $R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta)) \in [-a_1, a_1]$  and  $R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}^2_{1:M}}(\theta)) \in [-a_2, a_2]$  are two random variables with 0 mean, where  $a_1 = \sup\{R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta))\}$  and  $a_2 = \sup\{R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}^2_{1:M}}(\theta))\}$ .

Then, based on Markov's inequality and Hoeffding's lemma, we have

$$\begin{aligned} P\left(|R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}'_{1:M}}(\theta))| \leq \epsilon\right) &\geq 1 - \frac{2e^{\frac{1}{2}t^2 a_1^2}}{e^{t\epsilon}}, \\ P\left(|R_\tau(\theta) - (w_0 R_{\mathcal{D}_{1:M}}(\theta) + w_1 R_{\mathcal{D}^2_{1:M}}(\theta))| \leq \epsilon\right) &\geq 1 - \frac{2e^{\frac{1}{2}t^2 a_2^2}}{e^{t\epsilon}}. \end{aligned} \quad (8)$$

Since  $\forall(\mathbf{x}_i, y_i) \in \mathcal{D}_{s1}$ ,  $-\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^C) \leq R_{\mathcal{D}_{1:M}}(\theta)$  and  $\mathcal{D}_{s2} \subset \mathcal{D}_{s1} \subset \mathcal{D}_{1:M}$ , we have  $R_{\mathcal{D}'_{1:M}}(\theta) \leq R_{\mathcal{D}^2_{1:M}}(\theta)$ . Thus, we have  $\sup\{R_{\mathcal{D}'_{1:M}}(\theta) - R_{\mathcal{D}^2_{1:M}}(\theta)\} \leq 0$  since  $a_1 = a_2 + w_1 \sup\{R_{\mathcal{D}'_{1:M}}(\theta) - R_{\mathcal{D}^2_{1:M}}(\theta)\}$  and  $a_1 \leq a_2$ . Then, we can conclude that  $\delta_1 \leq \delta_2$ .  $\square$

<sup>1</sup>(Extended version of Markov's Inequality) Let  $x$  be a real-valued non-negative random variable and  $\varphi(\cdot)$  be a nondecreasing nonnegative function with  $\varphi(a) > 0$ . Then, for any  $\epsilon > 0$ ,  $P(x \geq \epsilon) \leq \frac{\mathbb{E}[\varphi(x)]}{\varphi(\epsilon)}$ .

<sup>2</sup>(Hoeffding's Lemma) Let  $x$  be a real-valued random variable with zero mean and  $p(x \in [a, b]) = 1$ . Then, for any  $z \in \mathbb{R}$ ,  $\mathbb{E}[e^{zx}] \leq \exp\left(\frac{1}{8}z^2(b-a)^2\right)$ .

Theorem 5.2 indicates that removing more data points that are easy to predict will obtain more stable NPs. Therefore, it is desirable to choose  $\mathcal{D}'_{1:M}$  (i.e.  $\mathcal{D}_{1:M} - \mathcal{D}_s$ ) as the whole set which is harder to predict than average, i.e. the whole set satisfying  $\forall(\mathbf{x}_i, y_i) \in \mathcal{D}'_{1:M}, -\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^c) \leq R_{\mathcal{D}'_{1:M}}(\theta)$ . Without loss of generality, we can extend Theorem 5.2 by considering several harder predicted sets to obtain a more stable solution by minimizing them all together. In this case, we need to prove that the stability of a model with  $K$  subsets is better than the model with  $K - 1$  subsets, as shown in Theorem 5.3.

**Theorem 5.3.** *Let  $\mathcal{D}_{1:M}$  ( $M \geq 2$ ) be a sampled meta-dataset of measure  $\mu_{N,\tau}$ . Let  $\mathcal{D}_{s1}, \mathcal{D}_{s2}, \dots, \mathcal{D}_{sK} \subset \mathcal{D}_{1:M}$  be  $K$  subsets and satisfy  $\forall(\mathbf{x}_i, y_i) \in \mathcal{D}_{sk} (k \in [K]), -\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^c) \leq R_{\mathcal{D}_{1:M}}(\theta)$ . Let  $\mathcal{D}'_{1:M} = \mathcal{D}_{1:M} - \mathcal{D}_{sk}$  for all  $k \in [K]$ , then for any  $\epsilon > 0$  and  $1 > w_k > 0$  for all  $k \in [K]$ , ( $w_0 + w_1 + \dots + w_K = 1$ ),  $w_0 R_{\mathcal{D}'_{1:M}}(\theta) + \sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}}(\theta)$  and  $(w_0 + w_K) R_{\mathcal{D}'_{1:M}}(\theta) + \sum_{k=1}^{K-1} w_k R_{\mathcal{D}_{1:M}}(\theta)$  are  $\delta_1$ -stable and  $\delta_2$ -stable, respectively, then  $\delta_1 \leq \delta_2$ .*

*Proof.* For simplicity, we denote  $w_0 R_{\mathcal{D}'_{1:M}}(\theta) + \sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}}(\theta)$  as  $R_1$  and  $(w_0 + w_K) R_{\mathcal{D}'_{1:M}}(\theta) + \sum_{k=1}^{K-1} w_k R_{\mathcal{D}_{1:M}}(\theta)$  as  $R_2$ . Let's assume that  $R_\tau(\theta) - R_1 \in [-a_1, a_1]$  and  $R_\tau(\theta) - R_2 \in [-a_2, a_2]$  are two random variables with 0 mean, where  $a_1 = \sup\{R_\tau(\theta) - R_1\}$  and  $a_2 = \sup\{R_\tau(\theta) - R_2\}$ . Then, based on Markov's inequality and Hoeffding's lemma, we have

$$P(|R_\tau(\theta) - R_1| \leq \epsilon) \geq 1 - \frac{2e^{\frac{1}{2}\epsilon^2 a_1^2}}{e^{\epsilon^2}}, \quad (|R_\tau(\theta) - R_2| \leq \epsilon) \geq 1 - \frac{2e^{\frac{1}{2}\epsilon^2 a_2^2}}{e^{\epsilon^2}}. \quad (9)$$

Similar to the proof of Theorem 5.3, we have  $R_{\mathcal{D}'_{1:M}}(\theta) \geq R_{\mathcal{D}_{1:M}}(\theta)$ , which indicates that  $\sup\{R_{\mathcal{D}'_{1:M}}(\theta) - R_{\mathcal{D}_{1:M}}(\theta)\} \leq 0$ . Since  $a_2 = a_1 + w_K \sup\{R_{\mathcal{D}'_{1:M}}(\theta) - R_{\mathcal{D}_{1:M}}(\theta)\}$ , we know that  $a_2 \leq a_1$ . Thus we can conclude that  $\frac{2e^{\frac{1}{2}\epsilon^2 a_2^2}}{e^{\epsilon^2}} \leq \frac{2e^{\frac{1}{2}\epsilon^2 a_1^2}}{e^{\epsilon^2}}$ , i.e.  $\delta_1 \leq \delta_2$ .  $\square$

Based on Theorem 5.3, we know that optimization on  $\mathcal{D}_{1:M}$  and more than one hard predictable subsets of  $\mathcal{D}_{1:M}$  can achieve more stable prediction. However, how to select data points that are difficult to predict from  $\mathcal{D}_{1:M}$  is still a challenging problem, especially, since we need to select  $K$  subsets.

## 5.1 The Proposed Solution

According to the analysis of stability in NPs, we propose a stable solution for NPs to achieve model stability with the aid of hard predictable subsets selection. Specifically, we introduce a solution to obtain those hard predictable subsets based on *only one* set of easily predicted data points which can be broken into four steps:

- (1) Selecting a existing NP model (e.g., CNP [9], NP[10]) and training it with meta-dataset  $\mathcal{D}_{1:M}$ ;
- (2) Selecting an easily predicted subset  $\mathcal{D}_s \subset \mathcal{D}_{1:M}$ , which satisfies that  $\forall(\mathbf{x}_i, y_i) \in \mathcal{D}_s, -\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^c) \leq R_{\mathcal{D}_{1:M}}(\theta)$ ;
- (3) Dividing  $\mathcal{D}_s$  into  $K$  non-overlapping subsets  $\mathcal{D}_{s1}, \mathcal{D}_{s2}, \dots, \mathcal{D}_{sK}$  and satisfying  $\cup_{k=1}^K \mathcal{D}_{sk} = \mathcal{D}_s$ ;
- (4) Defining  $K$  subsets that are difficult to predict, i.e.  $\mathcal{D}'_{1:M} = \mathcal{D}_{1:M} - \mathcal{D}_{sk}$  for all  $k \in [K]$ ;

Thus, a new extension of NPs is given as

$$\mathcal{L} = \arg \min_{\theta} w_0 R_{\mathcal{D}'_{1:M}}(\theta) + \sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}}(\theta), \quad (10)$$

where  $w_0, w_1, \dots, w_K$  indicate the contributions of each component and satisfy  $\sum_{k=0}^K w_k = 1$ .

The whole learning algorithm is given in Algorithm 1. From steps 1 to 12, we obtain  $K$  different hard predictable subsets, and the complexity of lines 1 to 12 is  $O(MN)$ . The complexity of line 11 is related to the applied NP models (such as CNP, NP, ANP, etc). Thus, the computational complexity of NPs with our stable solution is similar to the original NPs. As shown in Algorithm 1, we need to pre-train the base model to select samples. In fact, the pre-trained model can not only be used for sample selection but its model parameters can be used as the initialization of the stable version model. At this time, the training of the stable version can converge faster.

**Stability Guarantee** Here we give a theoretical guarantee of the proposed stable solution.

**Theorem 5.4.** *Let  $\mathcal{D}_{1:M}$  ( $M \geq 2$ ) be a sampled meta-dataset of measure  $\mu_{N,\tau}$ . Let  $\mathcal{D}_s \subset \mathcal{D}_{1:M}$  which satisfies that  $\forall(\mathbf{x}_i, y_i) \in \mathcal{D}_s, -\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^c) \leq R_{\mathcal{D}_{1:M}}(\theta)$ . By dividing  $\mathcal{D}_s$  into  $K$  subsets*

---

**Algorithm 1** Learning algorithm for stable NPs

---

**Input:** Meta-dataset  $\mathcal{D}_{1:M}$ ,  $\beta > 0.5$  is a predefined probability for data selection.  $\mathcal{D}_s = \emptyset$ .

- 1: Train a NP model with parameter  $\theta$  based on  $R_{\mathcal{D}_{1:M}}(\theta)$ ;
- 2: **for**  $(\mathbf{x}_i, y_i) \in \mathcal{D}_{1:M}$  **do**
- 3:     randomly generate split parameter  $\rho_i \in [0, 1]$ ;
- 4:     **if**  $((-\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^c) \leq R_{\mathcal{D}_{1:M}}(\theta)) \ \& \ \rho_i \leq \beta)$
- 5:     **or**  $((-\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^c) > R_{\mathcal{D}_{1:M}}(\theta)) \ \& \ \rho_i \leq 1 - \beta)$  **then**
- 6:          $\mathcal{D}_s \leftarrow \mathcal{D}_s \cup (\mathbf{x}_i, y_i)$ ;
- 7:     **end if**
- 8: **end for**
- 9: divide  $\mathcal{D}_s$  into  $\mathcal{D}_{s1}, \mathcal{D}_{s2}, \dots, \mathcal{D}_{sK}$  with  $\cup_{k=1}^K \mathcal{D}_{sk} = \mathcal{D}_s$ ;
- 10: **for**  $k = 1, 2, \dots, K$  **do**
- 11:      $\mathcal{D}_{1:M}^k = \mathcal{D}_{1:M} - \mathcal{D}_{sk}$ ;
- 12: **end for**
- 13: update parameters by optimizing  $\theta^* = \arg \min_\theta w_0 R_{\mathcal{D}_{1:M}}(\theta) + \sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}^k}(\theta)$ ;

**Output:** The learned optimal parameters  $\theta^*$ .

---

$\mathcal{D}_{s1}, \mathcal{D}_{s2}, \dots, \mathcal{D}_{sK}$  which satisfy that  $\cup_{k=1}^K \mathcal{D}_{sk} = \mathcal{D}_{1:M}$ . Let  $\mathcal{D}_{1:M}^0 = \mathcal{D}_{1:M} - \mathcal{D}_s$  and  $\mathcal{D}_{1:M}^k = \mathcal{D}_{1:M} - \mathcal{D}_{sk}$  for all  $k \in [K]$ , then for any  $\epsilon > 0$  and  $1 > w_k > 0$  for all  $k \in [K]$ , ( $w_0 + w_1 + \dots + w_K = 1$ ),  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + \sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}^k}(\theta)$  and  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + (1 - w_0) R_{\mathcal{D}_{1:M}^0}(\theta)$  are  $\delta_1$ -stable and  $\delta_2$ -stable, respectively, then  $\delta_1 \leq \delta_2$ .

*Proof.* For simplicity, we denote  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + \sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}^k}(\theta)$  as  $R_1$  and  $w_0 R_{\mathcal{D}_{1:M}}(\theta) + (1 - w_0) R_{\mathcal{D}_{1:M}^0}(\theta)$  as  $R_2$ . Let's assume that  $R_\tau(\theta) - R_1 \in [-a_1, a_1]$  and  $R_\tau(\theta) - R_2 \in [-a_2, a_2]$  are two random variables with 0 mean, where  $a_1 = \sup\{R_\tau(\theta) - R_1\}$  and  $a_2 = \sup\{R_\tau(\theta) - R_2\}$ . Then, based on Markov's inequality and Hoeffding's lemma, we have

$$P(|R_\tau(\theta) - R_1| \leq \epsilon) \geq 1 - \frac{2e^{\frac{1}{2}t^2 a_1^2}}{e^{t\epsilon}}, \quad P(|R_\tau(\theta) - R_2| \leq \epsilon) \geq 1 - \frac{2e^{\frac{1}{2}t^2 a_2^2}}{e^{t\epsilon}}. \quad (11)$$

$\forall k \in [K]$ ,  $\mathcal{D}_{sk} \subset \mathcal{D}_s$  and  $\forall (\mathbf{x}_i, y_i) \in \mathcal{D}_s$ ,  $-\log P_\theta(y_i|\mathbf{x}_i, \mathcal{D}^c) \leq R_{\mathcal{D}_{1:M}}(\theta)$ , we have  $R_{\mathcal{D}_{1:M}^k}(\theta) \leq R_{\mathcal{D}_{1:M}^0}(\theta)$ . By combining the above inequalities over all  $k \in [K]$ , we have

$$\sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}^k}(\theta) \leq \sum_{k=1}^K w_k R_{\mathcal{D}_{1:M}^0}(\theta) = (1 - w_0) R_{\mathcal{D}_{1:M}^0}(\theta). \quad (12)$$

Thus,  $\sup\{R_\tau(\theta) - R_1\} \leq \sup\{R_{1:M}(\theta) - R_2\}$ , i.e.  $a_1 \leq a_2$ . Thus we have  $\delta_1 \leq \delta_2$ .  $\square$

According to the above theorem, we can achieve model stability by selecting only one easily predicted subset.

## 6 Experiments

We started with learning predictive functions on synthetic datasets, and then high-dimensional tasks, e.g., system identification on physics engines, image completion, and Bayesian optimization, were performed to evaluate the properties of the NP-related models.

### 6.1 1D Regression

To verify the proposed stable solution, we combined the stable solution with different baseline NP classes (CNP [9], NP [10], ANP [14], ConvCNP [11], ConvNP [6], and their bootstrapping versions [18]) and compared them on 1D regression task. Among them, BCNP, BNP, BANP, BConvCNP, and NConvNP are recently proposed stable strategies for NPs with Bootstrap. Specifically, the stochastic process (SP) initializing with a 0 mean Gaussian Process (GP)  $y^{(0)} \sim GP(0, k(\cdot, \cdot))$  indexed in the interval  $x \in [-2.0, 2.0]$  were used to generate data, where the radial basis function kernel and Matern Kernel were adopted for model-data mismatch scenario. More detailed information can be obtained in the Appendix. We investigated the model performance in terms of different noise settings. We introduced Gaussian noise  $\mathcal{N}(0, 1)$  and added noise to different proportions of the data, such as  $\{0\%, 5\%, 10\%, 15\%\}$ . Table 1 lists the average log-likelihoods comparison in terms of different noise proportions. The best result is marked in bold. First, we can see that if we adopt the robust solution in baselines, the model achieves the best results on all the datasets, showing the effectiveness of the

Table 1: Average Log-likelihoods over all context and target points on realizations from Synthetic Stochastic Process on the different percent of added noise. Here we set the context size to 20. (Mean  $\pm$  Std). Note that adding ‘S’ before the original model name is a model with our stable solution.

Kernel	Method	Original		Noise(+5%)		Noise(+10%)		Noise(+15%)	
		context	target	context	target	context	target	context	target
RBF	CNP	0.8724 $\pm$ 0.008	0.4334 $\pm$ 0.007	0.8522 $\pm$ 0.005	0.4001 $\pm$ 0.010	0.8014 $\pm$ 0.006	0.3552 $\pm$ 0.004	0.7152 $\pm$ 0.006	0.2853 $\pm$ 0.005
	BCNP	0.9042 $\pm$ 0.009	0.4589 $\pm$ 0.006	0.8774 $\pm$ 0.006	0.4278 $\pm$ 0.008	0.8316 $\pm$ 0.006	0.3767 $\pm$ 0.005	0.7487 $\pm$ 0.007	0.3017 $\pm$ 0.006
	SCNP	<b>0.9255</b> $\pm$ 0.008	<b>0.4733</b> $\pm$ 0.004	<b>0.8935</b> $\pm$ 0.005	<b>0.4478</b> $\pm$ 0.006	<b>0.8517</b> $\pm$ 0.004	<b>0.3986</b> $\pm$ 0.005	<b>0.7621</b> $\pm$ 0.005	<b>0.3279</b> $\pm$ 0.006
	NP	0.8215 $\pm$ 0.004	0.3853 $\pm$ 0.005	0.8011 $\pm$ 0.004	0.3511 $\pm$ 0.006	0.7611 $\pm$ 0.005	0.3042 $\pm$ 0.008	0.6722 $\pm$ 0.005	0.2435 $\pm$ 0.007
	BNP	0.8722 $\pm$ 0.004	0.4211 $\pm$ 0.004	0.8321 $\pm$ 0.003	0.3876 $\pm$ 0.004	0.7922 $\pm$ 0.004	0.3389 $\pm$ 0.005	0.7189 $\pm$ 0.004	0.2776 $\pm$ 0.007
	SNP	<b>0.8955</b> $\pm$ 0.003	<b>0.4356</b> $\pm$ 0.004	<b>0.8567</b> $\pm$ 0.003	<b>0.4046</b> $\pm$ 0.005	<b>0.8165</b> $\pm$ 0.004	<b>0.3568</b> $\pm$ 0.006	<b>0.7356</b> $\pm$ 0.005	<b>0.2955</b> $\pm$ 0.006
	ANP	1.2563 $\pm$ 0.002	0.5763 $\pm$ 0.004	1.2245 $\pm$ 0.007	0.5347 $\pm$ 0.006	1.1742 $\pm$ 0.005	0.4871 $\pm$ 0.007	1.0982 $\pm$ 0.005	0.4151 $\pm$ 0.004
	BANP	1.2722 $\pm$ 0.004	0.5887 $\pm$ 0.006	1.2411 $\pm$ 0.005	0.5471 $\pm$ 0.005	1.1886 $\pm$ 0.006	0.4917 $\pm$ 0.006	1.0642 $\pm$ 0.006	0.4327 $\pm$ 0.005
	SANP	<b>1.2831</b> $\pm$ 0.000	<b>0.5994</b> $\pm$ 0.004	<b>1.2564</b> $\pm$ 0.004	<b>0.5578</b> $\pm$ 0.004	<b>1.2052</b> $\pm$ 0.004	<b>0.5025</b> $\pm$ 0.006	<b>1.1243</b> $\pm$ 0.005	<b>0.4356</b> $\pm$ 0.006
	ConvCNP	1.2631 $\pm$ 0.002	0.6421 $\pm$ 0.002	1.2333 $\pm$ 0.005	0.5415 $\pm$ 0.005	1.1827 $\pm$ 0.004	0.4936 $\pm$ 0.005	1.0241 $\pm$ 0.006	0.4262 $\pm$ 0.005
	BConvCNP	1.2761 $\pm$ 0.004	0.6531 $\pm$ 0.005	1.2476 $\pm$ 0.004	0.5533 $\pm$ 0.004	1.1931 $\pm$ 0.007	0.4986 $\pm$ 0.005	1.0716 $\pm$ 0.006	0.4396 $\pm$ 0.005
	SConvCNP	<b>1.3991</b> $\pm$ 0.001	<b>0.6793</b> $\pm$ 0.004	<b>1.2651</b> $\pm$ 0.003	<b>0.5623</b> $\pm$ 0.005	<b>1.2126</b> $\pm$ 0.004	<b>0.5096</b> $\pm$ 0.005	<b>1.1331</b> $\pm$ 0.006	<b>0.4461</b> $\pm$ 0.005
	ConvNP	1.2874 $\pm$ 0.003	0.6503 $\pm$ 0.004	1.2371 $\pm$ 0.006	0.5451 $\pm$ 0.006	1.1865 $\pm$ 0.004	0.4965 $\pm$ 0.006	0.9915 $\pm$ 0.005	0.4335 $\pm$ 0.006
	BConvNP	1.2922 $\pm$ 0.004	0.6627 $\pm$ 0.006	1.2505 $\pm$ 0.006	0.5583 $\pm$ 0.004	1.1971 $\pm$ 0.005	0.5025 $\pm$ 0.007	1.0731 $\pm$ 0.006	0.4436 $\pm$ 0.004
	SConvNP	<b>1.4036</b> $\pm$ 0.002	<b>0.6831</b> $\pm$ 0.003	<b>1.2671</b> $\pm$ 0.005	<b>0.5675</b> $\pm$ 0.004	<b>1.2188</b> $\pm$ 0.003	<b>0.5157</b> $\pm$ 0.005	<b>1.1389</b> $\pm$ 0.004	<b>0.4505</b> $\pm$ 0.005
Matern	CNP	0.8531 $\pm$ 0.005	0.2431 $\pm$ 0.010	0.8231 $\pm$ 0.005	0.2144 $\pm$ 0.010	0.7761 $\pm$ 0.008	0.1784 $\pm$ 0.007	0.7052 $\pm$ 0.005	0.1452 $\pm$ 0.006
	BCNP	0.8778 $\pm$ 0.005	0.2762 $\pm$ 0.009	0.8487 $\pm$ 0.006	0.2477 $\pm$ 0.009	0.8015 $\pm$ 0.007	0.2051 $\pm$ 0.007	0.7378 $\pm$ 0.005	0.1766 $\pm$ 0.006
	SCNP	<b>0.8963</b> $\pm$ 0.003	<b>0.2953</b> $\pm$ 0.006	<b>0.8689</b> $\pm$ 0.005	<b>0.2658</b> $\pm$ 0.007	<b>0.8268</b> $\pm$ 0.006	<b>0.2258</b> $\pm$ 0.006	<b>0.7567</b> $\pm$ 0.004	<b>0.1936</b> $\pm$ 0.005
	NP	0.7643 $\pm$ 0.015	0.2041 $\pm$ 0.015	0.7342 $\pm$ 0.002	0.1725 $\pm$ 0.008	0.6892 $\pm$ 0.004	0.1542 $\pm$ 0.006	0.6235 $\pm$ 0.008	0.1342 $\pm$ 0.007
	BNP	0.8156 $\pm$ 0.005	0.2689 $\pm$ 0.007	0.7789 $\pm$ 0.004	0.2215 $\pm$ 0.005	0.7421 $\pm$ 0.005	0.2117 $\pm$ 0.007	0.6715 $\pm$ 0.006	0.1828 $\pm$ 0.006
	SNP	<b>0.8368</b> $\pm$ 0.006	<b>0.2844</b> $\pm$ 0.005	<b>0.8036</b> $\pm$ 0.003	<b>0.2483</b> $\pm$ 0.003	<b>0.7635</b> $\pm$ 0.004	<b>0.2325</b> $\pm$ 0.006	<b>0.6973</b> $\pm$ 0.006	<b>0.2016</b> $\pm$ 0.005
	ANP	1.2421 $\pm$ 0.002	0.6366 $\pm$ 0.004	1.2115 $\pm$ 0.001	0.6001 $\pm$ 0.008	1.1784 $\pm$ 0.004	0.1622 $\pm$ 0.006	1.1252 $\pm$ 0.007	0.5274 $\pm$ 0.008
	BANP	1.3456 $\pm$ 0.003	0.6514 $\pm$ 0.005	1.3125 $\pm$ 0.005	0.6115 $\pm$ 0.002	1.2672 $\pm$ 0.004	0.1711 $\pm$ 0.005	1.2236 $\pm$ 0.006	0.5306 $\pm$ 0.006
	SANP	<b>1.3721</b> $\pm$ 0.002	<b>0.6653</b> $\pm$ 0.004	<b>1.3461</b> $\pm$ 0.003	<b>0.6256</b> $\pm$ 0.004	<b>1.3011</b> $\pm$ 0.003	<b>0.1782</b> $\pm$ 0.004	<b>1.2457</b> $\pm$ 0.005	<b>0.5356</b> $\pm$ 0.002
	ConvCNP	1.2515 $\pm$ 0.003	0.6418 $\pm$ 0.004	1.2226 $\pm$ 0.006	0.6085 $\pm$ 0.005	1.1832 $\pm$ 0.005	0.1871 $\pm$ 0.007	1.1326 $\pm$ 0.005	0.5351 $\pm$ 0.004
	BConvCNP	1.3527 $\pm$ 0.005	0.6616 $\pm$ 0.006	1.3252 $\pm$ 0.005	0.6235 $\pm$ 0.007	1.2767 $\pm$ 0.006	0.1952 $\pm$ 0.005	1.1315 $\pm$ 0.006	0.5417 $\pm$ 0.006
	SConvCNP	<b>1.3852</b> $\pm$ 0.003	<b>0.6731</b> $\pm$ 0.004	<b>1.3364</b> $\pm$ 0.004	<b>0.6335</b> $\pm$ 0.005	<b>1.2831</b> $\pm$ 0.003	<b>0.2037</b> $\pm$ 0.005	<b>1.1521</b> $\pm$ 0.005	<b>0.5557</b> $\pm$ 0.004
	ConvNP	1.2746 $\pm$ 0.002	0.6557 $\pm$ 0.005	1.2345 $\pm$ 0.007	0.6015 $\pm$ 0.005	1.1865 $\pm$ 0.006	0.1943 $\pm$ 0.005	1.1358 $\pm$ 0.005	0.5397 $\pm$ 0.003
	BConvNP	1.3356 $\pm$ 0.004	0.6787 $\pm$ 0.006	1.3305 $\pm$ 0.005	0.6383 $\pm$ 0.006	1.2851 $\pm$ 0.005	0.2015 $\pm$ 0.005	1.1415 $\pm$ 0.006	0.5338 $\pm$ 0.003
	SConvNP	<b>1.3878</b> $\pm$ 0.002	<b>0.6836</b> $\pm$ 0.004	<b>1.3435</b> $\pm$ 0.005	<b>0.6417</b> $\pm$ 0.004	<b>1.2866</b> $\pm$ 0.004	<b>0.2025</b> $\pm$ 0.005	<b>1.1521</b> $\pm$ 0.005	<b>0.5363</b> $\pm$ 0.006

stable solution. Besides, performances on all methods become less accurate in more complicated settings, while our solution has fewer effects. One interesting observation is that the improvements against the base model on CNP are less significant than NP and ANP. The possible reason is that CNP only predicts points out of the context set.

To investigate the model’s ability to address model-data mismatch scenarios, we conducted experiments on 1D regression tasks with Periodic kernel. Following the setting of BANP and similar noise settings of our previous kernels, we list the results on both original data and noise(+15) data in Table 2. In this model-data mismatch data, stable versions still significantly outperform their corresponding original versions.

Table 2: Experiments on 1D regression data with Periodic kernel.

Periodic	Original		Noise(+15%)	
	context	target	context	target
ANP	0.5730 $\pm$ 0.006	-4.2345 $\pm$ 0.005	0.3521 $\pm$ 0.005	-5.3211 $\pm$ 0.007
ConvCNP	0.5983 $\pm$ 0.006	-4.0215 $\pm$ 0.005	0.3658 $\pm$ 0.005	-4.5233 $\pm$ 0.008
ConvNP	0.6125 $\pm$ 0.005	-3.8952 $\pm$ 0.006	0.3756 $\pm$ 0.006	-4.3413 $\pm$ 0.008
BANP	0.6253 $\pm$ 0.003	-3.5413 $\pm$ 0.005	0.3651 $\pm$ 0.006	-4.2511 $\pm$ 0.015
BConvCNP	0.6342 $\pm$ 0.005	-3.4142 $\pm$ 0.004	0.3712 $\pm$ 0.004	-4.1750 $\pm$ 0.008
BConvNP	0.6355 $\pm$ 0.004	-3.3627 $\pm$ 0.005	0.3768 $\pm$ 0.008	-4.0116 $\pm$ 0.007
SANP	0.6315 $\pm$ 0.002	-3.3317 $\pm$ 0.004	0.3748 $\pm$ 0.004	-4.0515 $\pm$ 0.003
SConvCNP	0.6433 $\pm$ 0.002	-3.1515 $\pm$ 0.004	0.3866 $\pm$ 0.005	-3.9851 $\pm$ 0.004
SConvNP	0.6551 $\pm$ 0.000	-3.1062 $\pm$ 0.004	0.3981 $\pm$ 0.002	-3.8895 $\pm$ 0.004

## 6.2 Image Completion

Image completion can be regarded as a 2D function regression task and be interpreted as being generated from a stochastic process (since there are dependencies between pixel values). Following the setting in previous work [14], we trained the NPs on EMNIST [4] and  $32 \times 32$  CELEBA [23] using the standard train/test split with up to 200 context/target points at training. Detailed experiment settings are given in the Appendix. We evaluated average log-likelihoods over all points on realizations from image completion. Table 3 lists the comparisons between NPs with and without our stable solution in terms of original setting and noise setting, and the performance demonstrates the superiority of our stable solution.

## 6.3 System Identification on Physics Engines

The second synthetic experiment focuses on evaluating model dynamics on a classical simulator, Cart-Pole systems, which is detailed in [7, 27]. The Cart-Pole swing-up task is a standard benchmark for nonlinear control due to the non-linearity in the dynamics, and the requirement for nonlinear

Table 3: Average Log-likelihoods over all context and target points on EMNIST and CELEBA.

Dataset	Method	Original		Noise(+5%)		Noise(+10%)		Noise(+15%)	
		context	target	context	target	context	target	context	target
EMNIST	CNP	0.9522±0.023	0.7515±0.0015	0.8977±0.0016	0.6336±0.017	0.8242±0.0018	0.5784±0.009	0.6566±0.0017	0.5341±0.016
	BCNP	0.9678±0.010	0.8058±0.008	0.9015±0.008	0.6711±0.009	0.8415±0.007	0.6089±0.009	0.6788±0.006	0.5715±0.006
	SCNP	<b>0.9716±0.008</b>	<b>0.8343±0.006</b>	<b>0.9251±0.008</b>	<b>0.6971±0.007</b>	<b>0.8674±0.006</b>	<b>0.6343±0.007</b>	<b>0.6986±0.005</b>	<b>0.5877±0.005</b>
	NP	0.9678±0.004	0.7756±0.005	0.9011±0.009	0.6941±0.006	0.8544±0.009	0.6455±0.007	0.7034±0.009	0.5865±0.006
	BNP	0.9757±0.005	0.8358±0.005	0.8116±0.007	0.7625±0.006	0.8759±0.007	0.6773±0.007	0.7451±0.005	0.6237±0.005
	SNP	<b>0.9847±0.005</b>	<b>0.8562±0.006</b>	<b>0.8368±0.006</b>	<b>0.7844±0.005</b>	<b>0.8984±0.005</b>	<b>0.6984±0.005</b>	<b>0.7653±0.005</b>	<b>0.6456±0.004</b>
	ANP	1.1125±0.002	1.0321±0.004	0.9815±0.002	0.6366±0.006	0.9021±0.004	0.7053±0.008	0.8454±0.002	0.7034±0.005
	BANP	1.1355±0.003	1.0615±0.005	1.0236±0.002	0.6549±0.005	0.9155±0.004	0.7521±0.006	0.8612±0.003	0.7515±0.005
	SANP	<b>1.1531±0.000</b>	<b>1.0877±0.004</b>	<b>1.0421±0.002</b>	<b>0.6776±0.005</b>	<b>0.9321±0.002</b>	<b>0.7843±0.006</b>	<b>0.8732±0.003</b>	<b>0.7657±0.005</b>
	ConvCNP	1.1363±0.002	1.0461±0.004	1.0252±0.006	0.6448±0.005	0.9116±0.005	0.7115±0.006	0.8621±0.005	0.7246±0.004
	BConvCNP	1.1425±0.004	1.0787±0.006	1.0311±0.005	0.6626±0.005	0.9252±0.006	0.7617±0.006	0.8717±0.006	0.7627±0.005
	SConvCNP	<b>1.1631±0.003</b>	<b>1.0894±0.004</b>	<b>1.0563±0.004</b>	<b>0.6778±0.004</b>	<b>0.9356±0.004</b>	<b>0.7885±0.005</b>	<b>0.8813±0.005</b>	<b>0.7692±0.006</b>
	ConvNP	1.1415±0.002	1.0563±0.004	1.0286±0.006	0.6536±0.006	0.9168±0.005	0.7171±0.007	0.8675±0.005	0.7368±0.005
	BConvNP	1.1526±0.004	1.0837±0.005	1.0415±0.005	0.6684±0.005	0.9285±0.006	0.7762±0.006	0.8742±0.006	0.7727±0.005
	SConvNP	<b>1.1753±0.003</b>	<b>1.0934±0.004</b>	<b>1.0641±0.004</b>	<b>0.6837±0.005</b>	<b>0.9402±0.005</b>	<b>0.7925±0.006</b>	<b>0.8923±0.005</b>	<b>0.7853±0.005</b>
CELEBA	CNP	1.0323±0.016	0.7845±0.013	1.0177±0.016	0.7438±0.017	0.8956±0.009	0.7344±0.011	0.7677±0.012	0.6096±0.009
	BCNP	1.0452±0.009	0.8015±0.008	1.0275±0.009	0.7726±0.008	0.9351±0.006	0.8376±0.009	0.8015±0.010	0.6816±0.008
	SCNP	<b>1.0525±0.008</b>	<b>0.8243±0.006</b>	<b>1.0348±0.008</b>	<b>0.7868±0.006</b>	<b>0.9562±0.004</b>	<b>0.8545±0.008</b>	<b>0.8344±0.006</b>	<b>0.7045±0.005</b>
	NP	1.1333±0.004	0.8766±0.005	1.1043±0.015	0.8355±0.015	1.0034±0.008	0.8456±0.006	0.8935±0.009	0.6893±0.006
	BNP	1.1732±0.005	0.8901±0.006	1.1378±0.007	0.8678±0.006	1.0411±0.008	0.8711±0.005	0.9256±0.008	0.7671±0.006
	SNP	<b>1.1952±0.005</b>	<b>0.9062±0.006</b>	<b>1.1565±0.005</b>	<b>0.8846±0.005</b>	<b>1.0542±0.006</b>	<b>0.8956±0.004</b>	<b>0.9425±0.006</b>	<b>0.7985±0.005</b>
	ANP	1.1633±0.002	1.0163±0.004	1.1377±0.004	0.9866±0.006	1.0418±0.004	0.8845±0.006	0.9363±0.004	0.7346±0.008
	BANP	1.1751±0.002	1.0389±0.005	1.1488±0.004	1.0155±0.005	1.0602±0.005	0.9255±0.006	0.9489±0.004	0.8415±0.007
	SANP	<b>1.1854±0.000</b>	<b>1.0594±0.004</b>	<b>1.1685±0.002</b>	<b>1.0353±0.004</b>	<b>1.0772±0.002</b>	<b>0.9455±0.004</b>	<b>0.9655±0.003</b>	<b>0.8673±0.005</b>
	ConvCNP	1.1697±0.004	1.0366±0.004	1.1445±0.006	0.9947±0.006	1.0542±0.005	0.8971±0.007	0.9521±0.005	0.7563±0.006
	BConvCNP	1.1822±0.004	1.0467±0.004	1.1511±0.005	1.0271±0.005	1.0686±0.006	0.9317±0.005	0.9542±0.005	0.8527±0.005
	SConvCNP	<b>1.1889±0.003</b>	<b>1.0615±0.005</b>	<b>1.1753±0.004</b>	<b>1.0478±0.004</b>	<b>1.0752±0.004</b>	<b>0.9485±0.005</b>	<b>0.9693±0.007</b>	<b>0.8714±0.006</b>
	ConvNP	1.1767±0.004	1.0451±0.004	1.1485±0.005	1.0156±0.006	1.0594±0.005	0.9066±0.007	0.9573±0.004	0.7779±0.005
	BConvNP	1.1822±0.004	1.0555±0.004	1.1573±0.004	1.0351±0.005	1.0762±0.005	0.9461±0.003	0.9661±0.006	0.8636±0.003
	SConvNP	<b>1.1867±0.003</b>	<b>1.0668±0.004</b>	<b>1.1846±0.004</b>	<b>1.0487±0.005</b>	<b>1.0863±0.004</b>	<b>0.9511±0.005</b>	<b>0.9743±0.006</b>	<b>0.8836±0.003</b>

Table 4: Bayesian optimization experiments on data generated by different GP kernels.

Method	ANP	BANP	SANP	ConvCNP	BConvCNP	SConvCNP	ConvNP	BConvNP	SConvNP
RBF	0.1245±0.003	0.1341±0.003	<b>0.1142±0.002</b>	0.1215±0.002	0.1168±0.003	<b>0.1037±0.002</b>	0.1197±0.002	0.1156±0.003	<b>0.1015±0.003</b>
Matern	0.1518±0.003	0.1316±0.004	<b>0.1201±0.002</b>	0.1489±0.003	0.1301±0.002	<b>0.1216±0.002</b>	0.1446±0.002	0.1242±0.003	<b>0.1204±0.004</b>
Periodic	0.1892±0.002	0.1788±0.005	<b>0.1672±0.001</b>	0.1652±0.002	0.1526±0.004	<b>0.1487±0.003</b>	0.1611±0.002	0.1498±0.004	0.1446±0.002

controllers to successfully swing up and balance the pendulum. More detailed experimental settings are given in the Appendix.

For each configuration of the simulator including training and testing environments, we sampled 400 trajectories of the horizon as 10 steps using a random controller. During the testing process, 100 state transition pairs were randomly selected for each configuration of the environment, working as the maximum context points to identify the configuration of dynamics. Figure 2 shows the predictive Log-Likelihoods (LL) and Mean Average Error (MAE) on Cart-Pole State Transition Testing Dataset. We can see that our stable NPs achieve better likelihood and lower prediction error than the original ones. The variances on the stable one are consistently smaller than all original baselines.

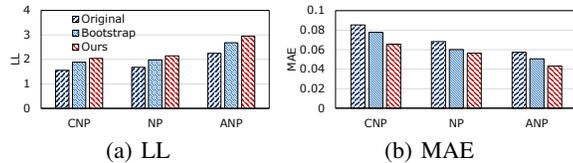


Figure 2: The predictive Log Likelihood (LL) and Mean Average Error (MAE) on Cart-Pole state transition testing dataset.

We can see that our stable NPs achieve better likelihood and lower prediction error than the original ones. The variances on the stable one are consistently smaller than all original baselines.

### 6.4 Bayesian Optimization

Following the setting in BANP [18], we conducted the Bayesian optimization experiment. Taking GP data with RBF, Matern, and Periodic prior functions as examples, we gave the results of ANP, ConvCNP, ConvNP, Bootstrapping versions, and our stable versions. To maintain consistent comparison, we standardized the initializations and normalized the results. We reported the best simple regret, which represents the difference between the current best observation and the global optimum. As shown in Table 4, we can see that our stable solutions consistently achieve lower regret than other NPs.

Table 5: Predator-prey model results.

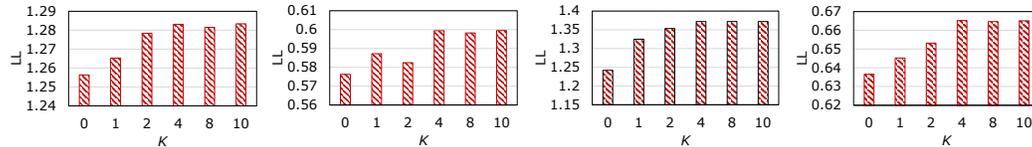
Method	ANP	BANP	SANP	ConvCNP	BConvCNP	SConvCNP	ConvNP	BConvNP	SConvNP
Simulated-context	2.5801 $\pm$ 0.003	2.5912 $\pm$ 0.002	<b>2.6127</b> $\pm$ 0.004	2.5912 $\pm$ 0.002	2.6068 $\pm$ 0.003	<b>2.6164</b> $\pm$ 0.002	2.5925 $\pm$ 0.002	2.6125 $\pm$ 0.005	<b>2.6231</b> $\pm$ 0.003
Simulated-target	1.8265 $\pm$ 0.003	1.8635 $\pm$ 0.004	<b>1.8844</b> $\pm$ 0.002	1.8352 $\pm$ 0.004	1.8524 $\pm$ 0.003	<b>1.9516</b> $\pm$ 0.002	1.9228 $\pm$ 0.005	1.9442 $\pm$ 0.003	<b>1.9662</b> $\pm$ 0.004
Real-context	1.7234 $\pm$ 0.002	1.8496 $\pm$ 0.005	<b>1.8412</b> $\pm$ 0.001	1.7956 $\pm$ 0.002	1.8026 $\pm$ 0.004	<b>1.8744</b> $\pm$ 0.003	1.8342 $\pm$ 0.002	1.8476 $\pm$ 0.005	<b>1.8796</b> $\pm$ 0.002
Real-target	-7.8042 $\pm$ 0.002	-5.4836 $\pm$ 0.004	<b>-5.2527</b> $\pm$ 0.001	-5.3414 $\pm$ 0.005	-5.1526 $\pm$ 0.004	<b>-5.3513</b> $\pm$ 0.004	-5.3155 $\pm$ 0.002	-5.2615 $\pm$ 0.004	<b>-5.2145</b> $\pm$ 0.003

### 6.5 Predator-prey Models

Following [18] and [11], we consider the Lotka–Volterra model [30], which is used to describe the evolution of predator–prey populations. We first trained the models using simulated data generated from a Lotka–Volterra model and tested them on real-world data (Hudson’s Bay hare-lynx data), which is quite different from the simulated data and can be considered as a mismatch scenario. Table 5 lists the results on both simulated and real data. Similar to the previous observation, our stable version still outperforms the original version. Among stable versions, SConvNP achieves the best performance.

### 6.6 Ablation Study

The key parameter in our stable solution is the number of hard predictable subsets  $K$ . Taking SANP as an example, we investigated the average log-likelihood in terms of different  $K$  on the 1D regression task, as shown in Figure 3. We can see that SANP performs better as  $K$  increases, reaches the best value at around  $K = 4$ , and then becomes stable in performance as  $K$  grows larger. As proved in Theorem 5.3, optimization on  $\mathcal{D}_{1:M}$  and more than one hard predictable subset of  $\mathcal{D}_{1:M}$  can achieve more stable prediction. We also conducted different experiments to explore the impact of different  $w_k$ . Taking the 1D regression task with RBF-GP data as an example, we set  $K = 3$  and different  $w_k$  for experiments, as shown in Table 6. It can be seen from the table below that when different weights are set, the model using a stable strategy is better than the original model, and when the weight is set to be equal, its performance is optimal. In addition, when the value of  $w_k (k \geq 1)$  is significantly different from  $w_0$ , such as  $(0.625, 0.125, 0.125, 0.125)$  and  $(0.0625, 0.3125, 0.3125, 0.3125)$ , its performance is more significantly reduced compared to  $(0.25, 0.25, 0.25, 0.25)$ , but it still has a significant improvement compared to the original non-stable model.



(a) RBF/context (b) RBF/Target (c) Matern/context (d) Matern/target  
Figure 3: Log-likelihood (SANP) comparisons with different  $K$  on 1D regression task.

Table 6: Log-likelihood comparisons with different  $w_k$  on 1D regression task.

Weights	(0.25,0.25,0.25,0.25)	(0.4,0.2,0.2,0.2)	(0.625,0.125,0.125,0.125)	(0.1,0.3,0.3,0.3)	(0.0625,0.3125,0.3125,0.3125)
SCNP	(0.9255, 0.4125)	(0.9127, 0.4019)	(0.9035, 0.3998)	(0.9149, 0.4086)	(0.8927, 0.3991)
SNP	(0.8955, 0.3925)	(0.8737, 0.3817)	(0.8716, 0.3809)	(0.8831, 0.3859)	(0.8657, 0.3775)
SANP	(1.2831, 0.5215)	(1.2776, 0.5187)	(1.2738, 0.5196)	(1.2791, 0.5203)	(1.2712, 0.5193)
SConvCNP	(1.3991, 0.5996)	(1.3916, 0.5933)	(1.3841, 0.5915)	(1.3934, 0.5946)	(1.3854, 0.5919)
SConvNP	(1.4036, 0.6015)	(1.3991, 0.6004)	(1.3931, 0.5992)	(1.4012, 0.6015)	(1.3957, 0.5995)

## 7 Conclusion and Future Work

In this paper, we provided theoretical guidelines for deriving stable solutions for NPs, which can obtain good generalization performance. Experiments demonstrated the proposed stable solution can help NPs to achieve more accurate and stable predictions. Although the theoretical analysis we give is based on regression models, it is still open to question whether this conclusion is appropriate for classification models. Therefore, we are interested in extending our theory, expecting it to apply to more different types of tasks.

## Acknowledgement

This work was partly supported by the National Natural Science Foundation of China under Grant 62176020; the National Key Research and Development Program (2020AAA0106800); the Joint Foundation of the Ministry of Education (8091B042235); the Beijing Natural Science Foundation under Grant L211016; the Fundamental Research Funds for the Central Universities (2019JBZ110); and Chinese Academy of Sciences (OEIP-O-202004).

## References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [2] Olivier Bousquet and André Elisseeff. Algorithmic stability and generalization performance. *Advances in Neural Information Processing Systems*, pages 196–202, 2001.
- [3] Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2:499–526, 2002.
- [4] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 2921–2926. IEEE, 2017.
- [5] Alec Farid and Anirudha Majumdar. Generalization bounds for meta-learning via pac-bayes and uniform stability. *Advances in Neural Information Processing Systems*, 34:2173–2186, 2021.
- [6] Andrew Foong, Wessel Bruinsma, Jonathan Gordon, Yann Dubois, James Requeima, and Richard Turner. Meta-learning stationary stochastic process prediction with convolutional neural processes. *Advances in Neural Information Processing Systems*, 33:8284–8295, 2020.
- [7] Yarin Gal, Rowan McAllister, and Carl Edward Rasmussen. Improving pilco with bayesian neural network dynamics models. In *Data-Efficient Machine Learning workshop, ICML*, volume 4, page 25, 2016.
- [8] Alexandre Galashov, Jonathan Schwarz, Hyunjik Kim, Marta Garnelo, David Saxton, Pushmeet Kohli, SM Eslami, and Yee Whye Teh. Meta-learning surrogate models for sequential decision making. *arXiv preprint arXiv:1903.11907*, 2019.
- [9] Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *International Conference on Machine Learning*, pages 1704–1713. PMLR, 2018.
- [10] Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018.
- [11] Jonathan Gordon, Wessel P Bruinsma, Andrew YK Foong, James Requeima, Yann Dubois, and Richard E Turner. Convolutional conditional neural processes. *arXiv preprint arXiv:1910.13556*, 2019.
- [12] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, pages 1225–1234. PMLR, 2016.
- [13] Makoto Kawano, Wataru Kumagai, Akiyoshi Sannai, Yusuke Iwasawa, and Yutaka Matsuo. Group equivariant conditional neural processes. *arXiv preprint arXiv:2102.08759*, 2021.
- [14] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. Attentive neural processes. *arXiv preprint arXiv:1901.05761*, 2019.
- [15] Mingyu Kim, Kyeongryeol Go, and Se-Young Yun. Neural processes with stochastic attention: Paying more attention to the context dataset. In *International Conference on Learning Representations*, 2022.

- [16] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [17] Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning using calibrated regression. In *International Conference on Machine Learning*, pages 2796–2804. PMLR, 2018.
- [18] Juho Lee, Yoonho Lee, Jungtaek Kim, Eunho Yang, Sung Ju Hwang, and Yee Whye Teh. Bootstrapping neural processes. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6606–6615. Curran Associates, Inc., 2020.
- [19] Dongsheng Li, Chao Chen, Qin Lv, Junchi Yan, Li Shang, and Stephen Chu. Low-rank matrix approximation with stability. In *International Conference on Machine Learning*, pages 295–303. PMLR, 2016.
- [20] Xixun Lin, Jia Wu, Chuan Zhou, Shirui Pan, Yanan Cao, and Bin Wang. Task-adaptive neural process for user cold-start recommendation. In *Proceedings of the Web Conference*, pages 1306–1316, 2021.
- [21] Huafeng Liu, Liping Jing, Dahai Yu, Mingjie Zhou, and Michael Ng. Learning intrinsic and extrinsic intentions for cold-start recommendation with neural stochastic processes. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 491–500, 2022.
- [22] Tongliang Liu, Dacheng Tao, Mingli Song, and Stephen J Maybank. Algorithm-dependent generalization bounds for multi-task learning. *IEEE transactions on pattern analysis and machine intelligence*, 39(2):227–241, 2016.
- [23] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [24] Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81):1–32, 2016.
- [25] Carl Edward Rasmussen. Gaussian processes in machine learning. In *Summer school on machine learning*, pages 63–71. Springer, 2003.
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [27] Qi Wang and Herke Van Hoof. Doubly stochastic variational inference for neural processes with hierarchical latent variables. In *International Conference on Machine Learning*, pages 10018–10028. PMLR, 2020.
- [28] Qi Wang and Herke van Hoof. Learning expressive meta-representations with mixture of expert neural processes. In *Advances in neural information processing systems*, 2022.
- [29] Andre Wibisono, Lorenzo Rosasco, and Tomaso Poggio. Sufficient conditions for uniform stability of regularization algorithms. 2009.
- [30] Darren J Wilkinson. *Stochastic modelling for systems biology*. CRC press, 2018.
- [31] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov, and Alexander Smola. Deep sets. *arXiv preprint arXiv:1703.06114*, 2017.
- [32] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

## A Inductive biases

Here, we revisit some properties, which would help us understand NPs. First, we give a concept of Permutation Invariant Function which is the basic property of stochastic process, e.g., NPs.

**Definition A.1.** (Permutation Invariant Function) A function  $f(\cdot) : \times_i^N \mathbb{R}^D \rightarrow \mathbb{R}^d$  mapping a set of data points  $\{\mathbf{x}_i\}_{i=1}^N$  is Permutation Invariant Function if

$$\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \rightarrow f = [f_1(\mathbf{x}_{\pi(1:N)}), f_2(\mathbf{x}_{\pi(1:N)}), \dots, f_d(\mathbf{x}_{\pi(1:N)})], \quad (13)$$

where  $x_i \in \mathbb{R}^D$  and the function output is a  $d$  dimensional vector. Operation  $\pi : [1, 2, \dots, N] \rightarrow [\pi_1, \pi_2, \dots, \pi_N]$  is a permutation set over the order of elements in the set.

**Definition A.2.** (Permutation Equivariant Function) A function  $f(\cdot) : \times_i^N \mathbb{R}^D \rightarrow \mathbb{R}^N$  mapping a set of data points  $\{\mathbf{x}_i\}_{i=1}^N$  is Permutation Invariant Function if

$$\mathbf{X}_\pi = [\mathbf{x}_{\pi_1}, \mathbf{x}_{\pi_2}, \dots, \mathbf{x}_{\pi_N}] \rightarrow f_\pi = \pi \circ f(\mathbf{x}_{1:N}), \quad (14)$$

where the function output contains  $N$  elements keeping the order of inputs.

Permutation Equivariant Function keeps the order of elements in the output consistent with that in the input under any permutation operation  $\pi$ . Permutation invariant functions are candidate functions for learning embeddings of a set or other order uncorrelated data structure  $\{\mathbf{x}_i\}_{i=1}^N$ , and the invariant property is easy to be verified. Here, we give a mean operation structure over the output

$$F(\mathbf{X}_{\pi(1:N)}) = \left( \frac{1}{N} \sum_{i=1}^N \phi_1(\mathbf{x}_i), \frac{1}{N} \sum_{i=1}^N \phi_2(\mathbf{x}_i), \dots, \frac{1}{N} \sum_{i=1}^N \phi_D(\mathbf{x}_i) \right) \quad (15)$$

## B Model Architecture

We show the architectural details of the CNP, NP, and ANP models used for the 1D and 2D function regression experiments. The neural process aims to learn a stochastic process (random function) mapping target features  $\mathbf{x}_i$  to prediction  $y_i$  given the context set  $\mathcal{D}^c$  as training data (a realization from the stochastic process), i.e., learning

$$\log P(\mathbf{y}^T | \mathbf{X}^T, \mathcal{D}^c) = \log P(\mathbf{y} | \mathbf{X}, \mathcal{D}^c) = \sum_{i=1}^N P(y_i | \mathbf{x}_i, \mathcal{D}^c). \quad (16)$$

Conditional neural process (CNP) [9] describes  $P(y_i | \mathbf{x}_i, \mathcal{D}^c)$  with a deterministic neural network taking  $\mathcal{D}^c$  to output the parameters of  $P(y_i | \mathbf{x}_i, \mathcal{D}^c)$ . CNP consists of an encoder  $f_{\text{enc}}(\cdot)$ , an aggregator  $f_{\text{agg}}(\cdot)$  and a decoder  $f_{\text{dec}}(\cdot)$ ; the encoder summarizes  $\mathcal{D}^c$  and  $\mathbf{x}_i$  into latent representations  $[\mathbf{r}_1, \dots, \mathbf{r}_{|\mathcal{C}|}] \in \mathbb{R}^{|\mathcal{C}| \times d}$  via permutation-invariant neural network [31], where  $d$  is the number of latent dimensions, and aggregator summarizes the encoded context features to a single representation  $\mathbf{r}^c$ , and decoder takes as input the aggregated representations  $\mathbf{r}^c$  and  $\mathbf{x}_i$  and output the single output-specific mean  $\mu_i$  and variance  $\sigma_i^2$  for the corresponding value of  $y_i$ .

$$\begin{aligned} \mathbf{r}_i &= f_{\text{enc}}(\mathbf{x}_i, y_i), \quad i \in \mathcal{C} \\ \mathbf{r}^c &= \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \mathbf{r}_i \\ \phi &= f_{\text{agg}}(\mathbf{r}^c) \\ (\mu_i, \sigma_i) &= f_{\text{dec}}(\phi, \mathbf{x}_i), \quad p(y_i | \mathbf{x}_i, \mathcal{D}^c) = \mathcal{N}(y_i; \mu_i, \sigma_i^2) \quad i \in \mathcal{T} \end{aligned} \quad (17)$$

where  $f_{\text{enc}}(\cdot)$  and  $f_{\text{dec}}(\cdot)$  are feed-forward neural networks. The decoder output  $\mu_i$  and variance  $\sigma_i^2$  are predicted mean and variance. We use *Gaussian* distribution  $\mathcal{N}(y_i; \mu_i, \sigma_i^2)$  as predictive distribution. CNP is trained to maximize the expected likelihood  $\mathbb{E}_{P(\mathcal{T})}[P(y_i | \mathbf{x}_i, \mathcal{D}^c)]$ .

Neural process [10] further models functional uncertainty using a global latent variable. Unlike CNP, which maps a context into a deterministic representation  $\tilde{\mathbf{r}}_i$ , NP encodes a context into a *Gaussian* latent variable  $\mathbf{z}$ , giving additional stochasticity in function construction. Following [14], we consider

an NP with both a deterministic path and latent path, where the deterministic path models the overall skeleton of the function  $\tilde{\mathbf{r}}_i$ , and the latent path models the functional uncertainty:

$$\begin{aligned}
 \mathbf{r}_i &= f_{\text{enc}}^{(1)}(\mathbf{x}_i, y_i), \quad i \in \mathcal{C} \\
 \mathbf{r}^{\mathcal{C}} &= \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \mathbf{r}_i \\
 \phi &= f_{\text{agg}}(\mathbf{r}) \\
 (\mu_z, \sigma_z) &= f_{\text{enc}}^{(2)}(\mathcal{D}^{\mathcal{C}}), \quad q(\mathbf{z}|\mathcal{D}^{\mathcal{C}}) = \mathcal{N}(\mathbf{z}; \mu_z, \sigma_z^2) \\
 (\mu_i, \sigma_i) &= f_{\text{dec}}(\phi, \mathbf{z}, \mathbf{x}_i), \quad p(y_i|\mathbf{x}_i, \mathbf{z}, \mathcal{D}^{\mathcal{C}}) = \mathcal{N}(y_i; \mu_i, \sigma_i^2) \quad i \in \mathcal{T}
 \end{aligned} \tag{18}$$

with  $f_{\text{enc}}^{(1)}(\cdot)$  and  $f_{\text{enc}}^{(2)}(\cdot)$  having the same structure as  $f_{\text{enc}}(\cdot)$  in Eq.(17). In this scenario, the conditional distribution is lower bounded as:

$$\log P(\mathbf{y}|\mathbf{X}, \mathcal{D}^{\mathcal{C}}) \geq \sum_{i=1}^N \mathbb{E}_{q(\mathbf{z}|\mathcal{D}^{\mathcal{C}})} \left[ \log \frac{P(y_i|\mathbf{x}_i, \mathbf{z}, \mathcal{D}^{\mathcal{C}}) P(\mathbf{z}|\mathcal{D}^{\mathcal{C}})}{q(\mathbf{z}|\mathbf{X}, \mathbf{y})} \right]. \tag{19}$$

We further approximate  $q(\mathbf{z}|\mathcal{D}^{\mathcal{C}}) \approx P(\mathbf{z}|\mathcal{D}^{\mathcal{C}})$  and train the model by maximizing this expected lower bound over tasks. Furthermore, ANP introduces attention mechanisms into NP to resolve the issue of under-fitting.

The architectural details of the CNP, NP, and ANP are the same as in [14]. Here we give the detailed architectures of the encoder and decoder of NPs.

### B.1 Encoder without attention

Encoder focuses on learning embeddings for each data point in the context set, and the basic component is multi-layer perceptron, which is defined by

$$\text{MLP}(l, d_{in}, d_h, d_{out}) = \text{LINEAR}(d_h, d_{out}) \circ \underbrace{(\text{RELU} \circ \text{LINEAR}(d_h, d_h) \circ \dots)}_{\times(l-2)} \circ \text{LINEAR}(d_h, d_{in}) \tag{20}$$

where  $l$  is the number of layers,  $d_{in}$ ,  $d_h$  and  $d_{out}$  are dimensionalities of inputs, hidden units and outputs. Here  $\text{RELU}(\cdot)$  is adapted as activation function.

The encoder in Vanilla CNP uses a deterministic encoder which focuses on learning embeddings for each data point in context set.

$$\begin{aligned}
 \mathbf{r}_i &= \text{MLP}(l_{e1}, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i]), \\
 \mathbf{r}^{\mathcal{C}} &= \sum_{i \in \mathcal{C}} \mathbf{r}_i, \quad \phi = \text{MLP}(l_{e2}, d_h, d_h)(\mathbf{r}^{\mathcal{C}})
 \end{aligned} \tag{21}$$

where  $d_x$  and  $d_y$  are the dimensionalities of  $\mathbf{x}_i$  and  $y_i$ .

To follow the encoder structure in NP, we introduce another encoder aligned with original deterministic encoder to permit the same number parameters, i.e.,

$$\begin{aligned}
 \mathbf{r}_i^{(1)} &= \text{MLP}(l_{e1}, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i]) \\
 \mathbf{r}_{\mathcal{C}}^{(1)} &= \sum_{i \in \mathcal{C}} \mathbf{r}_i^{(1)}, \quad \phi_1 = \text{MLP}(l_{e2}, d_h, d_h)(\mathbf{r}_{\mathcal{C}}^{(1)}) \\
 \mathbf{r}_i^{(2)} &= \text{MLP}(l_{e1}, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i]) \\
 \mathbf{r}_{\mathcal{C}}^{(2)} &= \sum_{i \in \mathcal{C}} \mathbf{r}_i^{(2)}, \quad \phi_2 = \text{MLP}(l_{e2}, d_h, d_h)(\mathbf{r}_{\mathcal{C}}^{(2)}) \\
 \phi &= [\phi_1, \phi_2]
 \end{aligned} \tag{22}$$

The encoder in NP contains a deterministic path and a latent path, i.e.,

$$\begin{aligned}
\mathbf{r}_i^{(1)} &= \text{MLP}(l_{de1}, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i]) \\
\mathbf{r}_C^{(1)} &= \sum_{i \in C} \mathbf{r}_i^{(1)}, \quad \phi = \text{MLP}(l_{de2}, d_h, d_h)(\mathbf{r}_C^{(1)}) \\
\mathbf{r}_i^{(2)} &= \text{MLP}(l_{la1}, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i]) \\
\mathbf{r}_C^{(2)} &= \sum_{i \in C} \mathbf{r}_i^{(2)}, \quad [\mu_z, \sigma'_z] = \text{MLP}(l_{la2}, d_h, d_h)(\mathbf{r}_C^{(2)}) \\
\sigma_z &= 0.1 + 0.9 \cdot \text{SIGMOID}(\sigma'_z), \quad \mathbf{z} \sim \mathcal{N}(\mu_z, \text{diag}(\sigma_z^2)).
\end{aligned} \tag{23}$$

In this case, the encoder outputs deterministic representation  $\phi$  and latent representation  $\mathbf{z}$ .

## B.2 Encoder with attention

The attention mechanism is widely used in NPs, Specifically, multi-head attention [26] is adapted, which is defined by

$$\begin{aligned}
\mathbf{Q}' &= \{\text{LINEAR}(d_q, d_{out})(\mathbf{q})\}_{\mathbf{q} \in \mathbf{Q}}, \quad \{\mathbf{Q}'_i\}_{i=1}^{n_{head}} = \text{SPLIT}(\mathbf{Q}', n_{head}), \\
\mathbf{K}' &= \{\text{LINEAR}(d_k, d_{out})(\mathbf{k})\}_{\mathbf{k} \in \mathbf{K}}, \quad \{\mathbf{K}'_i\}_{i=1}^{n_{head}} = \text{SPLIT}(\mathbf{K}', n_{head}), \\
\mathbf{V}' &= \{\text{LINEAR}(d_v, d_{out})(\mathbf{v})\}_{\mathbf{v} \in \mathbf{V}}, \quad \{\mathbf{V}'_i\}_{i=1}^{n_{head}} = \text{SPLIT}(\mathbf{V}', n_{head}), \\
\mathbf{H}_i &= \text{SOFTMAX}\left(\mathbf{Q}'_i(\mathbf{K}'_i)^\top / \sqrt{d_{out}}\right) \mathbf{V}'_i, \quad \mathbf{H} = \text{CONCAT}(\{\mathbf{H}_i\}_{i=1}^{n_{head}}) \\
\mathbf{H}' &= \text{LAYERNORM}(\mathbf{Q}' + \mathbf{H}) \\
\text{MHA}(d_{out})(\mathbf{Q}, \mathbf{K}, \mathbf{V}) &= \text{LAYERNORM}(\mathbf{H}' + \text{RELU}(\text{LINEAR}(d_{out}, d_{out})))
\end{aligned} \tag{24}$$

where  $d_q, d_v, d_k$  are the dimensionalities of query  $\mathbf{Q}$ , key  $\mathbf{K}$ , and value  $\mathbf{V}$ , respectively.  $n_{head}$  is the number of head. Here Layer normalization [1]  $\text{LAYERNORM}(\cdot)$  is adapted. It is easy to derive self-attention by setting  $\mathbf{Q} = \mathbf{K} = \mathbf{V}$ , i.e.,

$$\text{SA}(d_{out})(\mathbf{X}) = \text{MHA}(d_{out})(\mathbf{X}, \mathbf{X}, \mathbf{X}) \tag{25}$$

For CNP, the encoder with attention still contains two deterministic paths,

$$\begin{aligned}
f_{qk} &= \text{MLP}(l_{qk}, d_x, d_h, d_h) \\
\mathbf{Q} &= f_{qk}(\mathbf{x}_i), \quad i \in \mathcal{T} \\
\mathbf{K} &= \{f_{qk}(\mathbf{x}_i)\}, \quad i \in \mathcal{C} \\
\mathbf{V} &= \text{SA}(d_h) \cdot \{\text{MLP}(l_v, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i])\}_{i \in \mathcal{C}} \\
\phi_1 &= \text{MHA}(d_h)(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \\
\mathbf{H} &= \text{SA}(d_h) \cdot \{\text{RELU} \circ \text{MLP}(l_{e1}, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i])\}_{i \in \mathcal{C}} \\
\phi_2 &= \text{MLP}(l_e, d_h, d_h) \left( \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \mathbf{h}_i \right) \\
\phi &= [\phi_1, \phi_2]
\end{aligned} \tag{26}$$

Similarly, encoder with attention in NP contains a deterministic path and a latent path, i.e.,

$$\begin{aligned}
f_{qk} &= \text{MLP}(l_{qk}, d_x, d_h, d_h) \\
\mathbf{Q} &= f_{qk}(\mathbf{x}_i), \quad i \in \mathcal{T} \\
\mathbf{K} &= \{f_{qk}(\mathbf{x}_i)\}, \quad i \in \mathcal{C} \\
\mathbf{V} &= \text{SA}(d_h) \cdot \{\text{MLP}(l_v, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i])\}_{i \in \mathcal{C}} \\
\phi &= \text{MHA}(d_h)(\mathbf{Q}, \mathbf{K}, \mathbf{V})
\end{aligned} \tag{27}$$

and

$$\begin{aligned}
\mathbf{H} &= \text{SA}(d_h) (\{\text{RELU} \circ \text{MLP}(l_{e1}, d_x + d_y, d_h, d_h)([\mathbf{x}_i, y_i])\}_{i \in \mathcal{C}}) \\
[\mu_z, \sigma'_z] &= \text{MLP}(l_{la}, d_h, d_h) \left( \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \mathbf{h}_i \right) \\
\sigma_z &= 0.1 + 0.9 \cdot \text{SIGMOID}(\sigma'_z), \\
\mathbf{z} &\sim \mathcal{N}(\mu_z, \text{diag}(\sigma_z^2)).
\end{aligned} \tag{28}$$

### B.3 Decoder

The decoder focuses on predicting output for target points based on the encoder’s outputs  $\phi$ . For target point  $\{\mathbf{x}_i\}_{i \in \mathcal{T}}$ , the decoder of CNP is defined by

$$\begin{aligned}
[\mu_i, \sigma'_i] &= \text{MLP}(d_{dec}, 2d_h + d_x, d_h, 2d_y)[\phi, \mathbf{x}_i], \quad i \in \mathcal{T} \\
\sigma_i &= 0.1 + 0.9 \cdot \text{SOFTPLUS}(\sigma'_i) \\
y_i &\sim \mathcal{N}(\mu_i, \sigma_i)
\end{aligned} \tag{29}$$

Decoder of NP is defined by

$$\begin{aligned}
[\mu_i, \sigma'_i] &= \text{MLP}(d_{dec}, d_h + d_z + d_x, d_h, 2d_y)[\phi, \mathbf{x}_i, \mathbf{z}], \quad i \in \mathcal{T} \\
\sigma_i &= 0.1 + 0.9 \cdot \text{SOFTPLUS}(\sigma'_i) \\
y_i &\sim \mathcal{N}(\mu_i, \sigma_i)
\end{aligned} \tag{30}$$

## C Implementation Details and Experiments

For CNP [9], BCNP [18], and our SCNP, we apply the encoder with attention described in Eq (26) and decoder described in Eq (29). For NP [10], ANP [14], BNP [18], BANP [18] and our SNP and SANP models, we apply encoder with attention described in Eq (27) and (28), and decoder described in Eq (30).

### C.1 1D Regression

For synthetic 1D regression experiments, the neural architectures for CNP, NP, ANP, BCNP, BNP, BANP, and our SCNP/SNP/SANP refer to Appendix B. The number of hidden units is  $d_h = 128$  and latent representation  $d_z = 128$ . The number of layers are  $l_e = l_{de} = l_{la} = l_{qk} = l_v = 2$ .

We generate datasets for synthetic 1D regression. Specifically, the stochastic process (SP) initializes with a 0 mean Gaussian Process (GP)  $y^{(0)} \sim GP(0, k(\cdot, \cdot))$  indexed in the interval  $x \in [-2.0, 2.0]$ , where the radial basis function kernel  $k(x, x') = \sigma^2 \exp(-\|x - x'\|^2/2l^2)$  with  $s \sim U(0.1, 1.0)$  and  $\sigma \sim U(0.1, 0.6)$ . Furthermore, GP with Matern Kernel is adopted for model-data mismatch scenario, which is defined as  $k(x, x') = \sigma^2(1 + \sqrt{5}d/l + 5d^2/(3l^2)) \exp(-\sqrt{5}d/l)$  and  $d = \|x - x'\|$  with  $s \sim U(0.1, 1.0)$  and  $\sigma \sim U(0.1, 0.6)$ . For a fair comparison, we set the same data generation, training, and testing for all models.

We trained all models for 100,000 steps with each step computing updates with a batch containing 100 tasks. We used the Adam optimizer with an initial learning rate  $5 \cdot 10^{-4}$  and decayed the learning rate using Cosine annealing scheme for baselines. For SCNP/SNP/SANP, we set  $K = 3$ . The size of the context  $\mathcal{C}$  was drawn as  $|\mathcal{C}| \sim U(3, 200)$ . Testings were done for 3,000 batches with each batch containing 16 tasks (48,000 tasks in total).

We investigate the model stability from the size of the context set and the percent of added noise. First, we conduct experiments on different size of context set, i.e.,  $|\mathcal{C}| \in \{20, 50, 100, 200\}$ . Table 7 shows the Average Log-likelihoods performance comparison between different methods in terms of different context size. We can see that the performance becomes better with the increasing of  $|\mathcal{C}|$  and NPs with stable solution still achieve better performance. Second, we investigate the model performance in terms of different noise setting. Here we introduce Gaussian noise  $\mathcal{N}(0, 1)$  and add noise to different proportions of the data, such as  $\{0\%, 5\%, 10\%, 15\%\}$ . Table 8 lists the Average Log-likelihoods performance comparison in terms of different noise proportions.

Table 7: Average Log-likelihoods over all context and target points on realizations from Synthetic Stochastic Process on different size of context set. (Mean  $\pm$  Std).

Kernel	Method	20		50		100		200	
		context	target	context	target	context	target	context	target
RBF	CNP	0.8724 $\pm$ 0.008	0.4334 $\pm$ 0.007	0.9533 $\pm$ 0.005	0.4854 $\pm$ 0.010	1.1224 $\pm$ 0.005	0.5322 $\pm$ 0.007	0.1563 $\pm$ 0.004	0.5783 $\pm$ 0.004
	BCNP	0.9015 $\pm$ 0.009	0.4579 $\pm$ 0.007	0.9787 $\pm$ 0.007	0.5215 $\pm$ 0.008	1.1687 $\pm$ 0.007	0.5716 $\pm$ 0.009	0.1985 $\pm$ 0.005	0.6086 $\pm$ 0.005
	SCNP	<b>0.9255</b> $\pm$ 0.008	<b>0.4733</b> $\pm$ 0.004	<b>0.9944</b> $\pm$ 0.005	<b>0.5433</b> $\pm$ 0.006	<b>1.1833</b> $\pm$ 0.006	<b>0.5918</b> $\pm$ 0.005	<b>0.2111</b> $\pm$ 0.004	<b>0.6333</b> $\pm$ 0.004
	NP	0.8215 $\pm$ 0.004	0.3853 $\pm$ 0.005	0.9124 $\pm$ 0.006	0.4234 $\pm$ 0.003	1.0855 $\pm$ 0.003	0.4767 $\pm$ 0.005	1.1225 $\pm$ 0.002	0.5233 $\pm$ 0.004
	BNP	0.8714 $\pm$ 0.004	0.4122 $\pm$ 0.004	0.9712 $\pm$ 0.004	0.4718 $\pm$ 0.004	1.1426 $\pm$ 0.005	0.5269 $\pm$ 0.006	1.1716 $\pm$ 0.004	0.5698 $\pm$ 0.004
	SNP	<b>0.8955</b> $\pm$ 0.003	<b>0.4356</b> $\pm$ 0.004	<b>0.9866</b> $\pm$ 0.004	<b>0.4934</b> $\pm$ 0.005	<b>1.1637</b> $\pm$ 0.004	<b>0.5434</b> $\pm$ 0.004	<b>1.1958</b> $\pm$ 0.003	<b>0.5933</b> $\pm$ 0.003
Matern	ANP	1.2563 $\pm$ 0.002	0.5763 $\pm$ 0.004	1.3233 $\pm$ 0.002	0.6322 $\pm$ 0.003	1.4633 $\pm$ 0.003	0.6866 $\pm$ 0.006	1.4982 $\pm$ 0.006	0.7322 $\pm$ 0.004
	BANP	1.2715 $\pm$ 0.003	0.5878 $\pm$ 0.005	1.3325 $\pm$ 0.004	0.6465 $\pm$ 0.004	1.4778 $\pm$ 0.003	0.6915 $\pm$ 0.004	1.5115 $\pm$ 0.005	0.7436 $\pm$ 0.005
	SANP	<b>1.2831</b> $\pm$ 0.000	<b>0.5994</b> $\pm$ 0.004	<b>1.3452</b> $\pm$ 0.002	<b>0.6577</b> $\pm$ 0.003	<b>1.4898</b> $\pm$ 0.001	<b>0.7043</b> $\pm$ 0.002	<b>1.5285</b> $\pm$ 0.002	<b>0.7534</b> $\pm$ 0.005
	CNP	0.8531 $\pm$ 0.005	0.2431 $\pm$ 0.010	0.9123 $\pm$ 0.005	0.2984 $\pm$ 0.003	1.0522 $\pm$ 0.004	0.3542 $\pm$ 0.002	1.0984 $\pm$ 0.008	0.4022 $\pm$ 0.005
	BCNP	0.8765 $\pm$ 0.006	0.2788 $\pm$ 0.009	0.9411 $\pm$ 0.006	0.3266 $\pm$ 0.005	1.0752 $\pm$ 0.004	0.3762 $\pm$ 0.004	1.1245 $\pm$ 0.007	0.4326 $\pm$ 0.006
	SCNP	<b>0.8963</b> $\pm$ 0.003	<b>0.2953</b> $\pm$ 0.006	<b>0.9555</b> $\pm$ 0.004	<b>0.3467</b> $\pm$ 0.003	<b>1.0967</b> $\pm$ 0.003	<b>0.3967</b> $\pm$ 0.003	<b>1.1467</b> $\pm$ 0.008	<b>0.4556</b> $\pm$ 0.005
RBF	NP	0.7643 $\pm$ 0.015	0.2041 $\pm$ 0.015	0.8221 $\pm$ 0.002	0.2547 $\pm$ 0.003	0.9322 $\pm$ 0.003	0.3155 $\pm$ 0.004	1.0452 $\pm$ 0.005	0.3563 $\pm$ 0.005
	BNP	0.8052 $\pm$ 0.008	0.2651 $\pm$ 0.007	0.8678 $\pm$ 0.003	0.3163 $\pm$ 0.004	0.9672 $\pm$ 0.003	0.3656 $\pm$ 0.005	1.1052 $\pm$ 0.005	0.4015 $\pm$ 0.005
	SNP	<b>0.8368</b> $\pm$ 0.006	<b>0.2844</b> $\pm$ 0.005	<b>0.8956</b> $\pm$ 0.002	<b>0.3326</b> $\pm$ 0.004	<b>0.9959</b> $\pm$ 0.003	<b>0.3849</b> $\pm$ 0.004	<b>1.1215</b> $\pm$ 0.003	<b>0.4313</b> $\pm$ 0.004
	ANP	1.2421 $\pm$ 0.002	0.6366 $\pm$ 0.004	1.3022 $\pm$ 0.001	0.6881 $\pm$ 0.004	1.4211 $\pm$ 0.004	0.7331 $\pm$ 0.003	1.4631 $\pm$ 0.002	0.7753 $\pm$ 0.008
	BANP	1.3452 $\pm$ 0.007	0.6513 $\pm$ 0.004	1.4056 $\pm$ 0.003	0.7015 $\pm$ 0.004	1.4505 $\pm$ 0.004	0.7531 $\pm$ 0.005	1.4986 $\pm$ 0.004	0.7996 $\pm$ 0.006
	SANP	<b>1.3721</b> $\pm$ 0.002	<b>0.6653</b> $\pm$ 0.004	<b>1.4322</b> $\pm$ 0.003	<b>0.7126</b> $\pm$ 0.004	<b>1.4633</b> $\pm$ 0.004	<b>0.7644</b> $\pm$ 0.003	<b>1.5153</b> $\pm$ 0.005	<b>0.8125</b> $\pm$ 0.004

Table 8: Average Log-likelihoods over all context and target points on realizations from Synthetic Stochastic Process on different percent of added noise. Here we set the context size to 20. (Mean  $\pm$  Std). Note that adding ‘S’ before the original model name is a model with our stable solution.

Kernel	Method	Original		Noise(+5%)		Noise(+10%)		Noise(+15%)	
		context	target	context	target	context	target	context	target
RBF	CNP	0.8724 $\pm$ 0.008	0.4334 $\pm$ 0.007	0.8522 $\pm$ 0.005	0.4001 $\pm$ 0.010	0.8014 $\pm$ 0.006	0.3552 $\pm$ 0.004	0.7152 $\pm$ 0.006	0.2853 $\pm$ 0.005
	BCNP	0.9042 $\pm$ 0.009	0.4589 $\pm$ 0.006	0.8774 $\pm$ 0.006	0.4278 $\pm$ 0.008	0.8316 $\pm$ 0.006	0.3767 $\pm$ 0.005	0.7487 $\pm$ 0.007	0.3017 $\pm$ 0.006
	SCNP	<b>0.9255</b> $\pm$ 0.008	<b>0.4733</b> $\pm$ 0.004	<b>0.8935</b> $\pm$ 0.005	<b>0.4478</b> $\pm$ 0.006	<b>0.8517</b> $\pm$ 0.004	<b>0.3986</b> $\pm$ 0.005	<b>0.7621</b> $\pm$ 0.005	<b>0.3279</b> $\pm$ 0.006
	NP	0.8215 $\pm$ 0.004	0.3853 $\pm$ 0.005	0.8011 $\pm$ 0.004	0.3511 $\pm$ 0.006	0.7611 $\pm$ 0.005	0.3042 $\pm$ 0.008	0.6722 $\pm$ 0.005	0.2435 $\pm$ 0.007
	BNP	0.8722 $\pm$ 0.004	0.4211 $\pm$ 0.004	0.8321 $\pm$ 0.003	0.3876 $\pm$ 0.004	0.7922 $\pm$ 0.004	0.3389 $\pm$ 0.005	0.7189 $\pm$ 0.004	0.2776 $\pm$ 0.007
	SNP	<b>0.8955</b> $\pm$ 0.003	<b>0.4356</b> $\pm$ 0.004	<b>0.8567</b> $\pm$ 0.003	<b>0.4046</b> $\pm$ 0.005	<b>0.8165</b> $\pm$ 0.004	<b>0.3568</b> $\pm$ 0.006	<b>0.7356</b> $\pm$ 0.005	<b>0.2955</b> $\pm$ 0.006
Matern	ANP	1.2563 $\pm$ 0.002	0.5763 $\pm$ 0.004	1.2245 $\pm$ 0.007	0.5347 $\pm$ 0.006	0.1742 $\pm$ 0.005	0.4871 $\pm$ 0.007	0.9821 $\pm$ 0.005	0.4151 $\pm$ 0.004
	BANP	1.2722 $\pm$ 0.004	0.5887 $\pm$ 0.006	1.2411 $\pm$ 0.005	0.5471 $\pm$ 0.005	0.1886 $\pm$ 0.006	0.4917 $\pm$ 0.006	1.0642 $\pm$ 0.006	0.4327 $\pm$ 0.005
	SANP	<b>1.2831</b> $\pm$ 0.000	<b>0.5994</b> $\pm$ 0.004	<b>1.2564</b> $\pm$ 0.004	<b>0.5578</b> $\pm$ 0.004	<b>1.2052</b> $\pm$ 0.004	<b>0.5025</b> $\pm$ 0.006	<b>1.1243</b> $\pm$ 0.005	<b>0.4356</b> $\pm$ 0.006
	CNP	0.8531 $\pm$ 0.005	0.2431 $\pm$ 0.010	0.8231 $\pm$ 0.005	0.2144 $\pm$ 0.010	0.7761 $\pm$ 0.008	0.1784 $\pm$ 0.007	0.7052 $\pm$ 0.005	0.1452 $\pm$ 0.006
	BCNP	0.8778 $\pm$ 0.005	0.2762 $\pm$ 0.009	0.8487 $\pm$ 0.006	0.2477 $\pm$ 0.009	0.8015 $\pm$ 0.007	0.2051 $\pm$ 0.007	0.7378 $\pm$ 0.005	0.1766 $\pm$ 0.006
	SCNP	<b>0.8963</b> $\pm$ 0.003	<b>0.2953</b> $\pm$ 0.006	<b>0.8689</b> $\pm$ 0.005	<b>0.2658</b> $\pm$ 0.007	<b>0.8268</b> $\pm$ 0.006	<b>0.2258</b> $\pm$ 0.006	<b>0.7567</b> $\pm$ 0.004	<b>0.1936</b> $\pm$ 0.005
RBF	NP	0.7643 $\pm$ 0.015	0.2041 $\pm$ 0.015	0.7342 $\pm$ 0.002	0.1725 $\pm$ 0.008	0.6892 $\pm$ 0.004	0.1542 $\pm$ 0.006	0.6235 $\pm$ 0.008	0.1342 $\pm$ 0.007
	BNP	0.8156 $\pm$ 0.005	0.2689 $\pm$ 0.007	0.7789 $\pm$ 0.004	0.2215 $\pm$ 0.005	0.7421 $\pm$ 0.005	0.2117 $\pm$ 0.007	0.6715 $\pm$ 0.006	0.1828 $\pm$ 0.006
	SNP	<b>0.8368</b> $\pm$ 0.006	<b>0.2844</b> $\pm$ 0.005	<b>0.8036</b> $\pm$ 0.003	<b>0.2483</b> $\pm$ 0.003	<b>0.7635</b> $\pm$ 0.004	<b>0.2325</b> $\pm$ 0.006	<b>0.6973</b> $\pm$ 0.006	<b>0.2016</b> $\pm$ 0.005
	ANP	1.2421 $\pm$ 0.002	0.6366 $\pm$ 0.004	1.2115 $\pm$ 0.001	0.6001 $\pm$ 0.008	1.1784 $\pm$ 0.004	0.5622 $\pm$ 0.006	1.1252 $\pm$ 0.007	0.5274 $\pm$ 0.008
	BANP	1.3450 $\pm$ 0.003	0.6514 $\pm$ 0.005	1.3125 $\pm$ 0.005	0.6115 $\pm$ 0.002	1.2672 $\pm$ 0.004	0.5711 $\pm$ 0.005	1.2236 $\pm$ 0.006	0.5306 $\pm$ 0.006
	SANP	<b>1.3721</b> $\pm$ 0.002	<b>0.6653</b> $\pm$ 0.004	<b>1.3461</b> $\pm$ 0.003	<b>0.6256</b> $\pm$ 0.004	<b>1.3011</b> $\pm$ 0.003	<b>0.5782</b> $\pm$ 0.004	<b>1.2457</b> $\pm$ 0.005	<b>0.5356</b> $\pm$ 0.002

## C.2 System Identification on Physics Engines

The second synthetic experiment focuses on evaluating model dynamics on a classical simulator, Cart-Pole systems, which is detailed in [7, 27]. The Cart-Pole swing-up task is a standard benchmark for nonlinear control due to the non-linearity in the dynamics, and the requirement for nonlinear controllers to successfully swing up and balance the pendulum. A pendulum of length  $l$  is attached to a cart by a frictionless pivot. The system begins with the cart at position  $x_c = 0$  and the pendulum hanging down:  $\theta$ . The goal is to accelerate the cart by applying horizontal force  $u_t$  at each time-step  $t$  to invert and then stabilize the pendulum’s endpoint at the goal. There are some parameters that need to be known, such as cart mass  $m_c$ , pendulum mass  $m_p$ , acceleration of gravity  $g = 9.82m/s^2$ , time horizon  $T$ , time discretization  $\Delta t$  and ground friction coefficient  $f_c$ . In this case, the Cart-Pole swing-up task aims to forecast the transited state  $[x_c, \theta, x'_c, \theta']$  in time step  $t + 1$  based on the input as a state action pair  $[x_c, \theta, x'_c, \theta', a]$  in time step  $t$ .

For system identification task on physics engines, the neural architectures for CNP, NP, ANP, BANP and our RNP refer to Appendix B. The number of hidden units is  $d_h = 32$  and latent representation  $d_z = 32$ . The number of layers are  $l_e = l_{de} = l_{la} = l_{qk} = l_v = 2$ .

To generate a variety of trajectories under a random policy for this experiment, the mass  $m_c$  and the ground friction coefficient  $f_c$  are varied in the discrete choices  $m_c \in \{0.3, 0.4, 0.5, 0.6, 0.7\}$  and  $f_c \in \{0.06, 0.08, 0.1, 0.12\}$ . Each pair of  $[m_c, f_c]$  values specifies a dynamics environment, and we formulate all pairs of  $m_c \in \{0.3, 0.5, 0.7\}$  and  $f_c \in \{0.08, 0.12\}$  as training environments with the rest 16 pairs of configurations as the testing environments. For each configuration of the simulator including training and testing environments, we sample 400 trajectories of horizon as 10 steps using a random controller, and more details refer to Supplementary material. During the testing process,

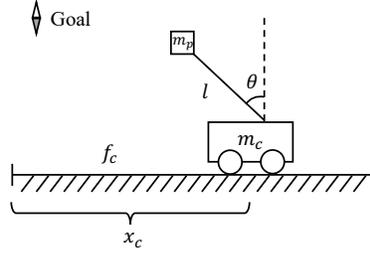


Figure 4: Cart-Pole Dynamical Systems. The cart and the pole are with masses  $m_c$  and  $m_p$ , and the length of the pole is  $l$ . And the configuration of the simulator is up to parameters of the cart-pole mass and the ground friction coefficient here with other hyper-parameters fixed in this experiment.

Table 9: Average Log-likelihoods over all context and target points on EMNIST and CELEBA.

Dataset	Method	Original		Noise(+5%)		Noise(+10%)		Noise(+15%)	
		context	target	context	target	context	target	context	target
EMNIST	CNP	0.9522 $\pm$ 0.023	0.7515 $\pm$ 0.0015	0.8977 $\pm$ 0.0016	0.6336 $\pm$ 0.017	0.8242 $\pm$ 0.0018	0.5784 $\pm$ 0.009	0.6566 $\pm$ 0.0017	0.5341 $\pm$ 0.016
	BCNP	0.9678 $\pm$ 0.010	0.8058 $\pm$ 0.008	0.9015 $\pm$ 0.008	0.6711 $\pm$ 0.009	0.8415 $\pm$ 0.007	0.6089 $\pm$ 0.009	0.6788 $\pm$ 0.006	0.5715 $\pm$ 0.006
	SCNP	<b>0.9716</b> $\pm$ 0.008	<b>0.8343</b> $\pm$ 0.006	<b>0.9251</b> $\pm$ 0.008	<b>0.6971</b> $\pm$ 0.007	<b>0.8674</b> $\pm$ 0.006	<b>0.6343</b> $\pm$ 0.007	<b>0.6986</b> $\pm$ 0.005	<b>0.5877</b> $\pm$ 0.005
	NP	0.9678 $\pm$ 0.004	0.7756 $\pm$ 0.005	0.9011 $\pm$ 0.009	0.6941 $\pm$ 0.006	0.8544 $\pm$ 0.009	0.6455 $\pm$ 0.007	0.7034 $\pm$ 0.009	0.5865 $\pm$ 0.006
	BNP	0.9757 $\pm$ 0.005	0.8358 $\pm$ 0.005	0.8116 $\pm$ 0.007	0.7625 $\pm$ 0.006	0.8759 $\pm$ 0.007	0.6773 $\pm$ 0.007	0.7451 $\pm$ 0.005	0.6237 $\pm$ 0.005
	SNP	<b>0.9847</b> $\pm$ 0.005	<b>0.8562</b> $\pm$ 0.006	<b>0.8368</b> $\pm$ 0.006	<b>0.7844</b> $\pm$ 0.005	<b>0.8984</b> $\pm$ 0.005	<b>0.6984</b> $\pm$ 0.005	<b>0.7653</b> $\pm$ 0.005	<b>0.6456</b> $\pm$ 0.004
	ANP	1.1125 $\pm$ 0.002	1.0321 $\pm$ 0.004	0.9815 $\pm$ 0.002	0.6366 $\pm$ 0.006	0.9021 $\pm$ 0.004	0.7053 $\pm$ 0.008	0.8454 $\pm$ 0.002	0.7034 $\pm$ 0.005
	BANP	1.1355 $\pm$ 0.003	1.0615 $\pm$ 0.005	1.0236 $\pm$ 0.002	0.6549 $\pm$ 0.005	0.9155 $\pm$ 0.004	0.7521 $\pm$ 0.006	0.8612 $\pm$ 0.003	0.7515 $\pm$ 0.005
	SANP	<b>1.1531</b> $\pm$ 0.000	<b>1.0877</b> $\pm$ 0.004	<b>1.0421</b> $\pm$ 0.002	<b>0.6776</b> $\pm$ 0.005	<b>0.9321</b> $\pm$ 0.002	<b>0.7843</b> $\pm$ 0.006	<b>0.8732</b> $\pm$ 0.003	<b>0.7657</b> $\pm$ 0.005
CELEBA	CNP	1.0323 $\pm$ 0.016	0.7845 $\pm$ 0.013	1.0177 $\pm$ 0.016	0.7438 $\pm$ 0.017	0.8956 $\pm$ 0.009	0.7344 $\pm$ 0.011	0.7677 $\pm$ 0.012	0.6096 $\pm$ 0.009
	BCNP	1.0452 $\pm$ 0.009	0.8015 $\pm$ 0.008	1.0275 $\pm$ 0.009	0.7726 $\pm$ 0.008	0.9351 $\pm$ 0.006	0.8376 $\pm$ 0.009	0.8015 $\pm$ 0.010	0.6816 $\pm$ 0.008
	SCNP	<b>1.0525</b> $\pm$ 0.008	<b>0.8243</b> $\pm$ 0.006	<b>1.0348</b> $\pm$ 0.008	<b>0.7869</b> $\pm$ 0.006	<b>0.9562</b> $\pm$ 0.004	<b>0.8545</b> $\pm$ 0.008	<b>0.8344</b> $\pm$ 0.006	<b>0.7045</b> $\pm$ 0.005
	NP	1.1333 $\pm$ 0.004	0.8766 $\pm$ 0.005	1.1043 $\pm$ 0.015	0.8355 $\pm$ 0.015	1.0034 $\pm$ 0.008	0.8456 $\pm$ 0.006	0.8935 $\pm$ 0.009	0.6893 $\pm$ 0.006
	BNP	1.1732 $\pm$ 0.005	0.8901 $\pm$ 0.006	1.1378 $\pm$ 0.007	0.8678 $\pm$ 0.006	1.0411 $\pm$ 0.008	0.8711 $\pm$ 0.005	0.9256 $\pm$ 0.008	0.7671 $\pm$ 0.006
	SNP	<b>1.1952</b> $\pm$ 0.005	<b>0.9062</b> $\pm$ 0.006	<b>1.1565</b> $\pm$ 0.005	<b>0.8846</b> $\pm$ 0.005	<b>1.0542</b> $\pm$ 0.006	<b>0.8956</b> $\pm$ 0.004	<b>0.9425</b> $\pm$ 0.006	<b>0.7985</b> $\pm$ 0.005
	ANP	1.2633 $\pm$ 0.002	1.0163 $\pm$ 0.004	1.1377 $\pm$ 0.004	0.9866 $\pm$ 0.006	1.0418 $\pm$ 0.004	0.8845 $\pm$ 0.006	0.9363 $\pm$ 0.004	0.7346 $\pm$ 0.008
	BANP	1.2751 $\pm$ 0.002	1.0389 $\pm$ 0.005	1.1488 $\pm$ 0.004	1.0155 $\pm$ 0.005	1.0602 $\pm$ 0.005	0.9255 $\pm$ 0.006	0.9489 $\pm$ 0.004	0.8415 $\pm$ 0.007
	SANP	<b>1.2854</b> $\pm$ 0.000	<b>1.0594</b> $\pm$ 0.004	<b>1.1685</b> $\pm$ 0.002	<b>1.0353</b> $\pm$ 0.004	<b>1.0772</b> $\pm$ 0.002	<b>0.9455</b> $\pm$ 0.004	<b>0.9655</b> $\pm$ 0.003	<b>0.8673</b> $\pm$ 0.005

100 state transition pairs are randomly selected for each configuration of the environment, working as the maximum context points to identify the configuration of dynamics.

### C.3 Image Completion

Analogous to the 1D experiments, we take random pixels of a given image at training as targets, and select a subset of this as contexts, again choosing the number of contexts and targets randomly ( $n \sim U[3, 200]$ ,  $m \sim n + U[0, 200 - n]$ ). The  $x_i$  are rescaled to  $[-1, 1]$  and the  $y_i$  are rescaled to  $[-0.5, 0.5]$ . We use a batch size of 16 for both EMNIST and CelebA, i.e. use 16 randomly selected images for each batch.

For image completion experiments on EMNIST and CelebA dataset, the neural architectures for CNP, NP, ANP, BCNP, BNP, BANP, and our SCNP, SNP, and SANP refer to Appendix B. The number of hidden unites is  $d_h = 128$  and latent representation  $d_z = 128$ . The number of layers are  $l_e = l_{de} = 4$ ,  $l_{ta} = l_{qk} = l_v = 5$ .  $h_{head} = 8$

### C.4 Uncertainty Measuring

Methods for reasoning under uncertainty are a key building block of accurate and reliable machine learning systems. We further analyze the learned models using the framework introduced in [17] to quantify uncertainty by investigate the calibration error and sharpness of the models. By assuming predictive distribution  $P_\theta(y_i|x_i, \mathcal{D}^c)$  as Gaussian distribution  $\mathcal{N}(y_i; \mu_i, \sigma_i^2)$ , we can get the probabilistic forecast  $F_i(x_i)$ . More formally,  $m$  confidence levels  $0 \leq p_1 < p_2 < \dots < p_m \leq 1$  are choose. For each threshold  $p_j$ , we compute the empirical frequency

$$\hat{p}_j = \frac{|y_i | F_i(y_i) \leq p_j, i \in \mathcal{T}|}{|\mathcal{T}|} \quad (31)$$

Table 10: Calibration error and sharpness of the models for 1D regression experiments (Mean  $\pm$  Std).

Method	RBF		Matern	
	CAL	SHA	CAL	SHA
CNP	0.0724 $\pm$ 0.008	0.0887 $\pm$ 0.007	0.0514 $\pm$ 0.005	0.0831 $\pm$ 0.010
BCNP	0.1015 $\pm$ 0.007	0.1151 $\pm$ 0.005	0.0724 $\pm$ 0.005	0.1152 $\pm$ 0.009
SCNP	0.1225 $\pm$ 0.005	0.1357 $\pm$ 0.006	0.0425 $\pm$ 0.005	0.1325 $\pm$ 0.008
NP	0.0615 $\pm$ 0.004	0.0715 $\pm$ 0.005	0.0343 $\pm$ 0.015	0.0717 $\pm$ 0.015
BNP	0.0871 $\pm$ 0.005	0.1052 $\pm$ 0.006	0.0325 $\pm$ 0.009	0.0715 $\pm$ 0.008
SNP	0.1155 $\pm$ 0.004	0.1257 $\pm$ 0.005	0.0347 $\pm$ 0.008	0.0718 $\pm$ 0.007
ANP	0.1532 $\pm$ 0.002	0.0616 $\pm$ 0.004	0.0921 $\pm$ 0.004	0.0871 $\pm$ 0.006
BANP	0.2353 $\pm$ 0.002	0.0689 $\pm$ 0.004	0.0752 $\pm$ 0.005	0.0741 $\pm$ 0.006
SANP	0.2633 $\pm$ 0.000	0.0741 $\pm$ 0.004	0.0415 $\pm$ 0.002	0.0667 $\pm$ 0.004

Table 11: Calibration error and sharpness of the models for system identification experiments (Mean  $\pm$  Std).

Method	CAL	SHA
CNP	0.0872 $\pm$ 0.008	0.0415 $\pm$ 0.007
BCNP	0.1051 $\pm$ 0.005	0.0765 $\pm$ 0.006
SCNP	0.1124 $\pm$ 0.005	0.0833 $\pm$ 0.005
NP	0.0821 $\pm$ 0.003	0.0581 $\pm$ 0.005
BNP	0.0952 $\pm$ 0.003	0.0616 $\pm$ 0.005
SNP	0.1001 $\pm$ 0.001	0.0668 $\pm$ 0.005
ANP	0.0863 $\pm$ 0.001	0.0673 $\pm$ 0.004
BANP	0.1235 $\pm$ 0.001	0.0815 $\pm$ 0.005
SANP	0.1431 $\pm$ 0.000	0.1094 $\pm$ 0.004

In this case, the calibration error is defined as a numerical score describing the quality of forecast calibration:

$$\text{CAL}(F_1, y_1, \dots, F_{|\mathcal{T}|}, y_{|\mathcal{T}|}) = \sum_{j=1}^m w_j \cdot (p_j - \hat{p}_j)^2 \quad (32)$$

here  $w_j$  is weight and we set  $w_j = 1$ .

Sharpness is measured by using the variance  $\text{var}(F_i) = \sigma_i^2$  of the random variable whose CDF is  $F_i$ . Low-variance predictions are tightly centered around one value. A sharpness score can be defined by

$$\text{SHA}(F_1, \dots, F_{|\mathcal{T}|}) = \frac{1}{|\mathcal{T}|} \sum_{i \in \mathcal{T}} \sigma_i^2 \quad (33)$$

We evaluated the CE and sharpness of CNP, NP, ANP, BCNP, NBP, BANP, and corresponding stable versions SCNP, SNP, and SANP trained in the experiments. Table 10, 11, 12 list the calibration error and sharpness score on 1D regression, system identification, and image completion tasks. In several settings, models with our stable solution can achieve better calibration and sharpness but work worse in some settings, such as the calibration error being worse than NP and ANP in terms of 1D regression tasks with Matern. The possible reason is that our method tends to produce conservative credible intervals, so become under-confident or less over-confident in some settings.

Table 12: Calibration error and sharpness of the models for image completion experiments (Mean  $\pm$  Std).

Method	EMNIST		CELEBA	
	CAL	SHA	CAL	SHA
CNP	0.0182 $\pm$ 0.002	0.0574 $\pm$ 0.002	0.0253 $\pm$ 0.002	0.0743 $\pm$ 0.001
BCNP	0.0415 $\pm$ 0.003	0.0716 $\pm$ 0.002	0.0412 $\pm$ 0.002	0.0981 $\pm$ 0.002
SCNP	0.0543 $\pm$ 0.001	0.0846 $\pm$ 0.002	0.0457 $\pm$ 0.002	0.1136 $\pm$ 0.002
NP	0.0163 $\pm$ 0.001	0.0671 $\pm$ 0.002	0.0261 $\pm$ 0.001	0.0711 $\pm$ 0.001
BNP	0.0352 $\pm$ 0.002	0.0815 $\pm$ 0.002	0.0463 $\pm$ 0.001	0.0981 $\pm$ 0.002
SNP	0.0446 $\pm$ 0.002	0.0918 $\pm$ 0.000	0.0532 $\pm$ 0.001	0.1046 $\pm$ 0.001
ANP	0.0156 $\pm$ 0.002	0.0656 $\pm$ 0.002	0.0261 $\pm$ 0.004	0.0815 $\pm$ 0.006
BANP	0.0412 $\pm$ 0.002	0.0871 $\pm$ 0.002	0.0513 $\pm$ 0.003	0.1125 $\pm$ 0.003
SANP	0.0558 $\pm$ 0.000	0.0954 $\pm$ 0.001	0.0632 $\pm$ 0.002	0.1265 $\pm$ 0.004

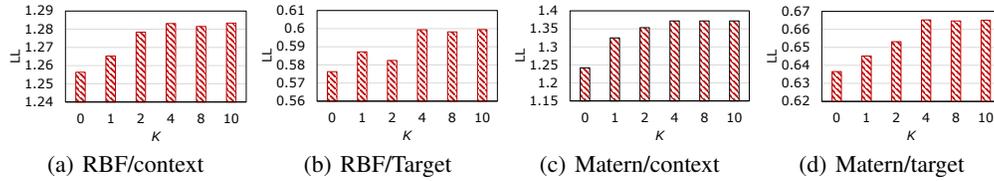


Figure 5: The Log-likelihood comparisons with different  $K$  on 1D regression task.

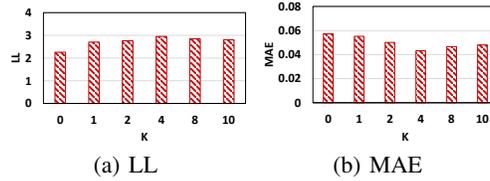


Figure 6: The LL and MAE comparisons with different  $K$  on system identification task.

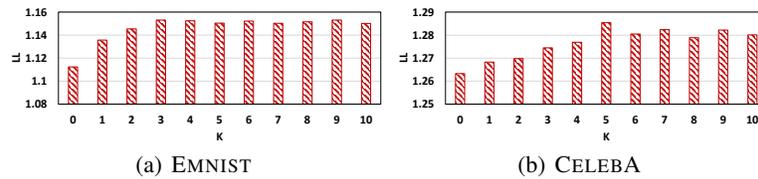


Figure 7: The LL comparisons with different  $K$  on image completion task.

### C.5 Ablation Study

The key parameter in our stable solution is the number of hard predictive subsets  $K$ . Taking SANP as an example, we investigated the average log-likelihood in terms of different  $K$  on 1D regression task, as shown in Figure 5. We can see that SANP performs better as  $K$  increases, reaches the best value at around  $K = 4$ , and then becomes stable in performance as  $K$  grows larger. As proved in Theorem 5.4 in the main manuscript, optimization on  $\mathcal{D}^c$  and more than one hard predicted subsets of  $\mathcal{D}^c$  can achieve more stable prediction. Similarly, we also conducted experiments to investigate the effect of  $K$  on System identification and image completion task and similar observations can be seen in Figure 6 and 7.