

1 Neural Processes

Neural Processes are an expansion of Gaussian Processes using Neural Networks introduced in [Gar+18b] they are trained on a context dataset which is. Neural processes are a meta-learning algorithm that can be used for few-shot learning. They are trained on a context dataset which is a set of input-output pairs $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$.

1.1 Meta Learning

Meta-Learning is a method of ML where the model is trained on a set of tasks and then tested on a new task. The goal is to learn a model that can learn new tasks quickly, i.e. ‘learning to learn’. These come under the terminology of few-shot learning where the model is trained on a small number of examples (more specifically, we can say N -way- K -shot learning where N is the number of classes and K is the number of examples per class).

Consider a set of datasets $\mathcal{D} = \{\mathcal{D}_i\}_{i=1}^D$ where each dataset \mathcal{D}_i is a set of input-output pairs $\mathcal{D}_i = \{(x_{ij}, y_{ij})\}_{j=1}^{N_i}$. The goal of meta-learning is to learn a model f that can learn a new task \mathcal{D}_{N+1} with a small number of examples, i.e. few-shot learning.

1.2 Architecture

A meta-dataset is split into two sets, the context set and the target set (also called the query set) where the sets are disjoint $\mathcal{D} = \mathcal{C} \cup \mathcal{T}$ and $\mathcal{C} \cap \mathcal{T} = \emptyset$. The model is trained on the context set and then tested on the target set to see how well it can generalize to new tasks. In essence, our task is to predict the outputs for the target conditioned on the training of the context set.

To achieve this we used a neural network (usually MLP) to encode the context set into embeddings.

$$\mathbf{r}(\mathcal{C}_i) = \text{Enc}_\theta(\mathcal{C}_i) = \text{Enc}_\theta(\{\mathbf{x}_i^{(c)}, \mathbf{y}_i^{(c)}\}) \quad (1.1)$$

Where \mathbf{r} is the embedding of the context set \mathcal{C}_i and θ are the parameters of the encoder. Then we aggregate the embeddings of the context sets to get a global representation of the dataset.

$$\mathbf{R}(\mathcal{C}) = \text{Agg}(\{\mathbf{r}(\mathcal{C}_i)\}_{i=1}^D) \quad (1.2)$$

Typically the aggregation function is a simple summation of the embeddings. The global representation \mathbf{R} is then used to condition the decoder to predict the outputs of the target set $\mathcal{T} = \{\mathbf{x}_i^{(t)}\}$ to give us a posterior distribution over the outputs $\mathbf{y}_i^{(t)}$.

$$p(\mathbf{y}_i^{(t)} | \mathbf{x}_i^{(t)}, \mathcal{C}) = \text{Dec}_\theta(\mathbf{x}_i^{(t)}, \mathbf{R}(\mathcal{C})) \quad (1.3)$$

The overall architecture is shown in Figure 1.1.

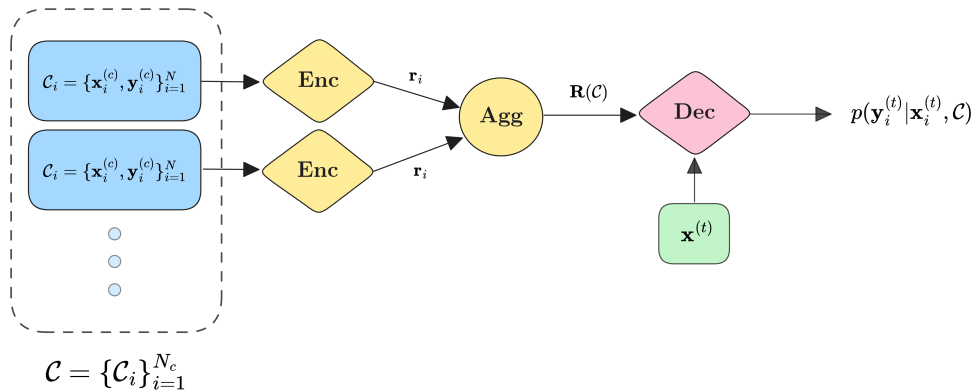


Figure 1.1: CNP Architecture

1.3 Neural Processes vs Gaussian Processes

From this, it is clear that NPs achieve a *prediction* map, π which maps the context set to the target set into a posterior distribution over the outputs.

$$\pi : \mathcal{C} \times \mathcal{X}^{(t)} \rightarrow \mathcal{P}(\mathcal{Y}) \quad (1.4)$$

This illustrates the similarity to Gaussian Processes (GP) [RW06] which perform the same underlying task of predicting the outputs of a target set conditioned on a context set. The difference is that GPs use a kernel function to compute the covariance between the context and target sets whereas NPs use a neural network to compute the covariance.

$$\text{GP} : \pi \rightarrow \mathcal{N}(m(\mathcal{C}, \mathbf{x}^{(t)}), k(\mathbf{x}^{(c)}, \mathbf{x}^{(t)})) \quad (1.5)$$

$$\text{NP} : \pi \rightarrow \text{Dec}_{\theta}(\mathbf{x}^{(t)}, \text{Agg}(\{\text{Enc}_{\theta}(\mathcal{C}_i)\}_{i=1}^D)) \quad (1.6)$$

1.4 Conditional Neural Processes

Conditional Neural Processes (CNPs) [Gar+18a] was one of the two original Neural Processes introduced by Garnelo et al. in 2018. Using the architecture framework described above, CNPs can be described as follows.

$$\text{Enc}_{\theta_e} = \text{MLP}_{\theta_e}$$

$$\text{Agg} = \text{Sum}$$

$$\text{Dec}_{\theta_d} = \text{MLP}_{\theta_d}$$

The encoding and summation is an implementation of the ‘DeepSet’ architecture [Zah+18] which is a neural network that is **Permutation Invariant** (PI) due to the summation operation. This means that the order of the input does not matter.

CNPs make the strong assumption that the posterior distribution *factorizes* over the target points. This means that the posterior distribution over the target points is independent of each other.

$$p(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \mathcal{C}) \stackrel{(a)}{=} \prod_{i=1}^{N_t} p(y_i^{(t)} | x_i^{(t)}, \mathbf{R}(\mathcal{C})) \quad (1.7)$$

$$\stackrel{(b)}{=} \prod_{i=1}^{N_t} \mathcal{N}(y_i^{(t)} | \mu_i, \sigma_i^2) \quad (1.8)$$

$$\stackrel{(c)}{=} \mathcal{N}(\mathbf{y}^{(t)} | \boldsymbol{\mu}(\mathbf{x}^{(t)}, \mathcal{C}), \boldsymbol{\Sigma}(\mathbf{x}^{(t)}, \mathcal{C})) \quad (1.9)$$

The benefit of this factorization assumption (a) is that the model can scale linearly with the number of target points with a tractable likelihood. However, this assumption means **CNPs are unable to generate coherent sample paths, they are only able to produce distributions over the target points.** Furthermore, we need to select a marginal likelihood for the distribution (b) which is usually a Heteroscedastic Gaussian Likelihood (Gaussian with a variance that varies with the input) [Gar+18a]. This also adds a further assumption as we have to select a likelihood for the distribution which may not be appropriate for the data we are modeling.

Since the product of Gaussians is a Gaussian (c) the model learns a mean and variance for each target point (though typically we learn the log variance to ensure it is positive), in this way, the model decoder can be interpreted as outputting a function of mean and variance for each target point.

As the likelihood is a Gaussian, the model can be trained using simple maximum likelihood estimation (MLE) by minimizing the negative log-likelihood (NLL) of the target points.

1.5 Latent Neural Processes

Neural Processes better named ‘Latent Neural Processes’ are an extension of CNPs that can generate coherent sample paths and are not restricted to a specific likelihood. They can do this by learning a latent representation of the context set \mathcal{C} which is then used to condition the decoder. We will call this latent representation \mathbf{z} (instead of \mathbf{R}) to avoid confusion with the non-latent global representation \mathbf{R} used in CNPs).

The encoder learns a *distribution* over the latent representation \mathbf{z} of the context set \mathcal{C} . Then the decoder learns a distribution over the target outputs $\mathbf{y}^{(t)}$ conditioned on the latent representation \mathbf{z} and the target inputs $\mathbf{x}^{(t)}$.

$$p(\mathbf{z}|\mathcal{C}) = \text{Enc}_\theta(\mathcal{C})$$

$$p(\mathbf{y}^{(t)}|\mathbf{x}^{(t)}, \mathbf{z}) = \text{Dec}_\theta(\mathbf{x}^{(t)}, \mathbf{z})$$

The full model is shown in Figure 1.2.

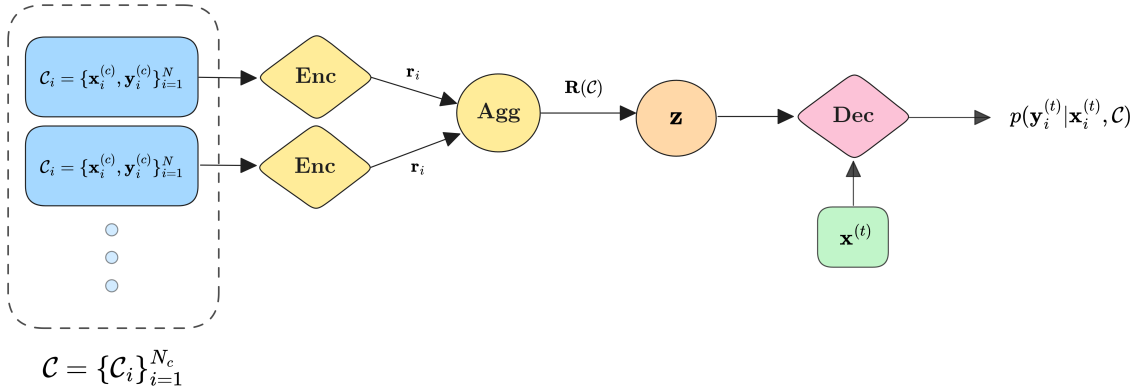


Figure 1.2: LNP Architecture

The likelihood of the target points is then given by the marginal likelihood of the latent representation \mathbf{z} .

$$p(\mathbf{y}^{(t)}|\mathbf{x}^{(t)}, \mathcal{C}) = \int p(\mathbf{z}|\mathcal{C})p(\mathbf{y}^{(t)}|\mathbf{x}^{(t)}, \mathbf{z})d\mathbf{z} \quad (1.10)$$

$$\stackrel{(a)}{=} \int p(\mathbf{z}|\mathcal{C}) \left(\prod_{i=1}^{N_t} p(y_i^{(t)}|x_i^{(t)}, \mathbf{z}) \right) d\mathbf{z} \quad (1.11)$$

$$\stackrel{(b)}{=} \int p(\mathbf{z}|\mathcal{C}) \left(\prod_{i=1}^{N_t} \mathcal{N}(y_i^{(t)}|\mu_i, \sigma_i^2) \right) d\mathbf{z} \quad (1.12)$$

We can see that we express the conditional distribution over the target points conditioned on the latent representation \mathbf{z} as a factorized (a) product of Gaussians (b), this may seem like we are making the same factorization assumption as CNPs. However, the difference is that we are not making this assumption conditioned on the latent variable instead of the context set. This integral models an *infinite mixture* of Gaussians which allows us to model *any* distribution over the target points. The downside is that the likelihood is intractable and we need to use approximate inference methods such as Variational Inference (VI) or Sample Estimation using Monte Carlo (MC) methods to train the model which typically leads to biased estimates of the likelihood with high variance, thus we require more samples to train the model.

References

- [Gar+18a] Marta Garnelo et al. *Conditional Neural Processes*. 2018. arXiv: [1807.01613 \[cs.LG\]](#).
- [Gar+18b] Marta Garnelo et al. *Neural Processes*. 2018. arXiv: [1807.01622 \[cs.LG\]](#).
- [RW06] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. Adaptive computation and machine learning. MIT Press, 2006, pp. I–XVIII, 1–248. ISBN: 026218253X.
- [Zah+18] Manzil Zaheer et al. *Deep Sets*. 2018. arXiv: [1703.06114 \[cs.LG\]](#).