

Bayesian Context Aggregation for Neural Processes

ICLR 2021

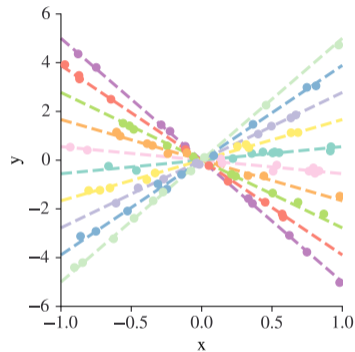
Michael Volpp^{1,2}, Fabian Flürenbrock¹, Lukas Grossberger¹, Christian Daniel¹,
Gerhard Neumann²

¹Bosch Center for Artificial Intelligence, Renningen, Germany

²Karlsruhe Institute of Technology, Karlsruhe, Germany

Probabilistic Regression as a Multi-task Learning Problem

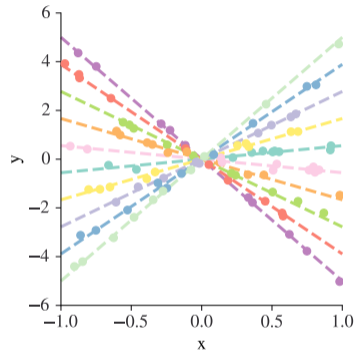
Family F of functions $f: \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ (“tasks”) with some form of shared structure



Probabilistic Regression as a Multi-task Learning Problem

Family F of functions $f: \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ ("tasks") with some form of shared structure

Noisy evaluations $D = \{(x_{\cdot;j}, y_{\cdot;j})\}_j$ with $y_{\cdot;j} = f_{\cdot}(x_{\cdot;j}) + \epsilon_{\cdot;j}$, $\epsilon_{\cdot;j} \sim N(0, \frac{2}{n})$

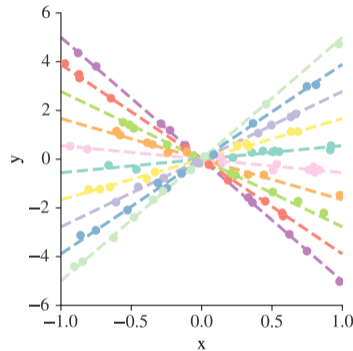


Probabilistic Regression as a Multi-task Learning Problem

Family F of functions $f: \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ ("tasks") with some form of shared structure

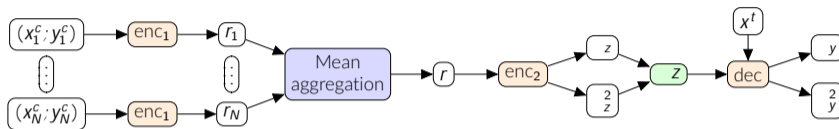
Noisy evaluations $D_i = \{(x_{\cdot;j}; y_{\cdot;j})\}_j$ with $y_{\cdot;j} = f_{\cdot}(x_{\cdot;j}) + \epsilon_j; \epsilon_j \sim \mathcal{N}(0; \frac{2}{n})$

Learn the predictive distribution $p(y^t | x^t; D^c)$ over target outputs, conditioned on a context set $D^c \subseteq D_i$.



The Neural Process¹

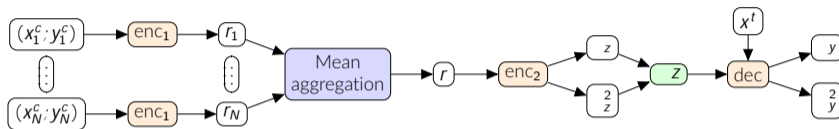
Neural network (NN)-based multi-task learning architecture



¹Garnelo et al., *Conditional Neural Processes*, ICML 2018, Garnelo et al. *Neural Processes*. arxiv/1807.01622 2018

The Neural Process¹

Neural network (NN)-based multi-task learning architecture

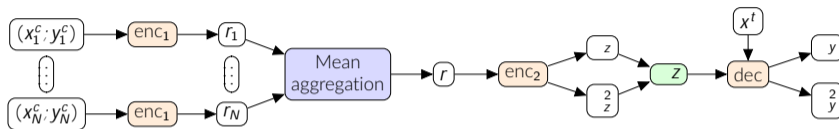


1. Infer latent representation $p(z | D^c)$ of the target function from $D^c = \{(x_n^c; y_n^c)\}_{n=1}^N$

¹Garnelo et al., *Conditional Neural Processes*, ICML 2018, Garnelo et al. *Neural Processes*. arxiv/1807.01622 2018

The Neural Process¹

Neural network (NN)-based multi-task learning architecture

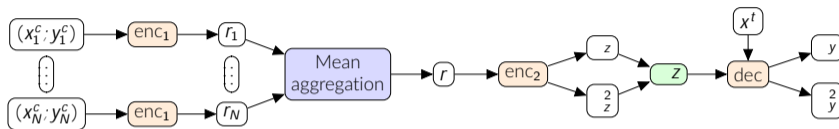


1. Infer latent representation $p(z | D^c)$ of the target function from $D^c = \{(x_n^c; y_n^c)\}_{n=1}^N$
 - a. Map each context tuple $(x_n^c; y_n^c)$ onto a latent observation $r_n = enc_1(x_n^c; y_n^c)$

¹Garnelo et al., *Conditional Neural Processes*, ICML 2018, Garnelo et al. *Neural Processes*. arxiv/1807.01622 2018

The Neural Process¹

Neural network (NN)-based multi-task learning architecture

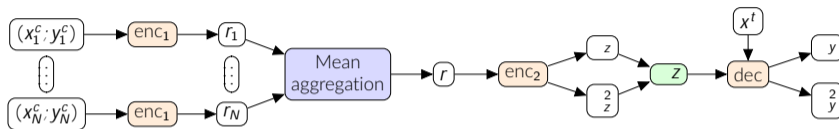


1. Infer latent representation $p(z | D^c)$ of the target function from $D^c = \{(x_n^c; y_n^c)\}_{n=1}^N$
 - a. Map each context tuple $(x_n^c; y_n^c)$ onto a latent observation $r_n = enc_1(x_n^c; y_n^c)$
 - b. Form an aggregated latent observation r using *mean aggregation*: $r = \frac{1}{N} \sum_{n=1}^N r_n$

¹Garnelo et al., *Conditional Neural Processes*, ICML 2018, Garnelo et al. *Neural Processes*. arxiv/1807.01622 2018

The Neural Process¹

Neural network (NN)-based multi-task learning architecture

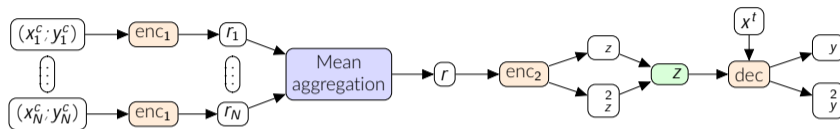


1. Infer latent representation $p(z | D^c)$ of the target function from $D^c = \{(x_n^c; y_n^c)\}_{n=1}^N$
 - a. Map each context tuple $(x_n^c; y_n^c)$ onto a latent observation $r_n = \text{enc}_1(x_n^c; y_n^c)$
 - b. Form an aggregated latent observation r using *mean aggregation*: $r = \frac{1}{N} \sum_{n=1}^N r_n$
 - c. Map r onto the parameters of the latent distribution: $(z; \bar{z}) = \text{enc}_2(r)$

¹Garnelo et al., *Conditional Neural Processes*, ICML 2018, Garnelo et al. *Neural Processes*. arxiv/1807.01622 2018

The Neural Process¹

Neural network (NN)-based multi-task learning architecture



1. Infer latent representation $p(z | D^c)$ of the target function from $D^c = \{(x_n^c; y_n^c)\}_{n=1}^N$
 - a. Map each context tuple $(x_n^c; y_n^c)$ onto a latent observation $r_n = enc_1(x_n^c; y_n^c)$
 - b. Form an aggregated latent observation r using *mean aggregation*: $r = \frac{1}{N} \sum_{n=1}^N r_n$
 - c. Map r onto the parameters of the latent distribution: $(z; \zeta) = enc_2(r)$
2. Map samples $z \sim p(z | D^c)$ onto a Gaussian output distribution: $(y; \hat{y}) = dec(z; x^t)$

¹Garnelo et al., *Conditional Neural Processes*, ICML 2018, Garnelo et al. *Neural Processes*. arxiv/1807.01622 2018

Context Aggregation and Task Ambiguity

Different areas in the $(x; y)$ -space can have different *task ambiguity* (TA)

Context Aggregation and Task Ambiguity

Different areas in the $(x; y)$ -space can have different *task ambiguity* (TA)

Context tuples $(x_n^c; y_n^c)$ with

high TA should have little influence on $p(z | D^c)$

low TA should have large influence on $p(z | D^c)$

Context Aggregation and Task Ambiguity

Different areas in the $(x; y)$ -space can have different *task ambiguity* (TA)

Context tuples $(x_n^c; y_n^c)$ with

high TA should have little influence on $p(z | D^c)$

low TA should have large influence on $p(z | D^c)$

Mean aggregation assigns the same weight

to all context tuples: $r = \frac{1}{N} \sum_{n=1}^N r_n$

Context Aggregation and Task Ambiguity

Different areas in the $(x; y)$ -space can have different *task ambiguity* (TA)

Context tuples $(x_n^c; y_n^c)$ with

high TA should have little influence on $p(z | D^c)$

low TA should have large influence on $p(z | D^c)$

Mean aggregation assigns the same weight

to all context tuples: $r = \frac{1}{N} \sum_{n=1}^N r_n$

How to efficiently incorporate task ambiguity into NP parameter inference?

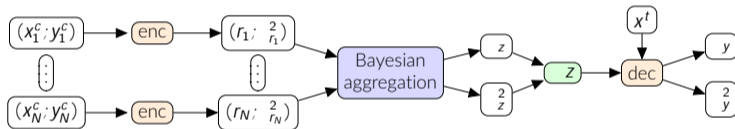
Bayesian Context Aggregation

*Context aggregation and parameter inference should be treated as one holistic mechanism!
Directly aggregate the context data into the statistical description of z !*

Bayesian Context Aggregation

*Context aggregation and parameter inference should be treated as one holistic mechanism!
Directly aggregate the context data into the statistical description of z !*

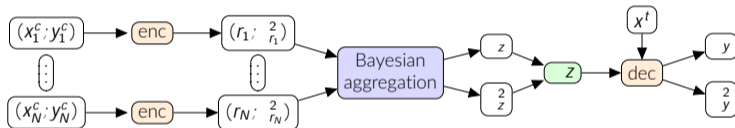
NP with our *Bayesian aggregation (BA)*:



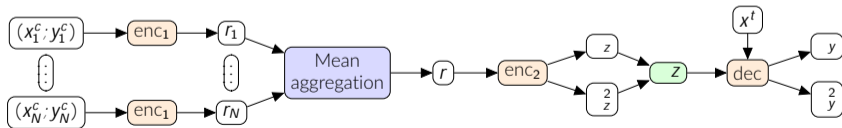
Bayesian Context Aggregation

*Context aggregation and parameter inference should be treated as one holistic mechanism!
Directly aggregate the context data into the statistical description of z !*

NP with our *Bayesian aggregation (BA)*:

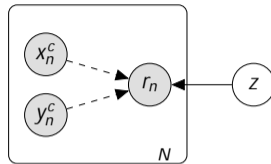


Compare: NP with traditional mean aggregation (MA):



Bayesian Context Aggregation

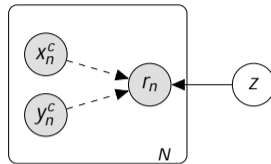
Context aggregation as Bayesian inference



Bayesian Context Aggregation

Context aggregation as Bayesian inference

Observation model: $p(r_n | z) = \mathcal{N}(r_n | z; \text{diag } \Sigma)$

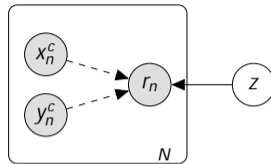


Bayesian Context Aggregation

Context aggregation as Bayesian inference

Observation model: $p(r_n | z) = \mathcal{N}(r_n | z; \text{diag } \frac{2}{r_n})$

Encoder learns: $(r_n; \frac{2}{r_n}) = \text{enc}(x_n^c; y_n^c)$



Bayesian Context Aggregation

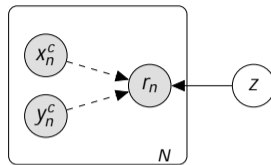
Context aggregation as Bayesian inference

Observation model: $p(r_n | z) = \mathcal{N}(r_n | z; \text{diag}(\frac{2}{r_n}))$

Encoder learns: $(r_n; \frac{2}{r_n}) = \text{enc}(x_n^c; y_n^c)$

Latent posterior: $p(z | \{r_n\}) = \mathcal{N}(z | z; \text{diag}(\frac{2}{z}))$

$$\frac{2}{z} = \prod_{n=1}^N \frac{2}{r_n}^{-1} ; \quad z = \prod_{n=1}^N \frac{z}{r_n}$$



(for each latent dim.)

Bayesian Context Aggregation

Context aggregation as Bayesian inference

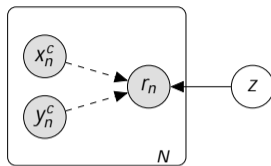
Observation model: $p(r_n | z) = \mathcal{N}(r_n | z; \text{diag}(\frac{2}{r_n}))$

Encoder learns: $(r_n; \frac{2}{r_n}) = \text{enc}(x_n^c; y_n^c)$

Latent posterior: $p(z | \{r_n\}) = \mathcal{N}(z | z; \text{diag}(\frac{2}{z}))$

$$\frac{2}{z} = \prod_{n=1}^N \frac{2}{r_n}^{-1}; \quad z = \prod_{n=1}^N \frac{z}{r_n}$$

r_n enters z with learned weight $\frac{2}{z} \frac{2}{r_n}$



(for each latent dim.)

Bayesian Context Aggregation

Context aggregation as Bayesian inference

Observation model: $p(r_n | z) = \mathcal{N}(r_n | z; \text{diag}(\frac{2}{r_n}))$

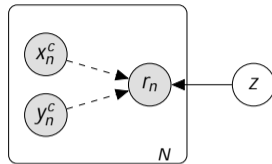
Encoder learns: $(r_n; \frac{2}{r_n}) = \text{enc}(x_n^c; y_n^c)$

Latent posterior: $p(z | \{r_n\}) = \mathcal{N}(z | z; \text{diag}(\frac{2}{z}))$

$$\frac{2}{z} = \prod_{n=1}^N \frac{2}{r_n}^{-1}; \quad z = \prod_{n=1}^N \frac{2}{r_n} r_n$$

r_n enters z with learned weight $\frac{2}{z} \frac{2}{r_n}$

Principled quantification of task ambiguity



(for each latent dim.)

Bayesian Context Aggregation

Context aggregation as Bayesian inference

Observation model: $p(r_n | z) = \mathcal{N}(r_n | z; \text{diag}(\frac{2}{r_n}))$

Encoder learns: $(r_n; \frac{2}{r_n}) = \text{enc}(x_n^c; y_n^c)$

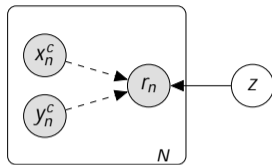
Latent posterior: $p(z | \{r_n\}) = \mathcal{N}(z | z; \text{diag}(\frac{2}{z}))$

$$\frac{2}{z} = \prod_{n=1}^N \frac{2}{r_n}^{-1}; \quad z = \prod_{n=1}^N \frac{2}{r_n}$$

r_n enters z with learned weight $\frac{2}{z} \frac{2}{r_n}$

Principled quantification of task ambiguity

Only marginal computational overhead



(for each latent dim.)

Bayesian Context Aggregation

Context aggregation as Bayesian inference

Observation model: $p(r_n | z) = \mathcal{N}(r_n | z; \text{diag}(\frac{2}{r_n}))$

Encoder learns: $(r_n; \frac{2}{r_n}) = \text{enc}(x_n^c; y_n^c)$

Latent posterior: $p(z | \{r_n\}) = \mathcal{N}(z | z; \text{diag}(\frac{2}{z}))$

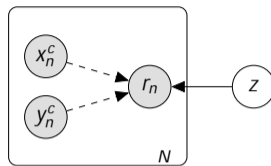
$$\frac{2}{z} = \prod_{n=1}^N \frac{2}{r_n}^{-1}; \quad z = \prod_{n=1}^N \frac{2}{r_n}$$

r_n enters z with learned weight $\frac{2}{z} \frac{2}{r_n}$

Principled quantification of task ambiguity

Only marginal computational overhead

Compatible with existing NP architectures



(for each latent dim.)

Experiments

	PB/det.		VI		MC	
	BA	MA (CNP)	BA	MA (LP-NP)	BA	MA
RBF GP	1:37 ± 0:15	0:94 ± 0:04	1:40 ± 0:04	0:45 ± 0:12	1:62 ± 0:05	1:07 ± 0:05
Weakly Periodic GP	1:13 ± 0:08	0:76 ± 0:02	0:89 ± 0:03	0:07 ± 0:14	1:30 ± 0:06	0:85 ± 0:04
Matern-5/2 GP	-0:50 ± 0:07	-0:68 ± 0:01	-0:79 ± 0:01	-1:09 ± 0:10	-0:33 ± 0:01	-0:90 ± 0:15
Furuta Dynamics	7:50 ± 0:27	7:06 ± 0:12	7:32 ± 0:18	5:57 ± 0:21	8:25 ± 0:33	7:55 ± 0:24

