

# Function-Space Variational Inference in the context of Implicit Processes

Machine Learning Group

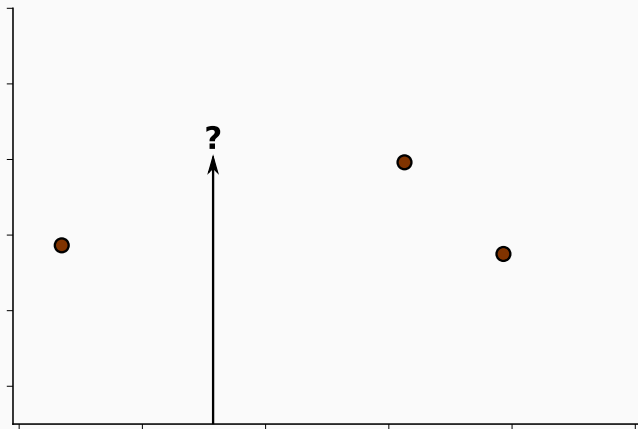
---

Luis Antonio Ortega Andrés

December 1, 2022

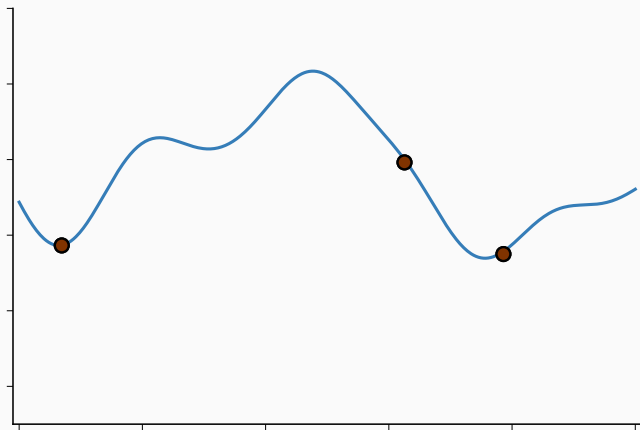
Autonomous University of Madrid

# Motivation



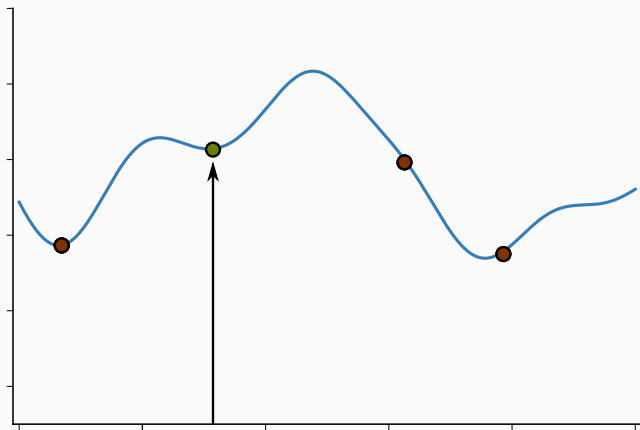
Objective: Make predictions over unknown points.

# Motivation



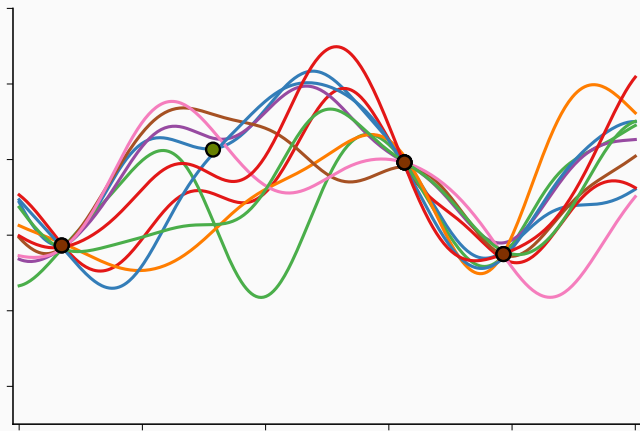
**Procedure:** Learn a function that explain the visible data.

# Motivation



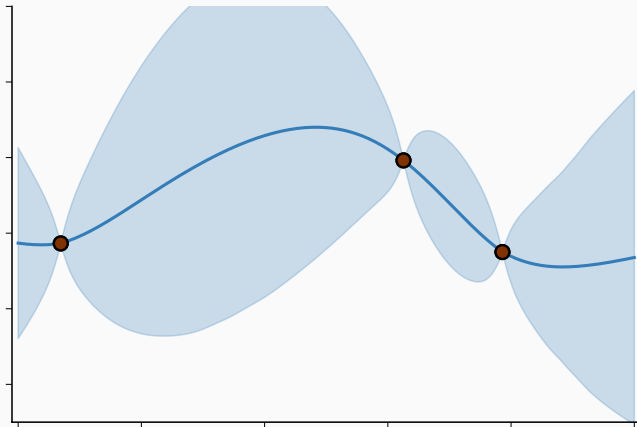
**Outcome:** The function can be used to predict.

# Motivation



Question: How confident is the prediction?

# Motivation



**Answer:** Bayesian approach.

1. Bayesian Supervised Learning
2. Variational Inference
3. Gaussian Processes
  - 3.1 Sparse Gaussian Processes
4. Implicit Processes
  - 4.1 Variational Implicit Processes
  - 4.2 Sparse Implicit Processes
  - 4.3 Linearized approximation
5. Sparse Linearized Implicit Processes
6. Experiments
7. Conclusions

# Bayesian Supervised Learning

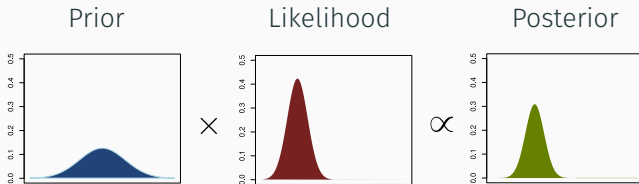
---



# Bayesian Supervised Learning

**Objective:** Learn an unknown function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^M$  given a set of observations  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \subset \mathbb{R}^D, \mathbf{y} = (y_1, \dots, y_N) \subset \mathbb{R}^M$ .

**Approach.** Consider a set of latent random variables  $\mathbf{z}$  that model the generation of the dataset  $P(\mathbf{y}|\mathbf{X}, \mathbf{z})$ .



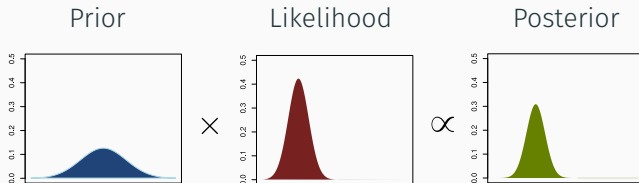
# Bayesian Supervised Learning

**Problem.** Prediction requires the posterior  $P(\mathbf{z}|\mathbf{X}, \mathbf{y})$ :

$$P(y_\star|\mathbf{x}_\star, \mathbf{X}, \mathbf{y}) = \int P(y_\star|\mathbf{x}_\star, \mathbf{z})P(\mathbf{z}|\mathbf{X}, \mathbf{y}) d\mathbf{z} ,$$

which is usually intractable due to the integral

$$P(\mathbf{z}|\mathbf{X}, \mathbf{y}) = \frac{P(\mathbf{y}, \mathbf{z}|\mathbf{X})}{\int P(\mathbf{y}, \mathbf{z}|\mathbf{X}) d\mathbf{z}} .$$



# Variational Inference

---

**Idea.** Approximate the posterior  $P(\mathbf{z}|\mathbf{X}, \mathbf{y})$  with a **simpler distribution**  $Q(\mathbf{z})$  and ensure that  $KL(Q(\mathbf{z}) \mid P(\mathbf{z}|\mathbf{X}, \mathbf{y}))$  is close to 0.

Formally,

$$\begin{aligned} Q^*(\mathbf{z}) &= \arg \min_{Q \in \mathcal{Q}} KL(Q(\mathbf{z}) \mid P(\mathbf{z}|\mathbf{X}, \mathbf{y})) \\ &= \arg \max_{Q \in \mathcal{Q}} \mathbb{E}_{Q(\mathbf{z})} [\log P(\mathbf{y}|\mathbf{X}, \mathbf{z})] - KL(Q(\mathbf{z}) \mid P(\mathbf{z})). \end{aligned}$$

**Idea.** Approximate the posterior  $P(\mathbf{z}|\mathbf{X}, \mathbf{y})$  with a **simpler distribution**  $Q(\mathbf{z})$  and ensure that  $KL(Q(\mathbf{z}) | P(\mathbf{z}|\mathbf{X}, \mathbf{y}))$  is close to 0.

Formally,

$$\begin{aligned} Q^*(\mathbf{z}) &= \arg \min_{Q \in \mathcal{Q}} KL(Q(\mathbf{z}) | P(\mathbf{z}|\mathbf{X}, \mathbf{y})) \\ &= \arg \max_{Q \in \mathcal{Q}} \underbrace{\mathbb{E}_{Q(\mathbf{z})} [\log P(\mathbf{y}|\mathbf{X}, \mathbf{z})] - KL(Q(\mathbf{z}) | P(\mathbf{z}))}_{\text{ELBO}} . \end{aligned}$$

**Idea.** Approximate the posterior  $P(\mathbf{z}|\mathbf{X}, \mathbf{y})$  with a **simpler distribution**  $Q(\mathbf{z})$  and ensure that  $KL(Q(\mathbf{z}) | P(\mathbf{z}|\mathbf{X}, \mathbf{y}))$  is close to 0.

Formally,

$$\begin{aligned} Q^*(\mathbf{z}) &= \arg \min_{Q \in \mathcal{Q}} KL(Q(\mathbf{z}) | P(\mathbf{z}|\mathbf{X}, \mathbf{y})) \\ &= \arg \max_{Q \in \mathcal{Q}} \underbrace{\mathbb{E}_{Q(\mathbf{z})} [\log P(\mathbf{y}|\mathbf{X}, \mathbf{z})]}_{\text{Data Fitting term}} - \underbrace{KL(Q(\mathbf{z}) | P(\mathbf{z}))}_{\text{Regularizer}} . \end{aligned}$$

# Gaussian Processes

---

# Gaussian Processes

A **Gaussian process** (GP) is a collection of random variables  $f(\cdot)$  such that any finite collection  $\mathbf{f} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$  follows a Gaussian distribution

$$\mathbf{f} = f(\mathbf{X}) = \left( f(\mathbf{x}_1), \dots, f(\mathbf{x}_N) \right) \sim \mathcal{N} \left( m(\mathbf{X}), \kappa(\mathbf{X}, \mathbf{X}) \right).$$

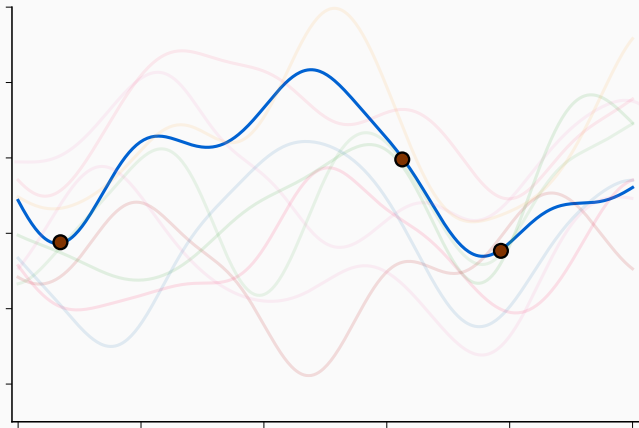
Gaussian processes place a **probability distribution over functions**.

Usually,

$$m(\mathbf{x}) = 0 \quad \text{and} \quad \kappa(\mathbf{x}_1, \mathbf{x}_2) = \sigma^2 \exp \left( -\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{2l^2} \right).$$



**Hypothesis:** The unknown function  $f$  is a (noisy) sample from a Gaussian process.



For a regression problem, using,

- the Gaussian process prior over functions,

$$P(\mathbf{f}) = \mathcal{N}\left(m(\mathbf{X}), \kappa(\mathbf{X}, \mathbf{X})\right),$$

- a suitable likelihood

$$P(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I}).$$

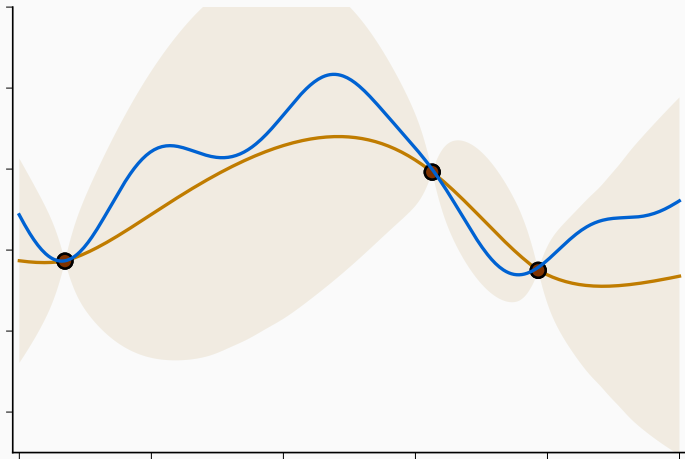
Predictive distribution is computable in closed form

$$P(f(\mathbf{x}^*)|\mathbf{y}, \mathbf{X}) = \mathcal{N}(\mu^*, \Sigma^*),$$

with

$$\begin{aligned}\mu^* &= m(\mathbf{x}^*) + \kappa(\mathbf{x}^*, \mathbf{X})(\kappa(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X})), \\ \Sigma^* &= \kappa(\mathbf{x}^*, \mathbf{x}^*) - \kappa(\mathbf{x}^*, \mathbf{X})(\kappa(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1}\kappa(\mathbf{X}, \mathbf{x}^*).\end{aligned}$$

$$\begin{aligned}\mu^* &= m(\mathbf{x}^*) + \kappa(\mathbf{x}^*, \mathbf{X})(\kappa(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X})), \\ \Sigma^* &= \kappa(\mathbf{x}^*, \mathbf{x}^*) - \kappa(\mathbf{x}^*, \mathbf{X})(\kappa(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1}\kappa(\mathbf{X}, \mathbf{x}^*).\end{aligned}$$



Hyper-parameters (kernel and likelihood) can be **optimized** using the marginal log likelihood:

$$\log P(\mathbf{y}) = -\frac{1}{2}\mathbf{y}^T(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{y} - \frac{1}{2}\log \det(\mathbf{K} + \sigma^2\mathbf{I}) - \frac{N}{2}\log 2\pi ,$$

with  $\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$ .

**Limitation.** Requires the computation of  $(\mathbf{K} + \sigma^2\mathbf{I})^{-1}$  which is  $\mathcal{O}(N^3)$ .

**Limitation.** Mini-batches cannot be used for optimization.

Hyper-parameters (kernel and likelihood) can be **optimized** using the marginal log likelihood:

$$\log P(\mathbf{y}) = -\frac{1}{2}\mathbf{y}^T(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{y} - \frac{1}{2}\log \det(\mathbf{K} + \sigma^2\mathbf{I}) - \frac{N}{2}\log 2\pi ,$$

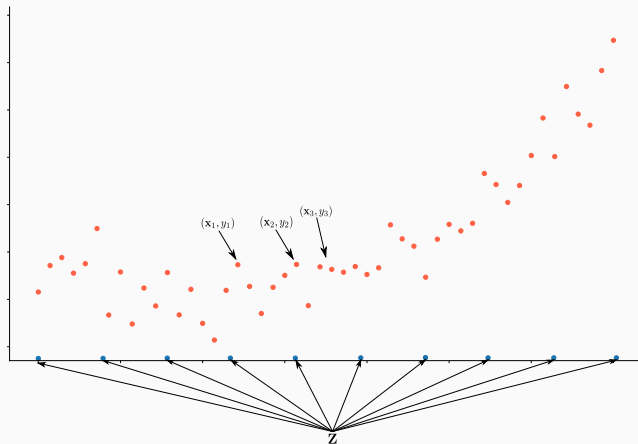
with  $\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$ .

**Limitation.** Requires the computation of  $(\mathbf{K} + \sigma^2\mathbf{I})^{-1}$  which is  $\mathcal{O}(N^3)$ .

**Limitation.** Mini-batches cannot be used for optimization.

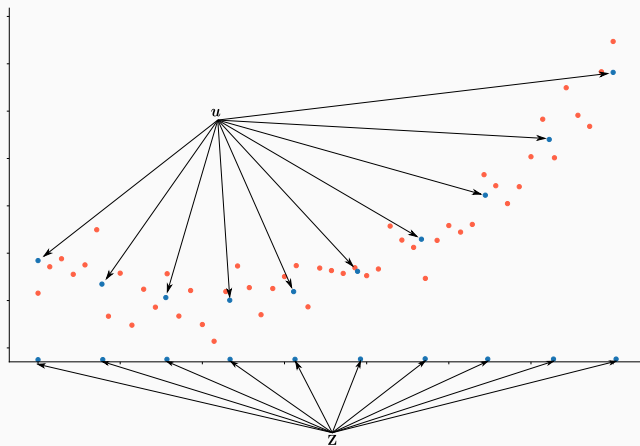
# Sparse Gaussian Processes

Initialize a set of **inducing locations**  $\mathbf{Z}$  with  $|\mathbf{Z}| < |\mathbf{X}|$ . For example  $\mathbf{Z}$  can be the centers of applying KMeans to  $\mathbf{X}$ .



# Sparse Gaussian Processes

Set a variational distribution over the **inducing points**  $\mathbf{u} = f(\mathbf{Z})$ , typically,  $Q(\mathbf{u}) = \mathcal{N}(\mathbf{m}, \mathbf{S})$ .

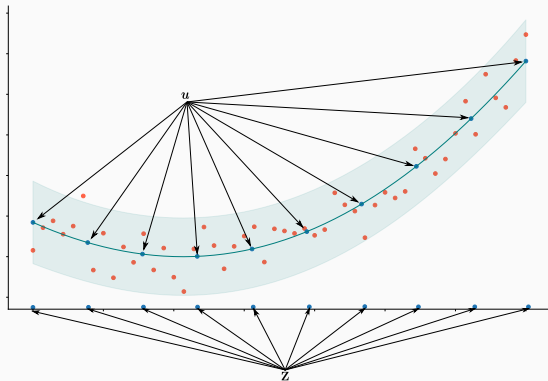


# Sparse Gaussian Processes

The approximated posterior predictive distribution is,

$$Q(f(\mathbf{x}_*)) = \int_{\mathbf{u}} P(f(\mathbf{x}_*)|\mathbf{u})Q(\mathbf{u}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

is Gaussian and suitable for making predictions.





Minimize the ELBO to optimize  $Q(\mathbf{u}, \mathbf{f})$ , with  $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))$ ,

$$\begin{aligned} Q^*(\mathbf{u}, \mathbf{f}) &= \arg \min_{Q \in \mathcal{Q}} KL\left(Q(\mathbf{u}, \mathbf{f}) \mid P(\mathbf{u}, \mathbf{f} \mid \mathbf{X}, \mathbf{y})\right) \\ &= \arg \max_{Q \in \mathcal{Q}} \mathbb{E}_{Q(\mathbf{f})} \left[ \log P(\mathbf{y} \mid \mathbf{X}, \mathbf{f}) \right] - KL\left(Q(\mathbf{u}) \mid P(\mathbf{u})\right). \end{aligned}$$

where

1.  $P(\mathbf{u})$  is given by the Gaussian Process distribution evaluated at  $\mathbf{Z}$ .
2. The Kullback-Leibler divergence is between Gaussian distributions.
3. The data-fitting term can be computed for Gaussian likelihoods or approximated by quadrature.

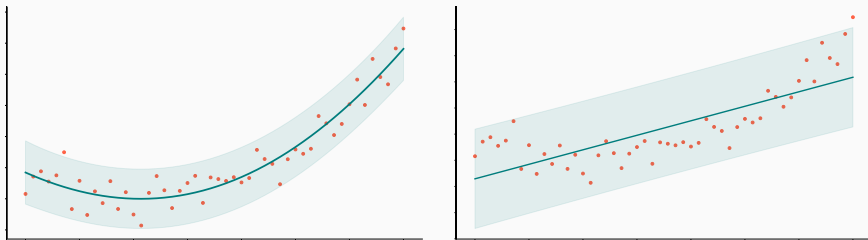
# Implicit Processes

---

# Implicit Processes

**Motivation:** Gaussian processes are limited by the parametric kernel family.

- Square exponential kernel:  $\kappa(\mathbf{x}_1, \mathbf{x}_2) = \sigma^2 \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{2l^2}\right)$ .
- Linear kernel:  $\kappa(\mathbf{x}_1, \mathbf{x}_2) = a\mathbf{x}_1^T \mathbf{x}_2 + b$ .



The Gaussian process prior over functions is **too restrictive**.

An **implicit stochastic process**<sup>1</sup> (IP) is a collection of random variables  $f(\cdot)$  such that any finite collection  $\mathbf{f} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$  is implicitly defined by the following generative process:

$$\mathbf{z} \sim P_{\mathbf{z}}(\mathbf{z}) \quad \text{and} \quad f(\mathbf{x}_n) = g_{\theta}(\mathbf{x}_n, \mathbf{z}), \quad \forall n = 1, \dots, N.$$

---

<sup>1</sup>Ma, C., Li, Y. & Hernandez-Lobato, J.M.. (2019). Variational Implicit Processes.

An **implicit stochastic process**<sup>1</sup> (IP) is a collection of random variables  $f(\cdot)$  such that any finite collection  $\mathbf{f} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$  is implicitly defined by the following generative process:

$$\mathbf{z} \sim P_{\mathbf{z}}(\mathbf{z}) \quad \text{and} \quad f(\mathbf{x}_n) = g_{\theta}(\mathbf{x}_n, \mathbf{z}), \quad \forall n = 1, \dots, N.$$

## Gaussian process

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \text{and} \quad f(\mathbf{x}_n) = \mathbf{L}(\mathbf{x}_n)^T \mathbf{z}, \quad \forall n = 1, \dots, N.$$

---

<sup>1</sup>Ma, C., Li, Y. & Hernandez-Lobato, J.M.. (2019). Variational Implicit Processes.

An **implicit stochastic process**<sup>1</sup> (IP) is a collection of random variables  $f(\cdot)$  such that any finite collection  $\mathbf{f} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$  is implicitly defined by the following generative process:

$$\mathbf{z} \sim P_{\mathbf{z}}(\mathbf{z}) \quad \text{and} \quad f(\mathbf{x}_n) = g_{\theta}(\mathbf{x}_n, \mathbf{z}), \quad \forall n = 1, \dots, N.$$

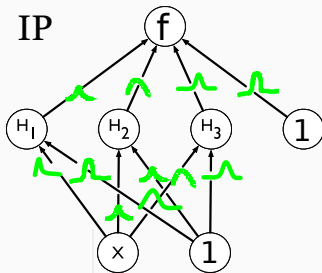
Bayesian Neural Networks.

$$(\mathbf{z}_1, \mathbf{z}_2) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\theta = (\mu_1, \mu_2, \sigma_1, \sigma_2)$$

$$\mathbf{h} = r((\mu_1 + \sigma_1 \mathbf{z}_1)^T \mathbf{x}_n).$$

$$g_{\theta}(\mathbf{x}_n, \mathbf{z}) = (\mu_2 + \sigma_2 \mathbf{z}_2)^T \mathbf{h}$$



<sup>1</sup>Ma, C., Li, Y. & Hernandez-Lobato, J.M.. (2019). Variational Implicit Processes.

# Variational Inference on Implicit Processes

## Problem statement:

1. The unknown target function **is a sample from an IP**, that is, an implicit distribution over stochastic processes  $P(f(\cdot))$ .
2. Given the set of observations  $(\mathbf{X}, \mathbf{y})$ , we aim to approximate the posterior distribution over functions  $P(f(\cdot)|\mathbf{X}, \mathbf{y})$ .

**Approach:** Use variational inference over the function-space distribution  $Q(f(\cdot))$ .

$$\begin{aligned} Q^*(f(\cdot)) &= \arg \min_{Q \in \mathcal{Q}} KL\left(Q(f(\cdot)) \mid P(f(\cdot)|\mathbf{X}, \mathbf{y})\right) \\ &= \arg \max_{Q \in \mathcal{Q}} \mathbb{E}_{Q(f)} \left[ \log P(\mathbf{y}|\mathbf{X}, f(\mathbf{X})) \right] - KL\left(Q(f(\cdot)) \mid P(f(\cdot))\right). \end{aligned}$$

**Approach:** Use variational inference over the function-space distribution

$$Q^*(f(\cdot)) = \arg \max_{Q \in \mathcal{Q}} \mathbb{E}_{Q(f)} \left[ \log P(\mathbf{y} | \mathbf{X}, f(\mathbf{X})) \right] - KL \left( Q(f(\cdot)) \mid P(f(\cdot)) \right).$$

**Difficulties:**

1. The prior  $P(f(\cdot))$  lacks a closed form.
2. The Kullback-Leibler divergence between stochastic processes is not well-defined.
3. The variational distribution  $Q(f(\cdot))$  must allow to compute or approximate by samples the data-fitting term.



**Variational Implicit Processes**<sup>2</sup>: Approximates the distribution over functions using a **linear combination of samples**.

**Sparse Implicit Processes**<sup>3</sup>: Uses **inducing points** for scalability and approximates the KL using an **external discriminator** (a Neural Network).

**Linearized approximation**<sup>4</sup>: Approximates the distribution over functions using a **linearization** of the BNN over the parameters.

---

<sup>2</sup>Ma, C., Li, Y. & Hernandez-Lobato, J.M.. (2019). Variational Implicit Processes.

<sup>3</sup>Rodríguez-Santana, S., Zaldivar, B. & Hernandez-Lobato, D.. (2022). Function-space Inference with Sparse Implicit Processes.

<sup>4</sup>Rudner, T., Chen, Z., Whye Y. & Gal, Y.. (2022). Tractable Function-Space Variational Inference in Bayesian Neural Networks.

# Variational Implicit Processes

Approximate  $P(f(\cdot))$  with a GP  $P_{\mathcal{GP}}(f(\cdot))$  based on samples  $f_1(\cdot), \dots, f_S(\cdot)$ .

Let

$$\hat{m}(\mathbf{x}) = \frac{1}{S} \sum_{s=1}^S f_s(\mathbf{x}), \quad \hat{\phi}(\mathbf{x}) = \frac{1}{\sqrt{S}} \left( f_1(\mathbf{x}) - \hat{m}(\mathbf{x}), \dots, f_S(\mathbf{x}) - \hat{m}(\mathbf{x}) \right)^T.$$

Then, setting a standard **Gaussian prior**  $P(\mathbf{a}) = \mathcal{N}(\mathbf{a}|\mathbf{0}, \mathbf{I})$ ,

$$\hat{f}(\mathbf{x}) = \hat{m}(\mathbf{x}) + \mathbf{a}^T \hat{\phi}(\mathbf{x}) \implies P_{\mathcal{GP}}(\hat{f}(\mathbf{x})) = \mathcal{N}(\hat{m}(\mathbf{x}), \hat{\phi}(\mathbf{x})^T \hat{\phi}(\mathbf{x})).$$

Using a variational distribution  $Q(\mathbf{a}) = \mathcal{N}(\mathbf{m}, \mathbf{S})$  induces a variational distribution over functions

$$Q(\hat{f}(\mathbf{x})) = \int_{\mathbf{a}} P(\hat{f}(\mathbf{x})|\mathbf{a}) Q(\mathbf{a}) = \mathcal{N}(\hat{m}(\mathbf{x}) + \hat{\phi}(\mathbf{x})^T \mathbf{m}, \hat{\phi}(\mathbf{x})^T \mathbf{S} \hat{\phi}(\mathbf{x})).$$

## Naming

$$\hat{\mathbf{f}} = (\hat{f}(\mathbf{x}_1), \dots, \hat{f}(\mathbf{x}_N)).$$

The ELBO is computed to minimize the KL divergence evaluated on  $\hat{\mathbf{f}}$  and  $\mathbf{a}$  rather than between stochastic processes

~~$$Q^*(f(\cdot)) = \arg \min_{Q \in \mathcal{Q}} KL(Q(f(\cdot)) \mid P(f(\cdot) \mid \mathbf{X}, \mathbf{y}))$$~~

$$Q^*(\hat{\mathbf{f}}, \mathbf{a}) = \arg \min_{Q \in \mathcal{Q}} KL(Q(\hat{\mathbf{f}}, \mathbf{a}) \mid P(\hat{\mathbf{f}}, \mathbf{a} \mid \mathbf{X}, \mathbf{y}))$$

$$= \arg \max_{Q \in \mathcal{Q}} \mathbb{E}_{Q(\hat{\mathbf{f}})} [\log P(\mathbf{y} \mid \mathbf{X}, \hat{\mathbf{f}})] - KL(Q(\mathbf{a}) \mid P(\mathbf{a})).$$

# Sparse Implicit Processes

1. Considers an **inducing points** approach, leading to the ELBO

$$\mathcal{L} = \mathbb{E}_{Q(\mathbf{f})} \left[ \log P(\mathbf{y}|\mathbf{X}, \mathbf{f}) \right] - KL\left(Q(\mathbf{u}) \mid P(\mathbf{u})\right).$$

2. The variational distribution of the inducing points  $Q(\mathbf{u})$  is an **IP**.
3. The Kullback-Leibler term is approximated using an **external discriminator**  $T$  (a neural network),

$$KL\left(Q(\mathbf{u}) \mid P(\mathbf{u})\right) \approx \mathbb{E}_{Q(\mathbf{u})}[T(\mathbf{u})].$$

4. The data-fitting term is **approximated using Monte-Carlo** samples of  $Q(\mathbf{u})$ ,

$$Q(\mathbf{f}) = \int_{\mathbf{u}} P(\mathbf{f}|\mathbf{u})Q(\mathbf{u}) \approx \frac{1}{S} \sum_{s=1}^S P(\mathbf{f}|\mathbf{u}_s),$$

where  $P(\mathbf{f}|\mathbf{u})$  is approximated as Gaussian with mean and covariances estimated empirically from the prior.

# Linearized approximation

Seeing the stochastic function  $f(\cdot, \boldsymbol{\theta})$  defined in terms of the stochastic parameters  $\boldsymbol{\theta}$  according to a distribution  $P_{\boldsymbol{\Theta}}$ , with

$$P_{\boldsymbol{\Theta}} = \mathcal{N}(\boldsymbol{m}, \boldsymbol{S}).$$

The stochastic function can be approximated with a **Taylor approximation of order 1 over the parameter space, centered on  $\boldsymbol{m}$ ,**

$$f(\cdot, \boldsymbol{\theta}) \approx \hat{f}(\cdot, \boldsymbol{\theta}) = f(\cdot, \boldsymbol{m}) + \mathcal{J}(\cdot, \boldsymbol{m})(\boldsymbol{\theta} - \boldsymbol{m}),$$

where

$$\mathcal{J}(\cdot, \boldsymbol{m}) = \frac{\partial f(\cdot, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}(\boldsymbol{m}).$$

As a result,  $\hat{f}(\cdot, \boldsymbol{\theta})$  is a **Gaussian process**,

$$\hat{P}(\hat{f}(\cdot, \boldsymbol{\theta})) = \mathcal{N}(f(\cdot, \boldsymbol{m}), \mathcal{J}(\cdot, \boldsymbol{m})\boldsymbol{S}\mathcal{J}(\cdot, \boldsymbol{m})^T).$$

$$\mathcal{L} = \mathbb{E}_{Q(f)} \left[ \log P(\mathbf{y} | \mathbf{X}, f(\mathbf{X})) \right] - KL \left( Q(f(\cdot)) \mid P(f(\cdot)) \right)$$

Using a variational distribution over the parameters  $Q_{\Theta}$ , **leads to a variational distribution over the stochastic functions**  $Q(f)$ , where samples can be easily taken.

The Kullback-Leibler divergence between stochastic processes can be defined using evaluations

$$KL \left( Q(f(\cdot)) \mid P(f(\cdot)) \right) = \sup_{\mathbf{C} \in 2^{\mathcal{X}}} KL \left( Q(f(\mathbf{C})) \mid P(f(\mathbf{C})) \right).$$

Therefore, **approximated empirically** using the linearized models on a set of **Context points**  $\{\mathbf{C}_1, \dots, \mathbf{C}_S\}$ ,

$$KL \left( Q(f(\cdot)) \mid P(f(\cdot)) \right) \approx \max_{\mathbf{C} \in \{\mathbf{C}_1, \dots, \mathbf{C}_S\}} KL \left( \hat{Q}(\hat{f}(\mathbf{C})) \mid \hat{P}(\hat{f}(\mathbf{C})) \right),$$

which is a KL between Gaussian distributions.

# Limitations

- **Variational implicit processes.** The linear approximation can be too strong,

$$f(\mathbf{x}) \approx \hat{m}(\mathbf{x}) + \mathbf{a}^T \hat{\phi}(\mathbf{x}).$$

- **Sparse Implicit Processes.** Relies on an external discriminator  $T$ , increasing the training time due to a double-loop training.

$$KL(Q(\mathbf{u}) \mid P(\mathbf{u})) \approx \mathbb{E}_{Q(\mathbf{u})}[T(\mathbf{u})].$$

- **Linearized approximation.** The set of **Context points**  $\{\mathbf{C}_1, \dots, \mathbf{C}_S\}$  must be defined by hand previously to any learning. With points both **in the training space and out of the training space** to ensure generalization.

$$KL(Q(f(\cdot)) \mid P(f(\cdot))) \approx \max_{\mathbf{C} \in \{\mathbf{C}_1, \dots, \mathbf{C}_S\}} KL(\hat{Q}(\hat{f}(\mathbf{C})) \mid \hat{P}(\hat{f}(\mathbf{C}))).$$

# Sparse Linearized Implicit Processes

---



# Linearized model with inducing points

We propose to use the **linearized model** along with the usage of **inducing points** to, simultaneously,

1. **Avoid using Context points** to approximate the KL divergence between stochastic processes
2. **Avoid using a discriminator** for the KL divergence between IPs.

**Features:**

1. The variational distribution over the inducing points  $Q(\mathbf{u})$  is Gaussian.
2. Both  $P(\mathbf{u})$  and  $P(\mathbf{f}|\mathbf{u})$  are approximated using the linearized model rather than samples from the IP.

How are  $P(\mathbf{u})$  and  $P(\mathbf{f}|\mathbf{u})$  approximated?

Consider the concatenation of the input features and the inducing locations  $(\mathbf{X}, \mathbf{Z})$ , and the linearized approximation of the prior evaluated on them,

$$\hat{f}((\mathbf{X}, \mathbf{Z}), \boldsymbol{\theta}) = f((\mathbf{X}, \mathbf{Z}), \mathbf{m}) + \mathcal{J}((\mathbf{X}, \mathbf{Z}), \mathbf{m})(\boldsymbol{\theta} - \mathbf{m}).$$

Then,  $\hat{f}((\mathbf{X}, \mathbf{Z}), \boldsymbol{\theta})$  is a Gaussian process,

$$\hat{P}(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left( \begin{pmatrix} f(\mathbf{X}, \mathbf{m}) \\ f(\mathbf{Z}, \mathbf{m}) \end{pmatrix}, \begin{pmatrix} \mathcal{J}(\mathbf{X}, \mathbf{m}) \mathbf{S} \mathcal{J}(\mathbf{X}, \mathbf{m})^T & \mathcal{J}(\mathbf{X}, \mathbf{m}) \mathbf{S} \mathcal{J}(\mathbf{Z}, \mathbf{m})^T \\ \mathcal{J}(\mathbf{Z}, \mathbf{m}) \mathbf{S} \mathcal{J}(\mathbf{X}, \mathbf{m})^T & \mathcal{J}(\mathbf{Z}, \mathbf{m}) \mathbf{S} \mathcal{J}(\mathbf{Z}, \mathbf{m})^T \end{pmatrix} \right).$$

where  $\hat{P}(\mathbf{f}|\mathbf{u})$  and  $\hat{P}(\mathbf{u})$  can be easily computed.

The variational posterior distribution can be computed in closed form to be Gaussian

$$Q(\mathbf{f}) = \int_{\mathbf{u}} \hat{P}(\mathbf{f}|\mathbf{u})Q(\mathbf{u}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) .$$

The ELBO can be easily computed for regression and approximated for classification

$$\mathcal{L} = \mathbb{E}_{Q(\mathbf{f})} \left[ \log P(\mathbf{y}|\mathbf{X}, \mathbf{f}) \right] - KL \left( Q(\mathbf{u}) \mid \hat{P}(\mathbf{u}) \right) .$$

# Experiments

---

# How good is the Taylor approximation?

$$f(\cdot, \boldsymbol{\theta}) \approx \hat{f}(\cdot, \boldsymbol{\theta}) = f(\cdot, \boldsymbol{m}) + \mathcal{J}(\cdot, \boldsymbol{m})(\boldsymbol{\theta} - \boldsymbol{m})$$

## Advantages:

1. The mean of the approximation is in the support of the function-space distribution.
2. Does not rely on taking samples of the prior, avoiding an hyperparameter that depends on the dimensionality of the data.

## Disadvantages:

1. The approximation is degenerate when  $\mathbb{E}_P[\boldsymbol{\theta}] = \boldsymbol{m} = \mathbf{0}$ .

Assume that

$$f(\cdot, (\mathbf{w}_1, \mathbf{w}_2)) = \mathbf{w}_2 \tanh(\mathbf{w}_1 \mathbf{x}), \quad \text{and} \quad \mathbb{E}_P[(\mathbf{w}_1, \mathbf{w}_2)] = (\mathbf{m}_1, \mathbf{m}_2) = (\mathbf{0}, \mathbf{0}).$$

then

$$\mathcal{J}(\mathbf{x}, \mathbf{m}) = \frac{\partial f(\mathbf{x}, (\mathbf{w}_1, \mathbf{w}_2))}{\partial (\mathbf{w}_1, \mathbf{w}_2)}(\mathbf{m}_1, \mathbf{m}_2) = \begin{pmatrix} \mathbf{m}_2(1 - \tanh(\mathbf{m}_1 \mathbf{x})^2) \mathbf{x} \\ \tanh(\mathbf{m}_1 \mathbf{x}) \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}$$

Meaning that

$$f(\cdot, \boldsymbol{\theta}) \approx \hat{f}(\cdot, \boldsymbol{\theta}) = f(\cdot, \mathbf{m}) + \cancel{\mathcal{J}(\cdot, \mathbf{m})(\boldsymbol{\theta} - \mathbf{m})}.$$

**However**, this can be solved using random initialization of the mean values of the prior distribution  $P_{\boldsymbol{\Theta}}$ .

# Is the Taylor GP close to the true GP?

We want to test if in cases where  $P(f(\cdot))$  is a GP, the Taylor approximation GP is a good approximation.

**Approach.** Create a Bayesian Neural Network whose implicit distribution equals that of a Gaussian Process.

1. **Squared exponential Kernel:** Single hidden layer BNN with cos activation and infinite width. Gaussian weights and uniform biases.
2. **Gaussian c.d.f activation:** Single hidden layer BNN with Gaussian c.d.f activation.

Let

$$f(\mathbf{x}) = \frac{1}{\sqrt{H}} \mathbf{w}_2^T \phi(\mathbf{w}_1^T \mathbf{x} + \mathbf{b}_1) + b_2,$$

where the dimensionality of the parameter vectors is  $H$

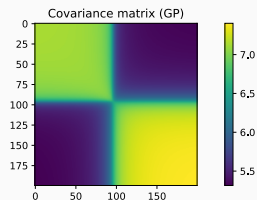
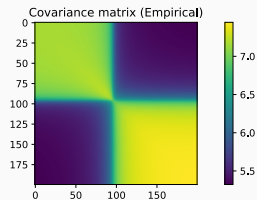
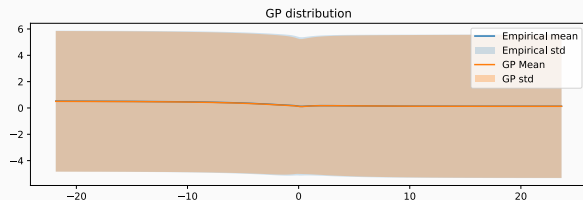
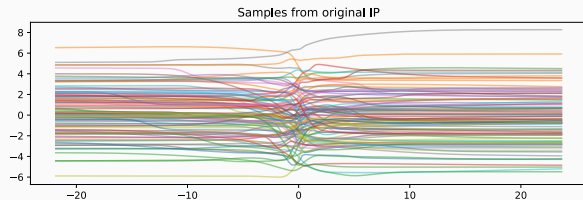
$$\begin{aligned} \mathbf{w}_1 &\sim \mathcal{N}(\mathbf{m}_{w_1}, \boldsymbol{\sigma}_{w_1} \mathbf{I}), & \mathbf{b}_1 &\sim \mathcal{N}(\mathbf{m}_{b_1}, \boldsymbol{\sigma}_{b_1} \mathbf{I}), \\ \mathbf{w}_2 &\sim \mathcal{N}(\mathbf{m}_{w_2}, \boldsymbol{\sigma}_{w_2} \mathbf{I}), & \mathbf{b}_2 &\sim \mathcal{N}(\mathbf{m}_{b_2}, \boldsymbol{\sigma}_{b_2} \mathbf{I}). \end{aligned}$$

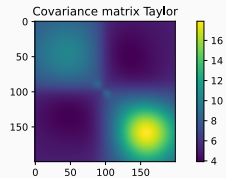
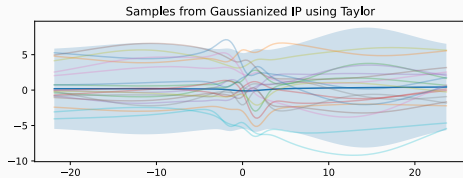
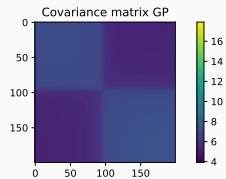
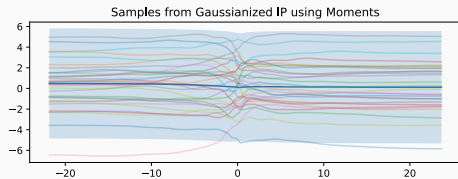
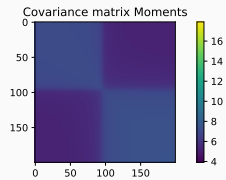
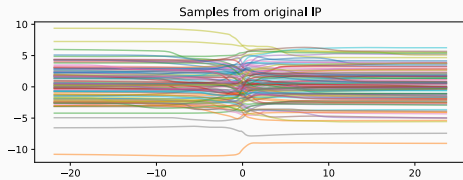
The distribution of  $f(\mathbf{x})$  tends to a GP when  $H \rightarrow \infty$ . The mean and covariance of  $f(\mathbf{x})$  can be computed using 1-dimensional quadrature.

In practise,  $H = 20$  is enough to approximate the GP.



# Showing that $H = 20$ is good enough





# Exact posterior distributions

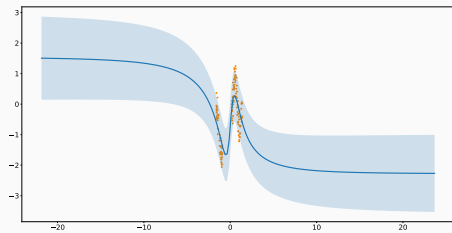
We are considering three main models based on the stochastic function

$$f(\mathbf{x}) = \frac{1}{\sqrt{H}} \mathbf{w}_2^T \phi(\mathbf{w}_1^T \mathbf{x} + \mathbf{b}_1) + b_2.$$

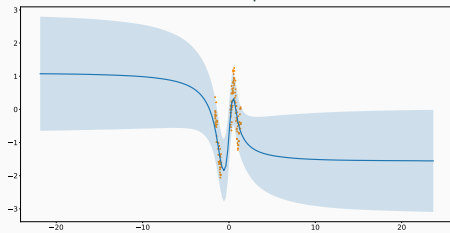
1. The **exact GP** of the stochastic function.
2. A GP where the mean and covariance are estimated using **samples** from the stochastic function.
3. A GP where the mean and covariance are the ones obtained from the **Taylor** approximation.

We are testing these three methods on a toy 1-D dataset where the exact GP posteriors can be computed.

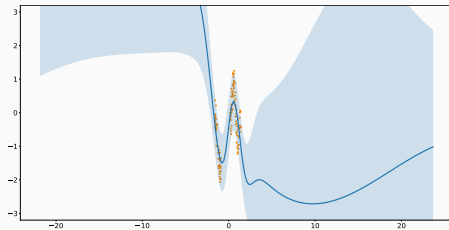
Theoretical GP



GP Samples



GP Taylor



## Conclusions

---

# Conclusions

- The Taylor approximation of order 1 is **not really close** to the exact distribution in the tested cases.
- **However**, preliminary results show that the predictive distribution close to the known points **is good enough** to encourage more research and testing.
- If the Taylor approximation results usable in practise, the proposed method would **avoid some of the problems of other function-space approaches**.
- The family of priors on which the Taylor approximation is good enough needs to be studied.

# Conclusions

- The Taylor approximation of order 1 is **not really close** to the exact distribution in the tested cases.
- **However**, preliminary results show that the predictive distribution close to the known points **is good enough** to encourage more research and testing.
- If the Taylor approximation results usable in practise, the proposed method would **avoid some of the problems of other function-space approaches**.
- The family of priors on which the Taylor approximation is good enough needs to be studied.

# Conclusions

- The Taylor approximation of order 1 is **not really close** to the exact distribution in the tested cases.
- **However**, preliminary results show that the predictive distribution close to the known points **is good enough** to encourage more research and testing.
- If the Taylor approximation results usable in practise, the proposed method would **avoid some of the problems of other function-space approaches**.
- The family of priors on which the Taylor approximation is good enough needs to be studied.



# Conclusions

- The Taylor approximation of order 1 is **not really close** to the exact distribution in the tested cases.
- **However**, preliminary results show that the predictive distribution close to the known points **is good enough** to encourage more research and testing.
- If the Taylor approximation results usable in practise, the proposed method would **avoid some of the problems of other function-space approaches**.
- The family of priors on which the Taylor approximation is good enough needs to be studied.

Thank you for your attention