



# Sparse GP Approximations

An Introduction to Pseudo-Input Methods

---

Markus Kaiser

[markus.kaiser@siemens.com](mailto:markus.kaiser@siemens.com)

5. April 2018

Siemens AG — CT RDA BAM LSY-DE  
Technical University of Munich

# Agenda

1. Standard GP Recap
2. Model Augmentation
3. SPGP: More Hyperparameters
4. SGPR: A Variational Approach
5. SVGP: Stochastic Optimization

## **Standard GP Recap**

---

# Gaussian Processes

## Definition

A **Gaussian Process (GP)** is a collection of random variables  $\{F_x\}$ , any finite subset of which has a joint Gaussian distribution.

- Extension of Gaussians to (infinite) function spaces
- $F_x$  models the function value  $f(x)$

## Mean and Kernel Functions

A GP is completely determined by two functions.

**Mean function**     $\mu_f(x) = \mathbb{E}[f(x)]$

**Kernel function**     $\mathcal{K}(x, x') = \text{cov}[f(x), f(x')]$

# GP Posterior

## GP predictive posterior

The **predictive posterior** for a function  $f \sim \mathcal{GP}(\mathbf{0}, \mathcal{K})$  with observations  $\mathbf{y} = f(\mathbf{X}) + \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I})$  is

$$p(f_* | X, y, X_*) = \mathcal{N}(f_* | \mu_*, \Sigma_*), \text{ with}$$

$$\mu_* = K_{*f} (K_{ff} + \sigma_n^2 \mathbf{I})^{-1} y$$

$$\Sigma_* = K_{**} - K_{*f} (K_{ff} + \sigma_n^2 \mathbf{I})^{-1} K_{f*}.$$

## GP predictive posterior

The **predictive posterior** for a function  $f \sim \mathcal{GP}(\mathbf{0}, \mathcal{K})$  with observations  $\mathbf{y} = f(\mathbf{X}) + \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I})$  is

$$p(f_* | X, y, X_*) = \mathcal{N}(f_* | \mu_*, \Sigma_*), \text{ with}$$

$$\mu_* = K_{*f} (K_{ff} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\Sigma_* = K_{**} - K_{*f} (K_{ff} + \sigma_n^2 \mathbf{I})^{-1} K_{f*}.$$

- We always need to consider the **whole dataset**
- The matrix inversion is  $\mathcal{O}(N^3)$  in time and  $\mathcal{O}(N^2)$  in space
- With precomputation we still have to pay  $\mathcal{O}(N^2)$  for  $\Sigma_*$

# GP Marginal Likelihood

## GP marginal likelihood

We usually use the **marginal likelihood**  $\mathcal{L}$  for model selection with maximum likelihood

$$\begin{aligned}\mathcal{L}^{\text{GP}}(\theta) &= -\log p(y \mid \theta) \\ &= -\log \int p(y \mid f, \theta) p(f \mid \theta) df \\ &= -\log \int \mathcal{N}(y \mid f, \sigma_n^2 I) \mathcal{N}(f \mid \mathbf{0}, K_{ff}) df \\ &= -\log \mathcal{N}(y \mid \mathbf{0}, K_{ff} + \sigma_n^2 I) \\ &= \frac{1}{2} \mathbf{y}^\top (K_{ff} + \sigma_n^2 I)^{-1} \mathbf{y} + \frac{1}{2} \log |K_{ff} + \sigma_n^2 I| + \frac{N}{2} \log(2\pi).\end{aligned}$$

# GP Marginal Likelihood

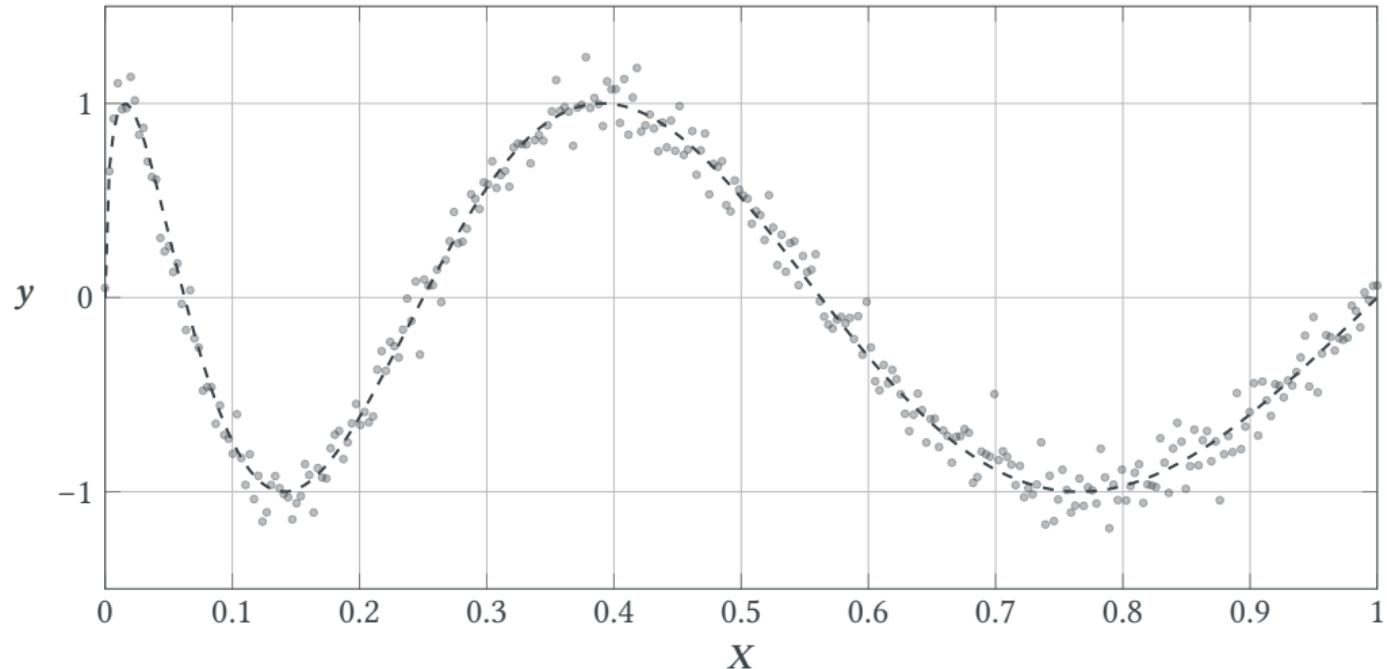
## GP marginal likelihood

We usually use the **marginal likelihood**  $\mathcal{L}$  for model selection with maximum likelihood

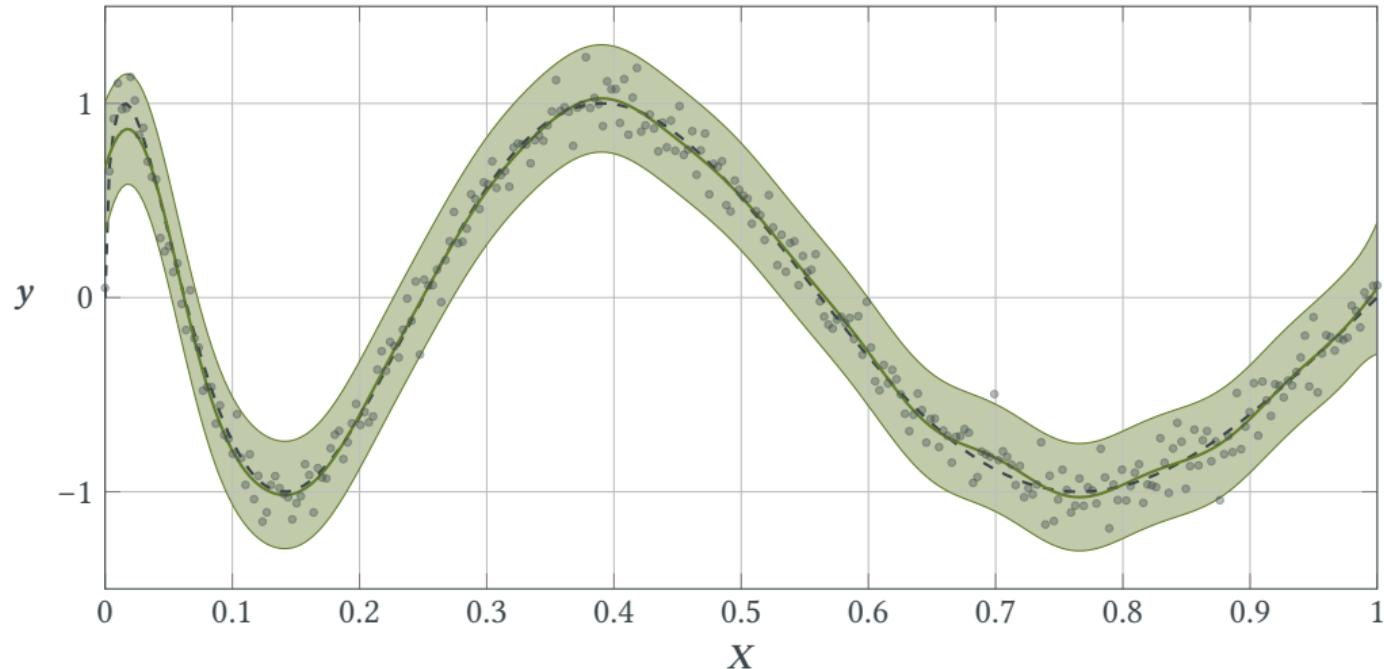
$$\begin{aligned}\mathcal{L}^{\text{GP}}(\theta) &= -\log p(y \mid \theta) \\ &= -\log \int p(y \mid f, \theta) p(f \mid \theta) df \\ &= -\log \int \mathcal{N}(y \mid f, \sigma_n^2 I) \mathcal{N}(f \mid \mathbf{0}, K_{ff}) df \\ &= -\log \mathcal{N}(y \mid \mathbf{0}, K_{ff} + \sigma_n^2 I) \\ &= \frac{1}{2} \mathbf{y}^\top (K_{ff} + \sigma_n^2 I)^{-1} \mathbf{y} + \frac{1}{2} \log |K_{ff} + \sigma_n^2 I| + \frac{N}{2} \log(2\pi).\end{aligned}$$

- Choose  $\theta^* \in \operatorname{argmin}_{\theta} \mathcal{L}(\theta)$
- $\mathcal{L}$  is also  $\mathcal{O}(N^3)$  in time and  $\mathcal{O}(N^2)$  in space

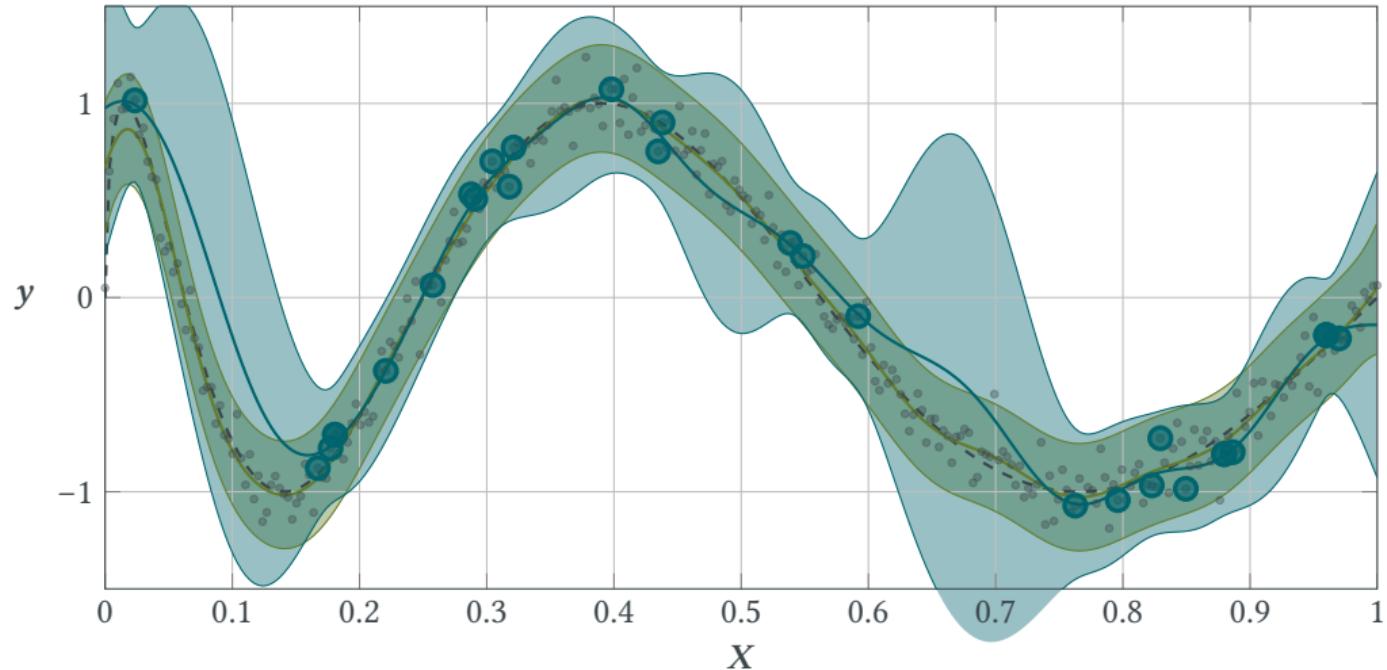
## Reducing the Data Set



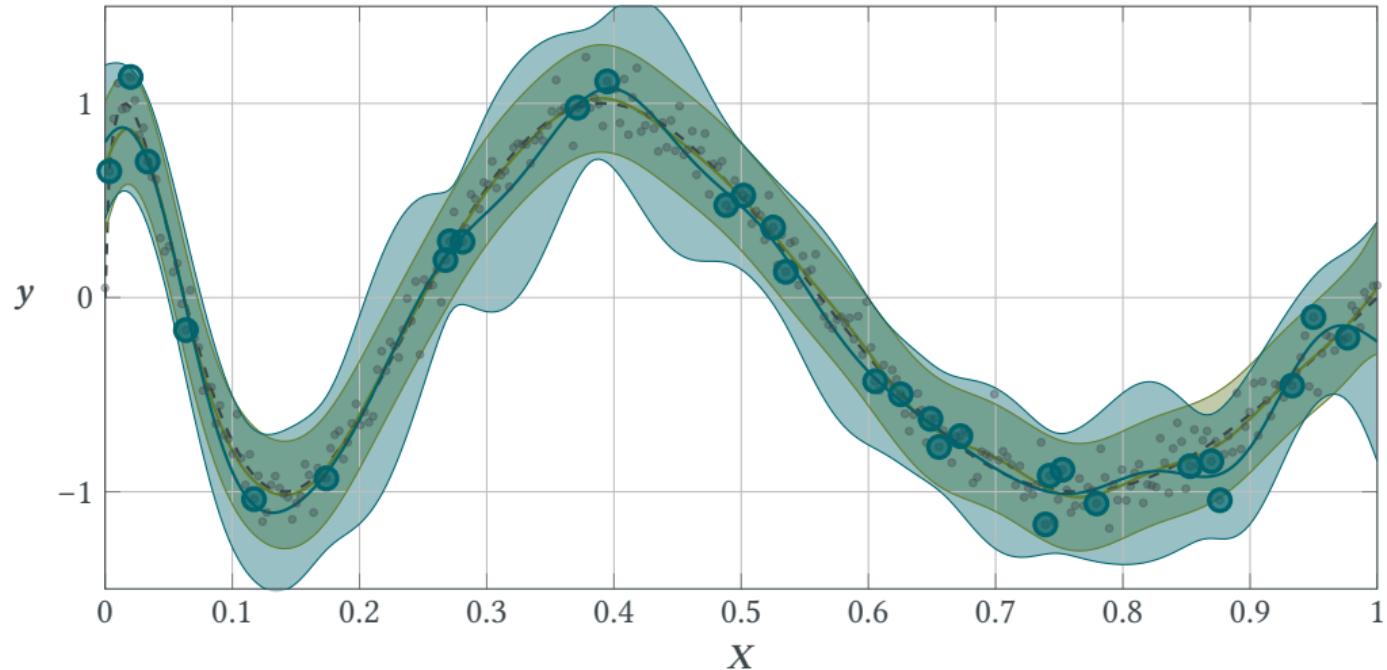
## Reducing the Data Set



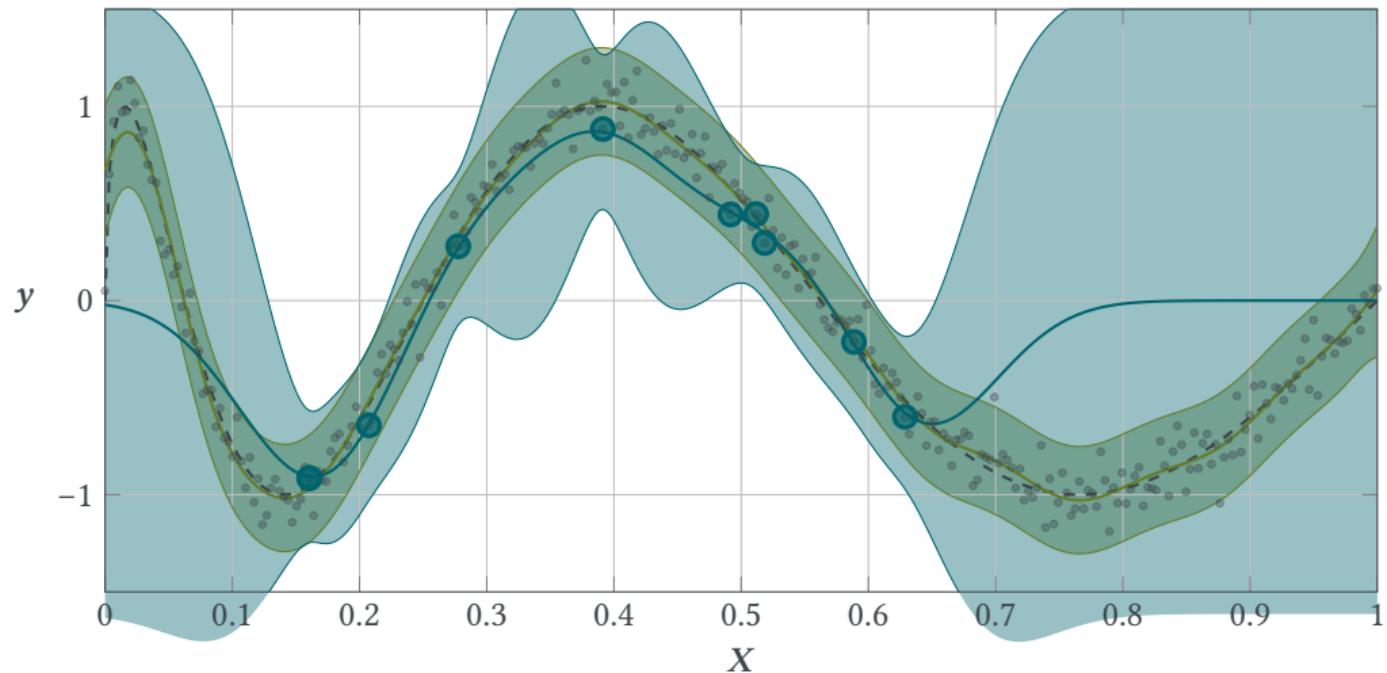
## Reducing the Data Set



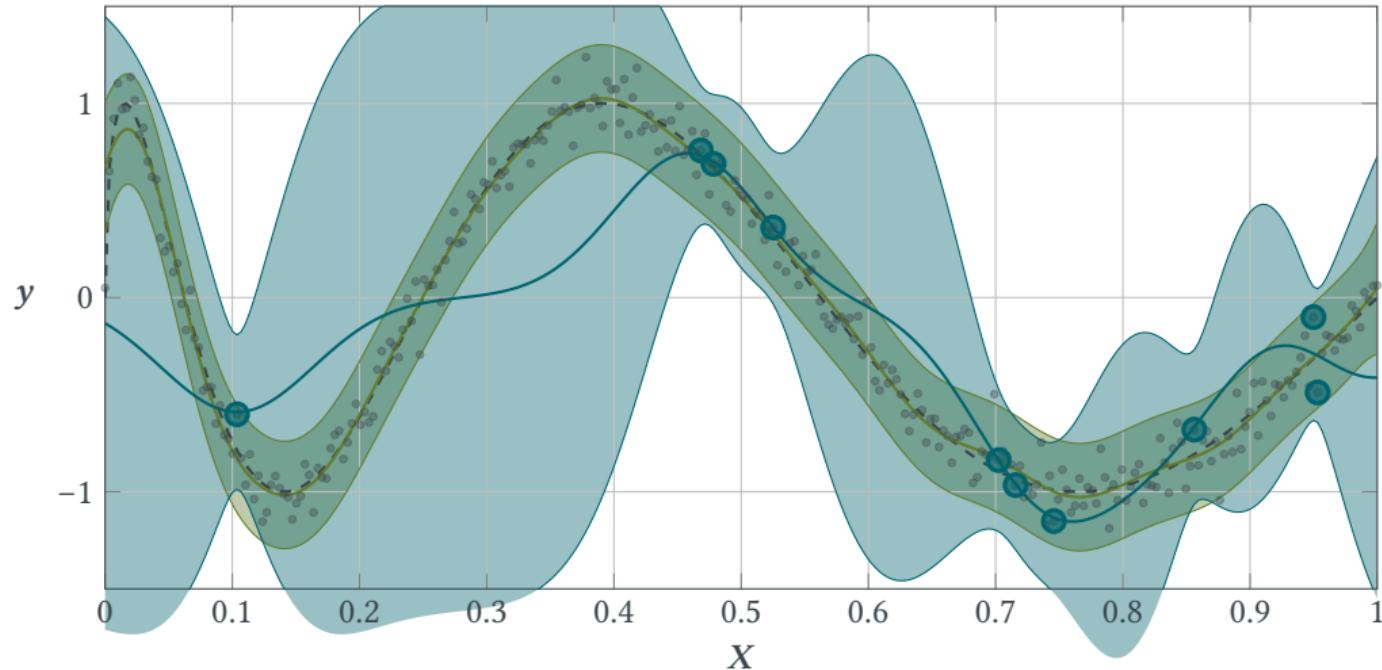
## Reducing the Data Set



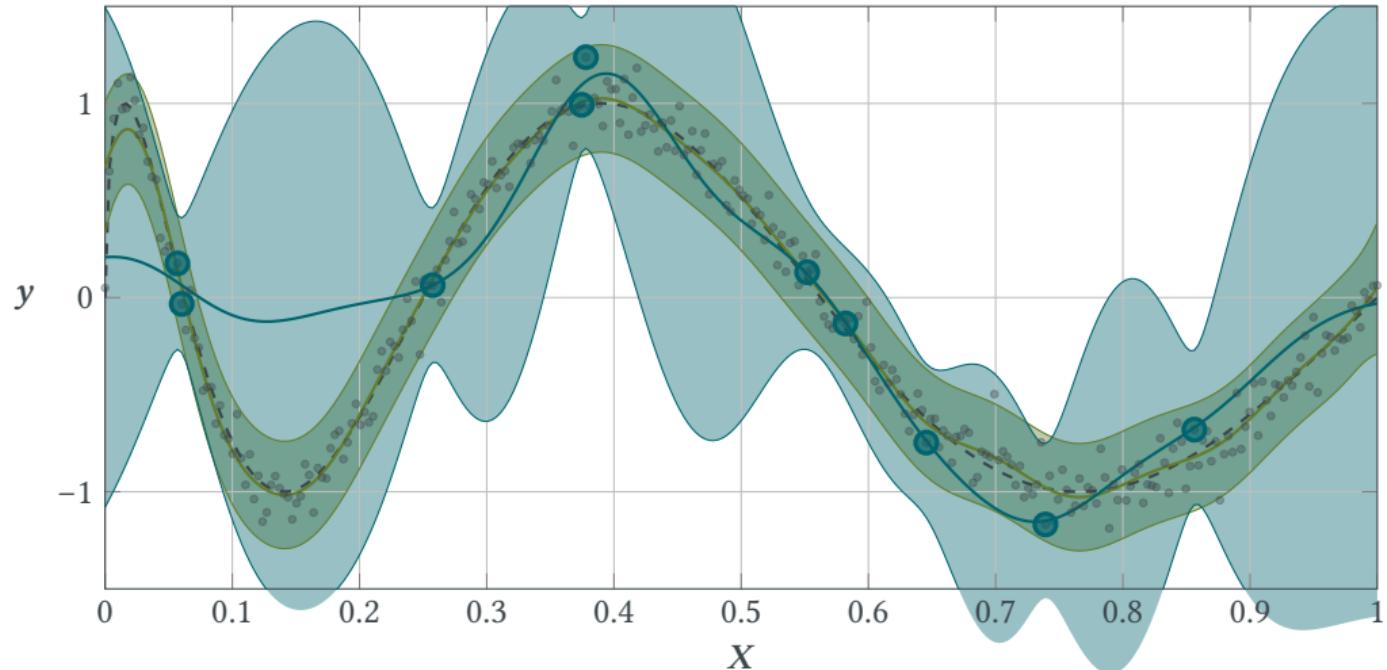
## Reducing the Data Set



## Reducing the Data Set



## Reducing the Data Set



Can we somehow not ignore most of the data?

## Model Augmentation

---

## Pseudo Data

- Keep  $N$  original observations  $\mathcal{D} = (X, y)$
- Assume  $M$  different **pseudo data points**  $\bar{\mathcal{D}} = (Z, u)$
- $u = f(Z)$  are **observations of the latent function**
- The  $\bar{\mathcal{D}}$  are also called **inducing points**

## Pseudo Data

- Keep  $N$  original observations  $\mathcal{D} = (X, y)$
- Assume  $M$  different **pseudo data points**  $\bar{\mathcal{D}} = (Z, u)$
- $u = f(Z)$  are **observations of the latent function**
- The  $\bar{\mathcal{D}}$  are also called **inducing points**

### Augmented GP Joint

Due to the consistency of GPs, the data and pseudo data are **jointly Gaussian**.

$$p(f, u) = \mathcal{N}\left(\begin{pmatrix} f \\ u \end{pmatrix} \middle| 0, \begin{pmatrix} K_{ff} & K_{fu} \\ K_{uf} & K_{uu} \end{pmatrix}\right)$$

- We usually drop the conditioning on  $X, Z$  and  $\theta$

## Conditioning on the Pseudo Data

### Augmented GP conditional likelihood

Because  $p(f, u)$  is Gaussian, we can calculate the **conditional likelihood** easily.

$$\begin{aligned} p(y | u) &= \int p(y | f) p(f | u) df \\ &= \int \mathcal{N}(y | f, \sigma_n^2 I) \mathcal{N}(f | K_{fu} K_{uu}^{-1} u, K_{ff} - K_{fu} K_{uu}^{-1} K_{uf}) df \\ &= \mathcal{N}(y | K_{fu} K_{uu}^{-1} u, K_{ff} - K_{fu} K_{uu}^{-1} K_{uf} + \sigma_n^2 I) \\ &= \mathcal{N}(y | K_{fu} K_{uu}^{-1} u, K_{ff} - Q_{ff} + \sigma_n^2 I) \end{aligned}$$

- We call  $Q_{ff}$  the low rank covariance

## Conditioning on the Pseudo Data

### Augmented GP conditional likelihood

Because  $p(f, u)$  is Gaussian, we can calculate the **conditional likelihood** easily.

$$\begin{aligned} p(y | u) &= \int p(y | f) p(f | u) df \\ &= \int \mathcal{N}(y | f, \sigma_n^2 I) \mathcal{N}(f | K_{fu} K_{uu}^{-1} u, K_{ff} - K_{fu} K_{uu}^{-1} K_{uf}) df \\ &= \mathcal{N}(y | K_{fu} K_{uu}^{-1} u, K_{ff} - K_{fu} K_{uu}^{-1} K_{uf} + \sigma_n^2 I) \\ &= \mathcal{N}(y | K_{fu} K_{uu}^{-1} u, K_{ff} - Q_{ff} + \sigma_n^2 I) \end{aligned}$$

- We call  $Q_{ff}$  the low rank covariance
- This is just the **predictive distribution** again
- Assuming iid data, we can calculate the likelihood in (at most)  $\mathcal{O}(N^2)$

But what use is it if we have no idea about  $\bar{D}$ ?

## **SPGP: More Hyperparameters**

---

## Sparse Pseudo-input Gaussian Process (SPGP)

We interpret the augmented model as a **new regression problem** arising from  $\bar{\mathcal{D}}$ .

- Assume **conditional independence of  $f$  (and  $y$ ) given  $u$**
- We can choose  $Z$  and only have to handle  $u$

## Sparse Pseudo-input Gaussian Process (SPGP)

We interpret the augmented model as a **new regression problem** arising from  $\bar{\mathcal{D}}$ .

- Assume **conditional independence of  $f$  (and  $y$ ) given  $u$**
- We can choose  $Z$  and only have to handle  $u$

$$\begin{aligned} p(y | u) &= \prod_n p(y_n | u) \\ &= \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} u, \text{diag}(K_{ff} - Q_{ff}) + \sigma_n^2 I) \end{aligned}$$

Using the **low rank covariance**

$$Q_{ab} := K_{au} K_{uu}^{-1} K_{ub}$$

# SPGP Marginal Likelihood

## SPGP Marginal Likelihood

$$\begin{aligned}\mathcal{L}^{\text{SPGP}}(\theta) &= -\log p(y \mid \theta) \\ &= -\log \int p(y \mid u) p(u) du \\ &= -\log \int \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} u, \text{diag}(K_{ff} - Q_{ff}) + \sigma_n^2 I) \mathcal{N}(u \mid 0, K_{uu}) du \\ &= -\log \mathcal{N}(y \mid 0, Q_{ff} + \text{diag}(K_{ff} - Q_{ff}) + \sigma_n^2 I)\end{aligned}$$

- We do not know  $u = f(Z)$ , so we **marginalize**
- Use the original **GP prior**  $p(u \mid \theta) = \mathcal{N}(u \mid 0, K_{uu})$
- **Joint optimization of  $Z$**  with other hyperparameters
- Time-complexity is  $\mathcal{O}(NM^2)$ , **linear** in the data

## Interpretation of the SPGP

$$\mathcal{L}^{\text{GP}}(\theta) = -\log \mathcal{N}(y \mid 0, \mathbf{K}_{ff} + \sigma_n^2 \mathbf{I})$$

$$\mathcal{L}^{\text{SPGP}}(\theta) = -\log \mathcal{N}(y \mid 0, \mathbf{Q}_{ff} + \text{diag}(\mathbf{K}_{ff} - \mathbf{Q}_{ff}) + \sigma_n^2 \mathbf{I})$$

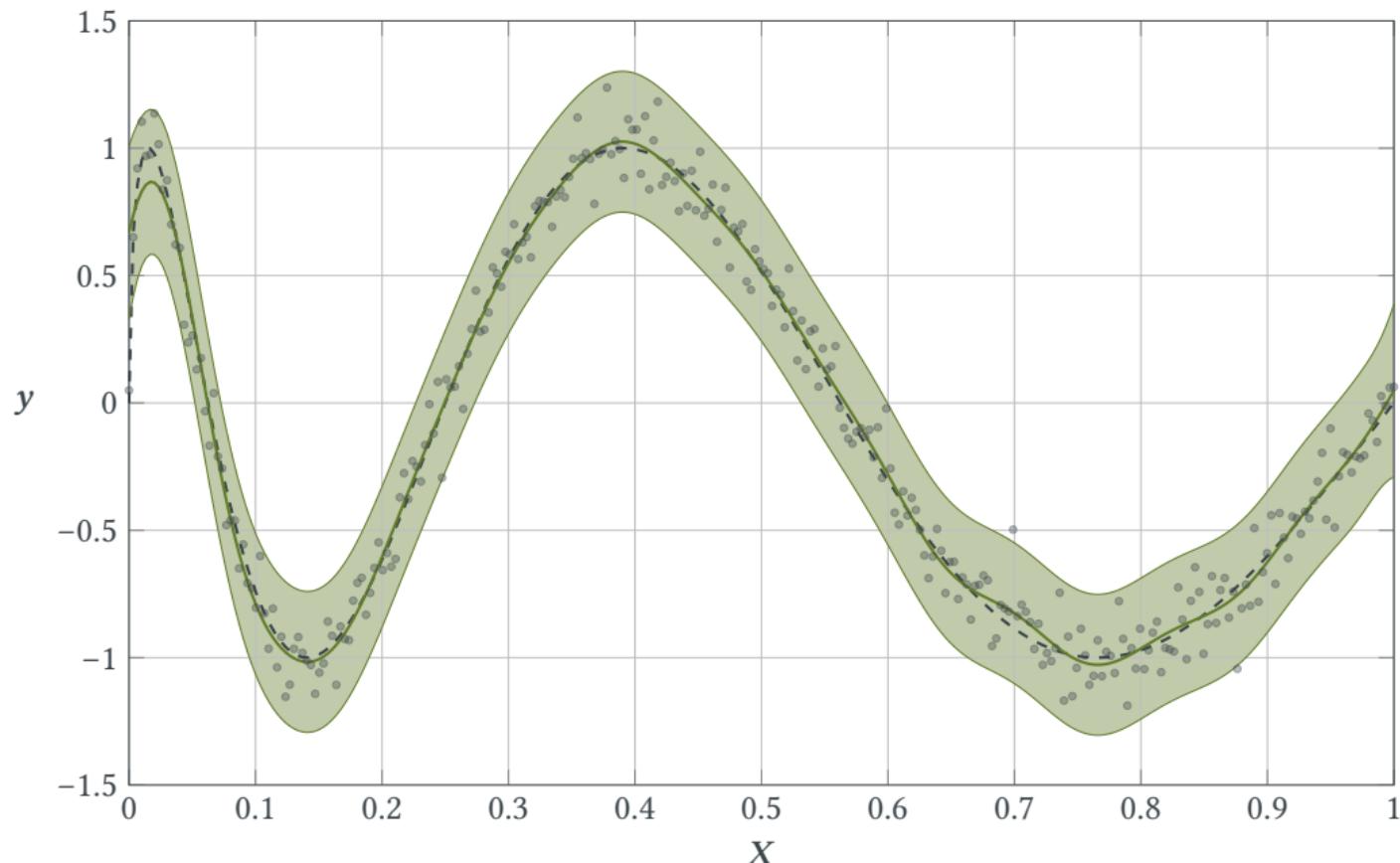
### SPGP kernel

The SPGP derived from  $\mathcal{GP}(0, \mathcal{K})$  is a new Gaussian Process  $\mathcal{GP}(0, \mathcal{K}^{\text{SPGP}})$  with

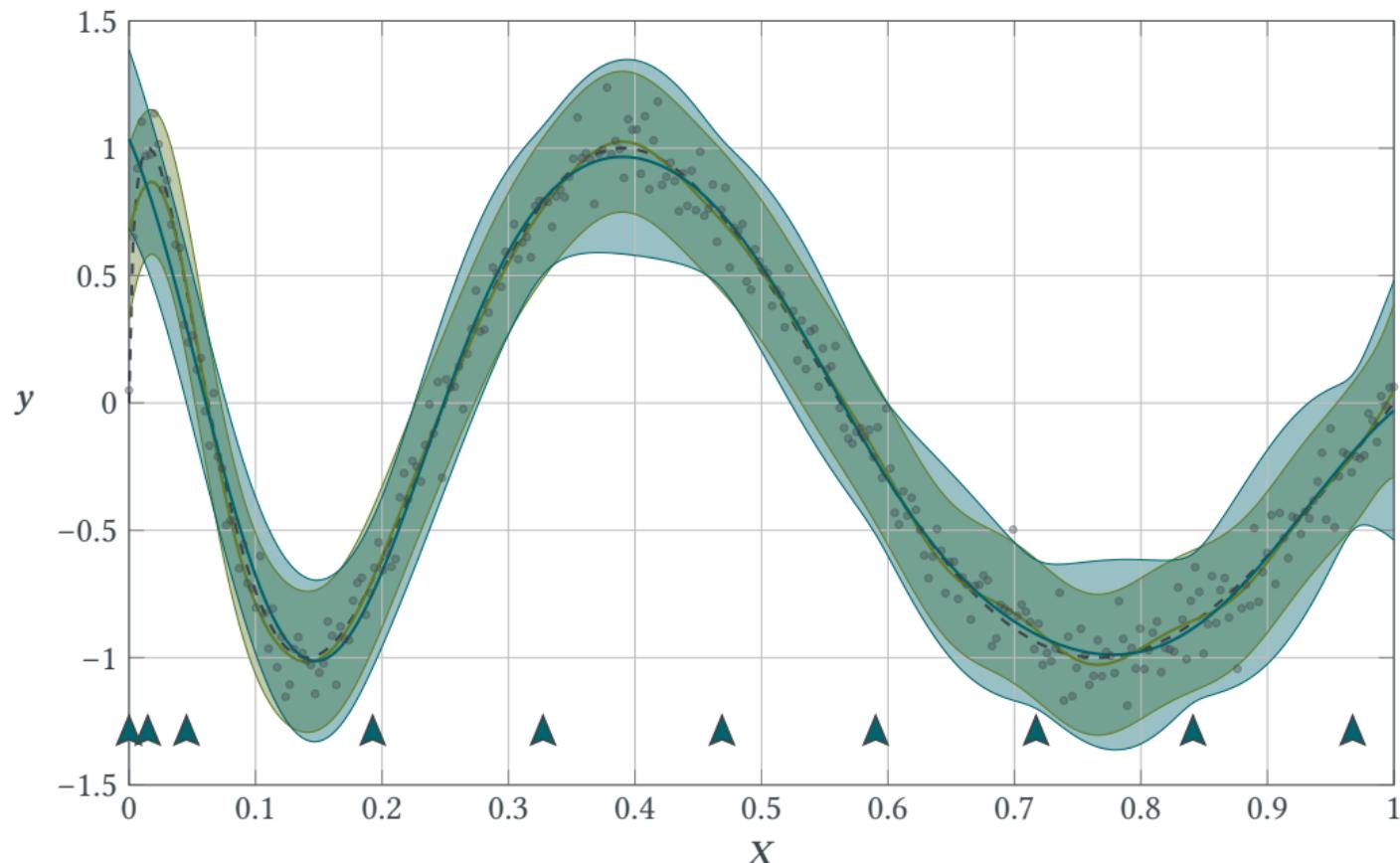
$$\begin{aligned}\mathcal{K}^{\text{SPGP}}(x, x') &= \mathcal{Q}(x, x') + \delta_{xx'} (\mathcal{K}(x, x') - \mathcal{Q}(x, x')) \\ \mathcal{Q}(x, x') &= K_{xu} K_{uu}^{-1} K_{ux'}\end{aligned}$$

- $\mathcal{K}^{\text{SPGP}}$  is heavily dependent on  $Z$  and no longer interpretable
- It has a potentially **huge number of (hyper-)parameters**...
- ...which makes the model prone to **overfitting**

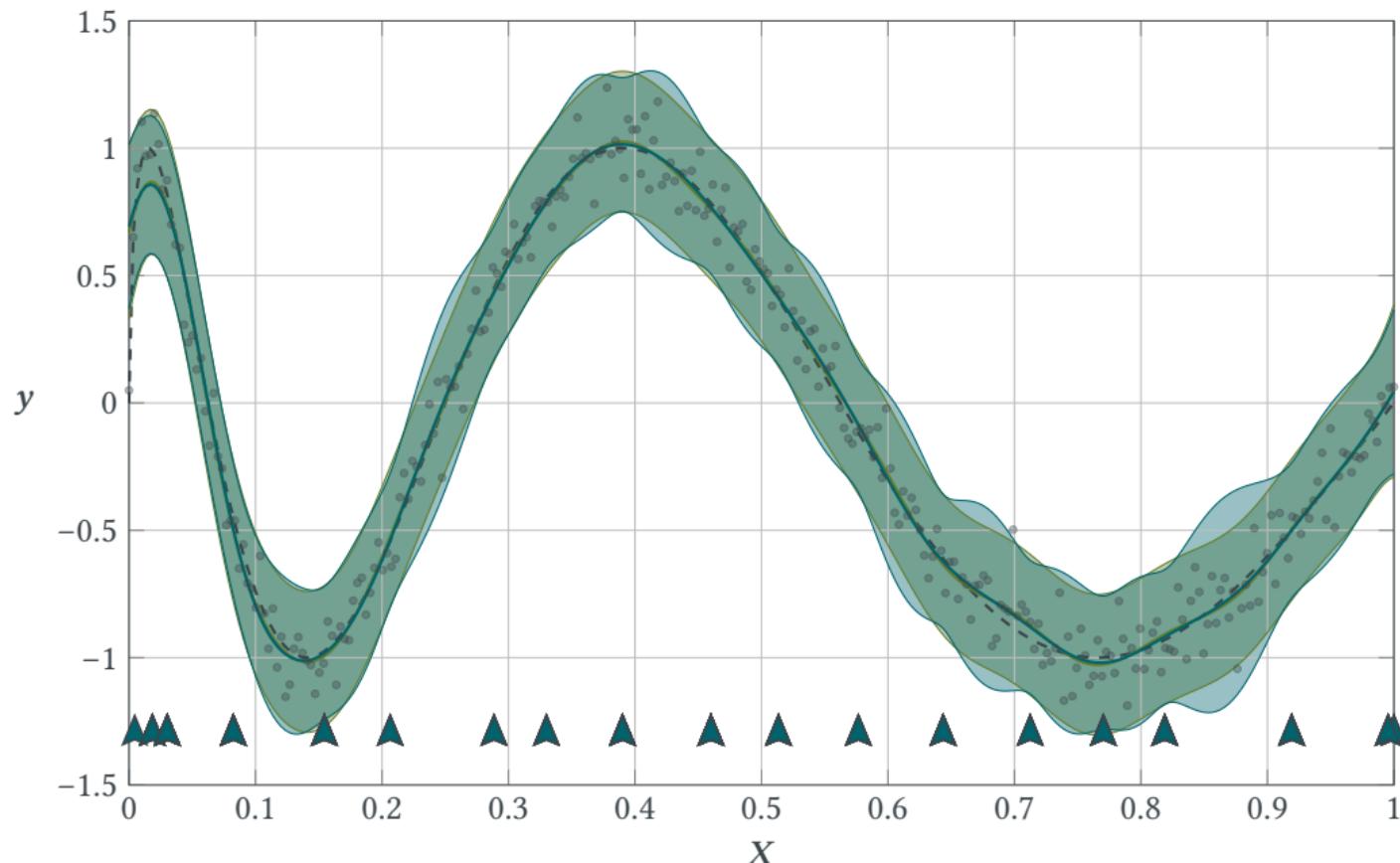
# SPGP Showcase



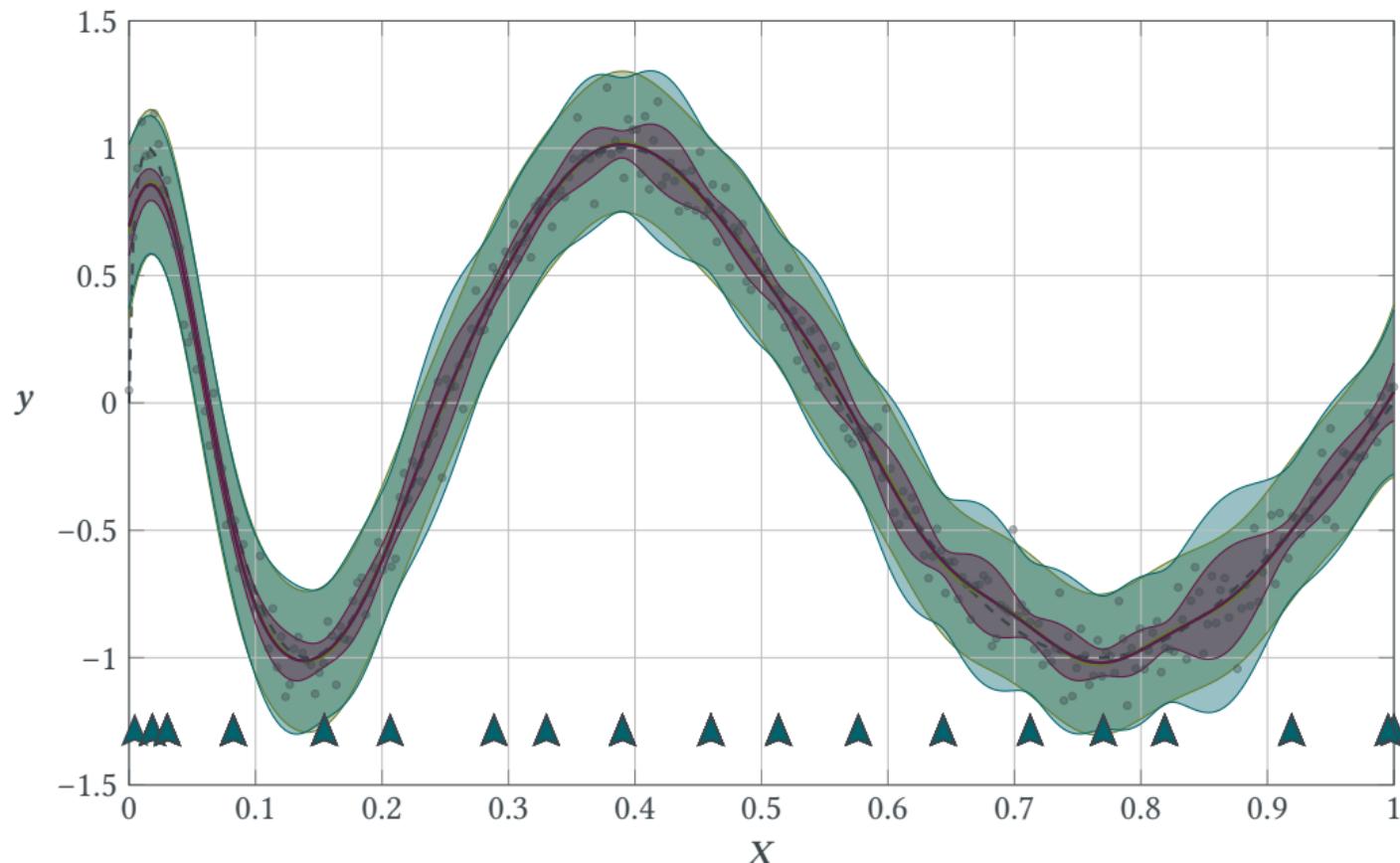
# SPGP Showcase



# SPGP Showcase



# SPGP Showcase



## **SGPR: A Variational Approach**

---

## Sparse Gaussian Process Regression (SGPR)

We want to variationally approximate the **original Gaussian Process**.

- Derive a **variational lower bound** on the original marginal likelihood  $\mathcal{L}^{\text{GP}}$
- Consider a **custom variational scheme**  $q(f, u) = p(f | u) q(u)$

## Sparse Gaussian Process Regression (SGPR)

We want to variationally approximate the **original Gaussian Process**.

- Derive a **variational lower bound** on the original marginal likelihood  $\mathcal{L}^{\text{GP}}$
- Consider a **custom variational scheme**  $q(f, u) = p(f | u) q(u)$
- Inspect the augmented predictive posterior

$$p(f^* | y) = \int p(f^* | f, u) p(f, u | y) df du$$

- **$u$  is optimal** if  $f^*$  and  $f$  are **independent given  $u$**

$$p(f^* | f, u) = p(f^* | u)$$

$$p(f, u | y) = p(f | u) p(u | y)$$

## Sparse Gaussian Process Regression (SGPR)

We want to variationally approximate the **original Gaussian Process**.

- Derive a **variational lower bound** on the original marginal likelihood  $\mathcal{L}^{\text{GP}}$
- Consider a **custom variational scheme**  $q(f, u) = p(f | u) q(u)$
- Inspect the augmented predictive posterior

$$p(f^* | y) = \int p(f^* | f, u) p(f, u | y) df du$$

- **$u$  is optimal** if  $f^*$  and  $f$  are **independent given  $u$**

$$p(f^* | f, u) = p(f^* | u)$$

$$p(f, u | y) = p(f | u) p(u | y)$$

- We will approximate this situation with

$$q(f^*) = \int p(f^* | u) p(f | u) q(u) df du = \int p(f^* | u) q(u) du$$

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(u)) &= -\log p(y | \theta, Z, u) \\ &= -\log \int p(y | f) p(f | u) p(u) df du \\ &= -\log \int q(f, u) \frac{p(y | f) p(f | u) p(u)}{q(f, u)} df du \\ &\geq -\int q(f, u) \log \frac{p(y | f) p(f | u) p(u)}{q(f, u)} df du \\ &= -\int p(f | u) q(u) \log \frac{p(y | f) p(f | u) p(u)}{p(f | u) q(u)} df du \\ &= -\int p(f | u) q(u) \log \frac{p(y | f) p(u)}{q(u)} df du \\ &= -\int q(u) \left( \int p(f | u) \log p(y | f) df + \log \frac{p(u)}{q(u)} \right) du\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(u)) &= -\log p(y \mid \theta, Z, u) \\&= -\log \int p(y \mid f) p(f \mid u) p(u) df du \\&= -\log \int q(f, u) \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&\geq -\int q(f, u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&= -\int p(f \mid u) q(u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{p(f \mid u) q(u)} df du \\&= -\int p(f \mid u) q(u) \log \frac{p(y \mid f) p(u)}{q(u)} df du \\&= -\int q(u) \left( \int p(f \mid u) \log p(y \mid f) df + \log \frac{p(u)}{q(u)} \right) du\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(u)) &= -\log p(y \mid \theta, Z, u) \\&= -\log \int p(y \mid f) p(f \mid u) p(u) df du \\&= -\log \int q(f, u) \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&\geq - \int q(f, u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{p(f \mid u) q(u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(u)}{q(u)} df du \\&= - \int q(u) \left( \int p(f \mid u) \log p(y \mid f) df + \log \frac{p(u)}{q(u)} \right) du\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(u)) &= -\log p(y \mid \theta, Z, u) \\&= -\log \int p(y \mid f) p(f \mid u) p(u) df du \\&= -\log \int q(f, u) \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&\geq - \int q(f, u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{p(f \mid u) q(u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(u)}{q(u)} df du \\&= - \int q(u) \left( \int p(f \mid u) \log p(y \mid f) df + \log \frac{p(u)}{q(u)} \right) du\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(u)) &= -\log p(y \mid \theta, Z, u) \\&= -\log \int p(y \mid f) p(f \mid u) p(u) df du \\&= -\log \int q(f, u) \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&\geq - \int q(f, u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{p(f \mid u) q(u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(u)}{q(u)} df du \\&= - \int q(u) \left( \int p(f \mid u) \log p(y \mid f) df + \log \frac{p(u)}{q(u)} \right) du\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(u)) &= -\log p(y \mid \theta, Z, u) \\&= -\log \int p(y \mid f) p(f \mid u) p(u) df du \\&= -\log \int q(f, u) \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&\geq - \int q(f, u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{q(f, u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(f \mid u) p(u)}{p(f \mid u) q(u)} df du \\&= - \int p(f \mid u) q(u) \log \frac{p(y \mid f) p(u)}{q(u)} df du \\&= - \int q(u) \left( \int p(f \mid u) \log p(y \mid f) df + \log \frac{p(u)}{q(u)} \right) du\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathbb{E}_{p(f|\boldsymbol{u})}[\log p(\boldsymbol{y}|f)] &= \int p(f|\boldsymbol{u}) \log p(\boldsymbol{y}|f) df \\ &= \int \log \mathcal{N}(y|f, \sigma_n^2 \mathbf{I}) \mathcal{N}(f|K_{fu} K_{uu}^{-1} \boldsymbol{u}, K_{ff} - Q_{ff}) df \\ &= \log \mathcal{N}(y|K_{fu} K_{uu}^{-1} \boldsymbol{u}, \sigma_n^2 \mathbf{I}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &= \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathbb{E}_{p(f|\boldsymbol{u})}[\log p(\boldsymbol{y}|f)] &= \int p(f|\boldsymbol{u}) \log p(\boldsymbol{y}|f) df \\ &= \int \log \mathcal{N}(\boldsymbol{y} | f, \sigma_n^2 \mathbf{I}) \mathcal{N}(f | \mathbf{K}_{fu} \mathbf{K}_{uu}^{-1} \boldsymbol{u}, \mathbf{K}_{ff} - \mathbf{Q}_{ff}) df \\ &= \log \mathcal{N}(\boldsymbol{y} | \mathbf{K}_{fu} \mathbf{K}_{uu}^{-1} \boldsymbol{u}, \sigma_n^2 \mathbf{I}) - \frac{1}{2\sigma_n^2} \text{tr}(\mathbf{K}_{ff} - \mathbf{Q}_{ff}) \\ &= \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(\mathbf{K}_{ff} - \mathbf{Q}_{ff})\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathbb{E}_{p(f|\boldsymbol{u})}[\log p(\boldsymbol{y}|f)] &= \int p(f|\boldsymbol{u}) \log p(\boldsymbol{y}|f) df \\ &= \int \log \mathcal{N}(\boldsymbol{y} | f, \sigma_n^2 \mathbf{I}) \mathcal{N}(f | K_{fu} K_{uu}^{-1} \boldsymbol{u}, K_{ff} - Q_{ff}) df \\ &= \log \mathcal{N}(\boldsymbol{y} | K_{fu} K_{uu}^{-1} \boldsymbol{u}, \sigma_n^2 \mathbf{I}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &= \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathbb{E}_{p(f|\boldsymbol{u})}[\log p(\boldsymbol{y}|f)] &= \int p(f|\boldsymbol{u}) \log p(\boldsymbol{y}|f) df \\ &= \int \log \mathcal{N}(\boldsymbol{y} | f, \sigma_n^2 \mathbf{I}) \mathcal{N}(f | K_{fu} K_{uu}^{-1} \boldsymbol{u}, K_{ff} - Q_{ff}) df \\ &= \log \mathcal{N}(\boldsymbol{y} | K_{fu} K_{uu}^{-1} \boldsymbol{u}, \sigma_n^2 \mathbf{I}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &= \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(\boldsymbol{u})) &= - \int q(\boldsymbol{u}) \left( \int p(f \mid \boldsymbol{u}) \log p(y \mid f) df + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \left( \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \log \frac{G(\boldsymbol{u}) p(\boldsymbol{u})}{q(\boldsymbol{u})} d\boldsymbol{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

- Phew!
- We have not specified  $q(\boldsymbol{u})$ . How should we choose it?

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(\boldsymbol{u})) &= - \int q(\boldsymbol{u}) \left( \int p(f \mid \boldsymbol{u}) \log p(y \mid f) df + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \left( \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(\mathbf{K}_{ff} - \mathbf{Q}_{ff}) + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \log \frac{G(\boldsymbol{u}) p(\boldsymbol{u})}{q(\boldsymbol{u})} d\boldsymbol{u} + \frac{1}{2\sigma_n^2} \text{tr}(\mathbf{K}_{ff} - \mathbf{Q}_{ff})\end{aligned}$$

- Phew!
- We have not specified  $q(\boldsymbol{u})$ . How should we choose it?

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(\boldsymbol{u})) &= - \int q(\boldsymbol{u}) \left( \int p(f \mid \boldsymbol{u}) \log p(y \mid f) df + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \left( \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \log \frac{G(\boldsymbol{u}) p(\boldsymbol{u})}{q(\boldsymbol{u})} d\boldsymbol{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

- Phew!
- We have not specified  $q(\boldsymbol{u})$ . How should we choose it?

## SGPR variational bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z, q(\boldsymbol{u})) &= - \int q(\boldsymbol{u}) \left( \int p(f \mid \boldsymbol{u}) \log p(y \mid f) df + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \left( \log G(\boldsymbol{u}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) + \log \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} \right) d\boldsymbol{u} \\ &= - \int q(\boldsymbol{u}) \log \frac{G(\boldsymbol{u}) p(\boldsymbol{u})}{q(\boldsymbol{u})} d\boldsymbol{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

- Phew!
- We have not specified  $q(\boldsymbol{u})$ . How should we choose it?

## The optimal $q(\mathbf{u})$

$$\begin{aligned}\min_{q(\mathbf{u})} \mathcal{L}^{\text{SGPR}}(\theta, Z, q(\mathbf{u})) &= \min_{q(\mathbf{u})} - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &= \min_{q(\mathbf{u})} - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} \\ &= \min_{q(\mathbf{u})} \text{KL}(q(\mathbf{u}) \| G(\mathbf{u}) p(\mathbf{u}))\end{aligned}$$

## The optimal $q(\mathbf{u})$

$$\begin{aligned}\min_{q(\mathbf{u})} \mathcal{L}^{\text{SGPR}}(\theta, Z, q(\mathbf{u})) &= \min_{q(\mathbf{u})} - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &= \min_{q(\mathbf{u})} - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} \\ &= \min_{q(\mathbf{u})} \text{KL}(q(\mathbf{u}) \| G(\mathbf{u}) p(\mathbf{u}))\end{aligned}$$

- Minimize the divergence by choosing

$$q^*(\mathbf{u}) \propto G(\mathbf{u}) p(\mathbf{u}) = \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} \mathbf{u}, \sigma_n^2 \mathbf{I}) \mathcal{N}(\mathbf{u} \mid \mathbf{0}, K_{uu})$$

## The optimal $q(\mathbf{u})$

$$\begin{aligned}\min_{q(\mathbf{u})} \mathcal{L}^{\text{SGPR}}(\theta, Z, q(\mathbf{u})) &= \min_{q(\mathbf{u})} - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &= \min_{q(\mathbf{u})} - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} \\ &= \min_{q(\mathbf{u})} \text{KL}(q(\mathbf{u}) \| G(\mathbf{u}) p(\mathbf{u}))\end{aligned}$$

- Minimize the divergence by choosing

$$q^*(\mathbf{u}) \propto G(\mathbf{u}) p(\mathbf{u}) = \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} \mathbf{u}, \sigma_n^2 \mathbf{I}) \mathcal{N}(\mathbf{u} \mid \mathbf{0}, K_{uu})$$

- This is Gaussian again

$$q^*(\mathbf{u}) = \frac{G(\mathbf{u}) p(\mathbf{u})}{\int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u}} = \mathcal{N}(\mathbf{u} \mid \boldsymbol{\mu}_u, \boldsymbol{\Sigma}_u)$$

$$\boldsymbol{\mu}_u = \sigma_n^{-2} K_{mm} B^{-1} \mathbf{K}_{mn} \mathbf{y} \quad \boldsymbol{\Sigma}_u = K_{mm} B^{-1} K_{mm} \quad B = K_{mm} + \sigma_n^{-2} \mathbf{K}_{mn} \mathbf{K}_{nm}$$

(Slight abuse of notation)

## SGPR variational bound

- Resubstituting  $q^*$  into  $\mathcal{L}^{\text{SGPR}}$  yields the final bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z) &\geq \mathcal{L}^{\text{SGPR}}(\theta, Z, q^*(u)) \\&= - \int q^*(u) \log \frac{G(u) p(u)}{q^*(u)} du + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \frac{G(u) p(u)}{\int G(u) p(u) du} \log \frac{G(u) p(u)}{\int G(u) p(u) du} du + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \int G(u) p(u) du + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \mathcal{N}(y | 0, Q_{ff} + \sigma_n^2 I) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SGPR variational bound

- Resubstituting  $q^*$  into  $\mathcal{L}^{\text{SGPR}}$  yields the final bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z) &\geq \mathcal{L}^{\text{SGPR}}(\theta, Z, q^*(\mathbf{u})) \\&= - \int q^*(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q^*(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \frac{G(\mathbf{u}) p(\mathbf{u})}{\int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u}} \log \frac{G(\mathbf{u}) p(\mathbf{u})}{\int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u}} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \mathcal{N}(y | 0, Q_{ff} + \sigma_n^2 I) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SGPR variational bound

- Resubstituting  $q^*$  into  $\mathcal{L}^{\text{SGPR}}$  yields the final bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z) &\geq \mathcal{L}^{\text{SGPR}}(\theta, Z, q^*(\mathbf{u})) \\&= - \int q^*(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q^*(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \frac{G(\mathbf{u}) p(\mathbf{u})}{\int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u}} \log \frac{G(\mathbf{u}) p(\mathbf{u})}{\frac{G(\mathbf{u}) p(\mathbf{u})}{\int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u}}} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \mathcal{N}(y \mid 0, Q_{ff} + \sigma_n^2 I) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SGPR variational bound

- Resubstituting  $q^*$  into  $\mathcal{L}^{\text{SGPR}}$  yields the final bound

$$\begin{aligned}\mathcal{L}^{\text{SGPR}}(\theta, Z) &\geq \mathcal{L}^{\text{SGPR}}(\theta, Z, q^*(\mathbf{u})) \\&= - \int q^*(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q^*(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \frac{G(\mathbf{u}) p(\mathbf{u})}{\int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u}} \log \frac{G(\mathbf{u}) p(\mathbf{u})}{\int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u}} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \int G(\mathbf{u}) p(\mathbf{u}) d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \log \mathcal{N}(\mathbf{y} \mid \mathbf{0}, Q_{ff} + \sigma_n^2 \mathbf{I}) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## Interpretation of SGPR

$$\mathcal{L}^{\text{GP}}(\theta) = -\log \mathcal{N}(y \mid \mathbf{0}, K_{ff} + \sigma_n^2 \mathbf{I})$$

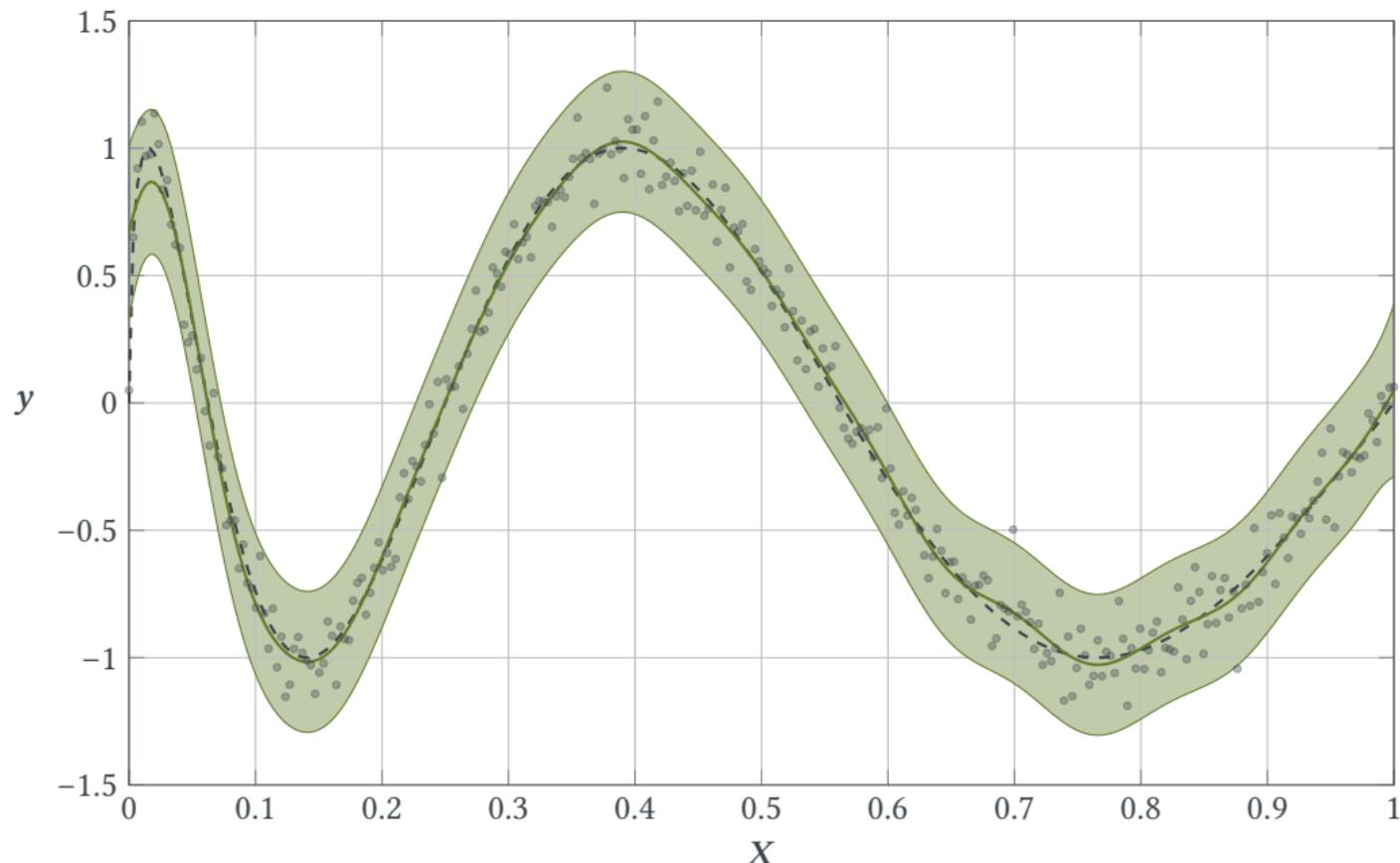
$$\mathcal{L}^{\text{SPGP}}(\theta) = -\log \mathcal{N}(y \mid \mathbf{0}, Q_{ff} + \text{diag}(K_{ff} - Q_{ff}) + \sigma_n^2 \mathbf{I})$$

### SGPR variational bound

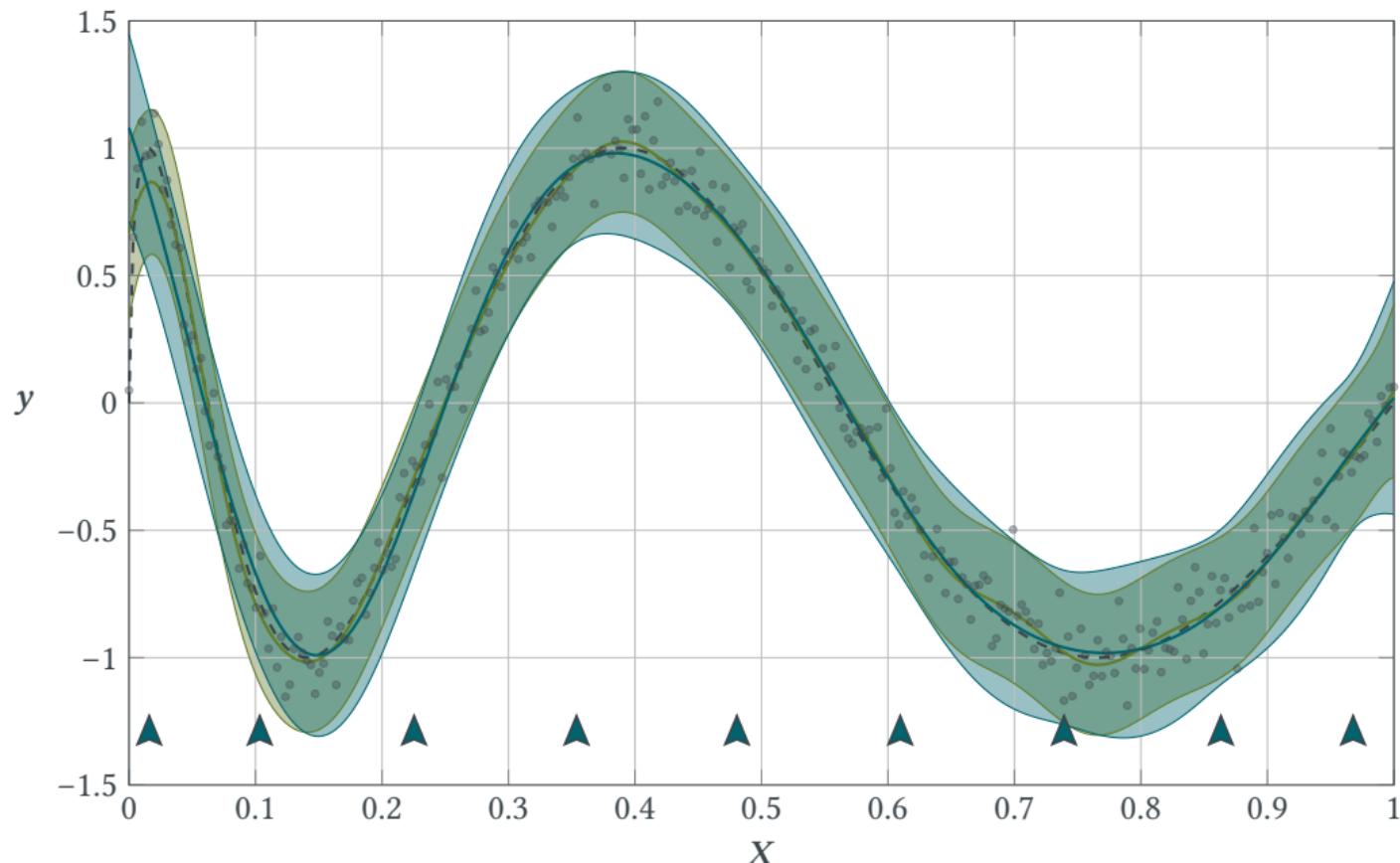
$$\mathcal{L}^{\text{SGPR}}(\theta, Z) \geq -\log \mathcal{N}(y \mid \mathbf{0}, Q_{ff} + \sigma_n^2 \mathbf{I}) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})$$

- This likelihood is  $\mathcal{O}(NM^2)$ , same with approximate predictions
- It does **not factorize** along the data, which makes stochastic optimization impossible
- The bound is tight if  $\bar{\mathcal{D}} = \mathcal{D}$
- Choosing good  $\bar{\mathcal{D}}$  does only make the bound tighter and cannot lead to overfitting

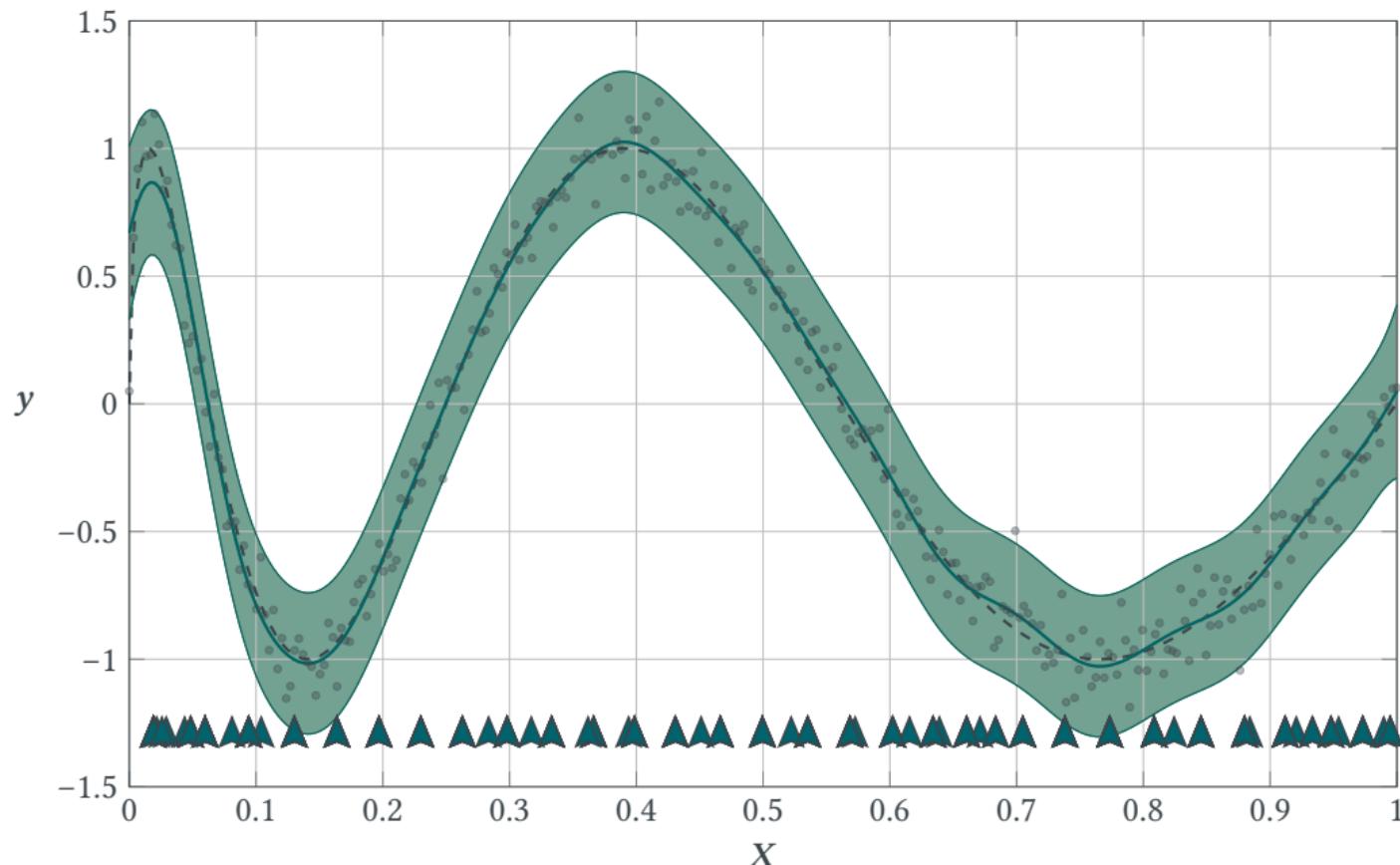
## SGPR Showcase



# SGPR Showcase



## SGPR Showcase



## **SVGP: Stochastic Optimization**

---

## Scalable Variational Gaussian Process (SVGP)

We want to variationally approximate the **original Gaussian Process**.

- Use the same strategy for the **variational lower bound** as SGPR
- Consider a **custom variational scheme**  $q(f, u) = p(f | u) q(u)$
- Instead of choosing  $q^*(u)$ , optimize  $q(u) = \mathcal{N}(u | m, S)$

## Scalable Variational Gaussian Process (SVGP)

We want to variationally approximate the **original Gaussian Process**.

- Use the same strategy for the **variational lower bound** as SGPR
- Consider a **custom variational scheme**  $q(f, u) = p(f | u) q(u)$
- Instead of choosing  $q^*(u)$ , optimize  $q(u) = \mathcal{N}(u | m, S)$

- $q^*$  is optimal but expensive to evaluate
- Instead we find a cheaper but worse approximation via optimization
- This hurts the bound, but does not alter the model
- $S$  has **quadratically many parameters**

## SVGP variational bound

$$\begin{aligned}\mathcal{L}^{\text{SVGP}}(\theta, Z, \mathbf{m}, S) &\geq - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \log G(\mathbf{u}) q(\mathbf{u}) d\mathbf{u} + \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \mathbb{E}_{q(\mathbf{u})}[\log G(\mathbf{u})] + \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SVGP variational bound

$$\begin{aligned}\mathcal{L}^{\text{SVGP}}(\theta, Z, \mathbf{m}, S) &\geq - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \log G(\mathbf{u}) q(\mathbf{u}) d\mathbf{u} + \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \mathbb{E}_{q(\mathbf{u})} [\log G(\mathbf{u})] + \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SVGP variational bound

$$\begin{aligned}\mathcal{L}^{\text{SVGP}}(\theta, Z, \mathbf{m}, S) &\geq - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \log G(\mathbf{u}) q(\mathbf{u}) d\mathbf{u} + \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \mathbb{E}_{q(u)}[\log G(u)] + \text{KL}(q(u) \| p(u)) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SVGP variational bound

$$\begin{aligned}\mathcal{L}^{\text{SVGP}}(\theta, Z, \mathbf{m}, S) &\geq - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int q(\mathbf{u}) \log \frac{G(\mathbf{u}) p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \int \log G(\mathbf{u}) q(\mathbf{u}) d\mathbf{u} + \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\&= - \mathbb{E}_{q(\mathbf{u})}[\log G(\mathbf{u})] + \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})\end{aligned}$$

## SVGP variational bound

$$\begin{aligned}\mathbb{E}_{q(\boldsymbol{u})}[\log G(\boldsymbol{u})] &= \int \log G(\boldsymbol{u}) q(\boldsymbol{u}) d\boldsymbol{u} \\ &= \int \log \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} \boldsymbol{u}, \sigma_n^2 \mathbf{I}) \mathcal{N}(\boldsymbol{u} \mid \boldsymbol{m}, S) d\boldsymbol{u} \\ &= \log \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} \boldsymbol{m}, \sigma_n^2 \mathbf{I}) - \frac{1}{2\sigma_n^2} \text{tr}(K_{fu} K_{uu}^{-1} S K_{uu}^{-1} K_{uf})\end{aligned}$$

## SVGP variational bound

$$\begin{aligned}\mathcal{L}^{\text{SVGP}}(\theta, Z, \mathbf{m}, S) &\geq -\mathbb{E}_{q(\mathbf{u})}[\log G(\mathbf{u})] + \text{KL}(q(\mathbf{u}) \parallel p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &= -\log \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} \mathbf{m}, \sigma_n^2 \mathbf{I}) + \text{KL}(q(\mathbf{u}) \parallel p(\mathbf{u})) \\ &\quad + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) + \frac{1}{2\sigma_n^2} \text{tr}(K_{fu} K_{uu}^{-1} S K_{uu}^{-1} K_{uf})\end{aligned}$$

## SVGP interpretation

$$\mathcal{L}^{\text{GP}}(\theta) = -\log \mathcal{N}(y \mid \mathbf{0}, K_{ff} + \sigma_n^2 \mathbf{I})$$

$$\mathcal{L}^{\text{SPGP}}(\theta) = -\log \mathcal{N}(y \mid \mathbf{0}, Q_{ff} + \text{diag}(K_{ff} - Q_{ff}) + \sigma_n^2 \mathbf{I})$$

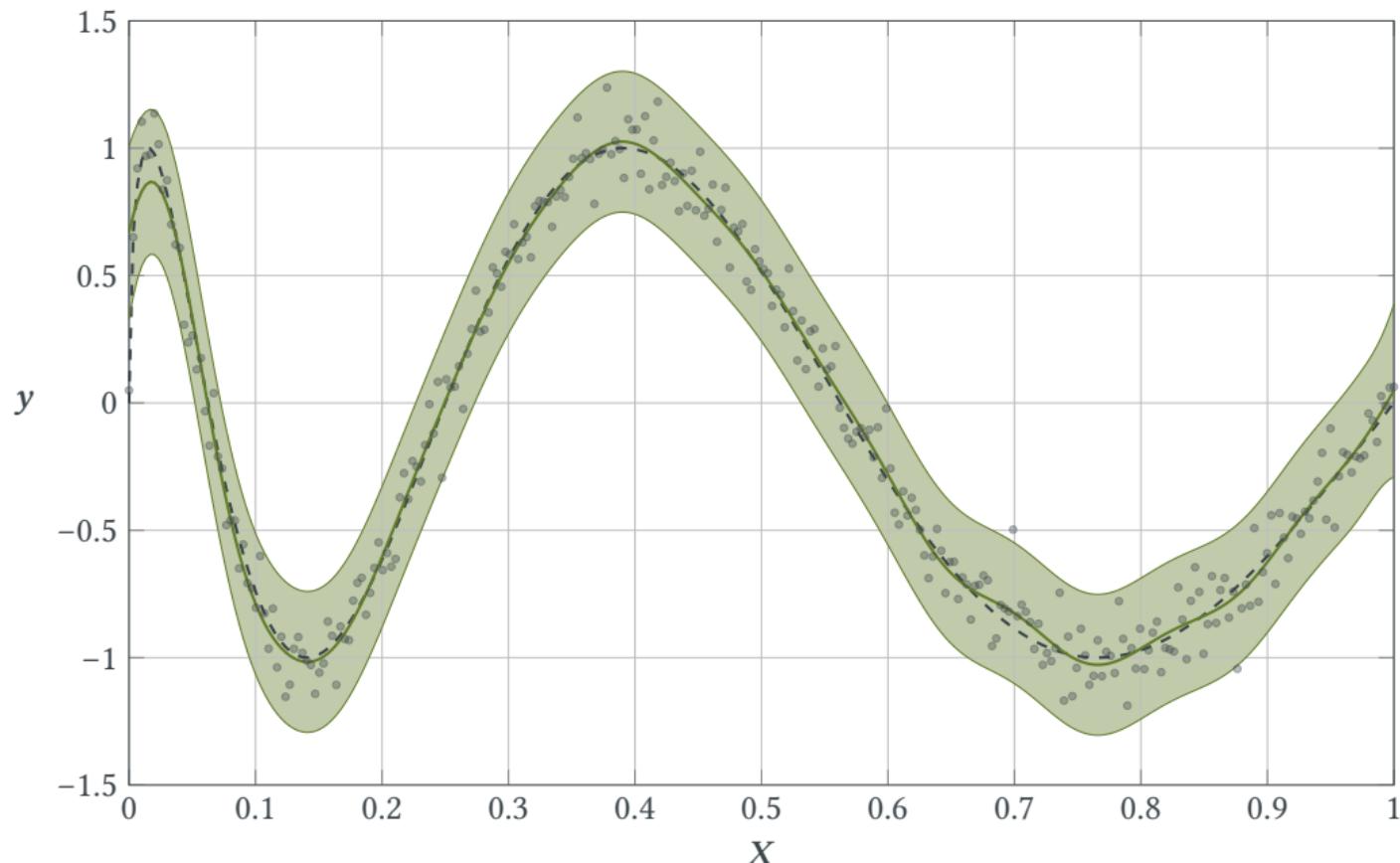
$$\mathcal{L}^{\text{SGPR}}(\theta, Z) \geq -\log \mathcal{N}(y \mid \mathbf{0}, Q_{ff} + \sigma_n^2 \mathbf{I}) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff})$$

### SVGP variational bound

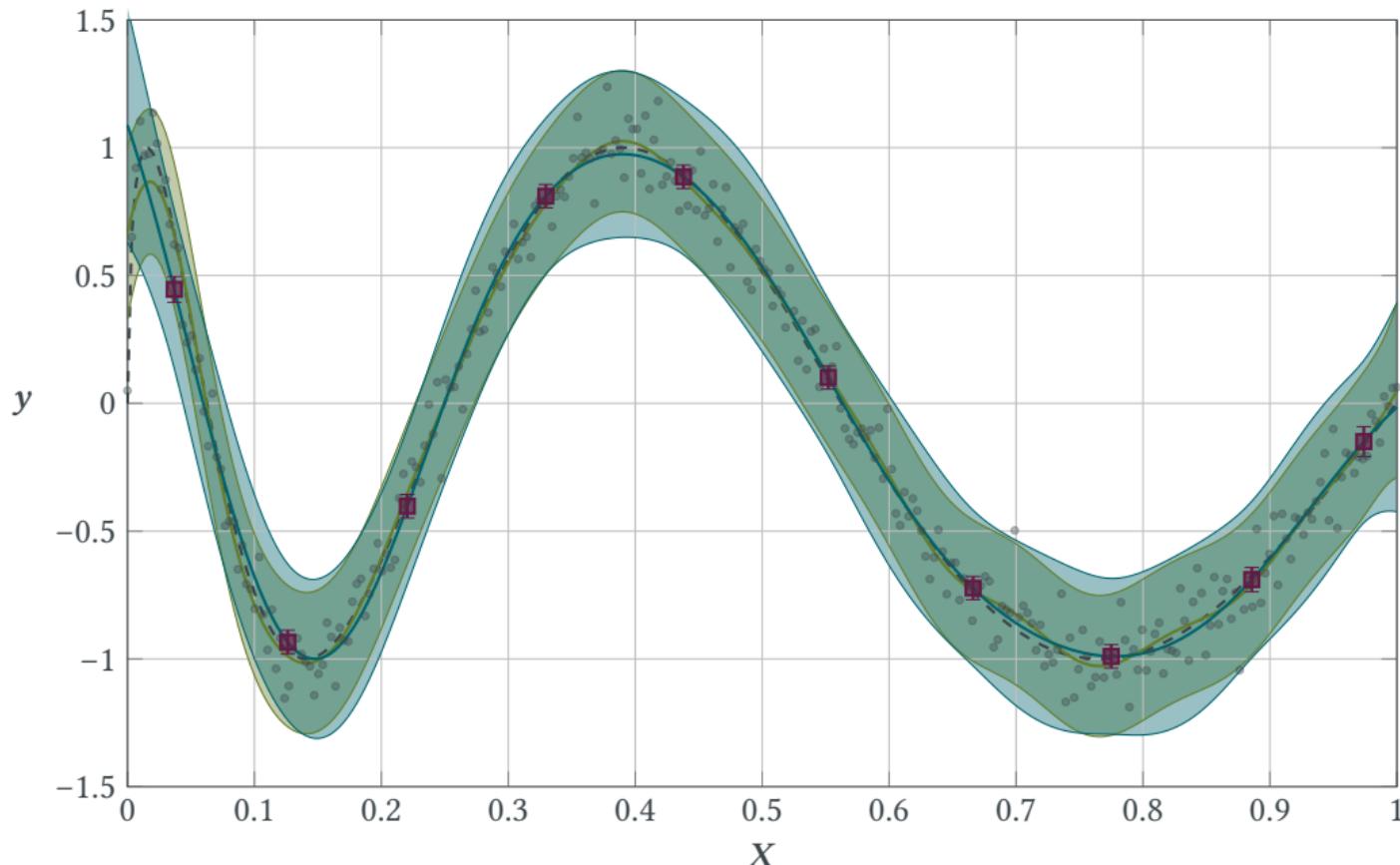
$$\begin{aligned} \mathcal{L}^{\text{SVGP}}(\theta, Z, \mathbf{m}, S) &\geq -\log \mathcal{N}(y \mid K_{fu} K_{uu}^{-1} \mathbf{m}, \sigma_n^2 \mathbf{I}) + \frac{1}{2\sigma_n^2} \text{tr}(K_{ff} - Q_{ff}) \\ &\quad + \text{KL}(q(\mathbf{u}) \parallel p(\mathbf{u})) + \frac{1}{2\sigma_n^2} \text{tr}(K_{fu} K_{uu}^{-1} S K_{uu}^{-1} K_{uf}) \end{aligned}$$

- The complete likelihood is still  $\mathcal{O}(NM^2)$
- But it **factorizes** along the data - we can do stochastic optimization!
- A single prediction is only  $\mathcal{O}(M^2)$  **and we do not need the data**

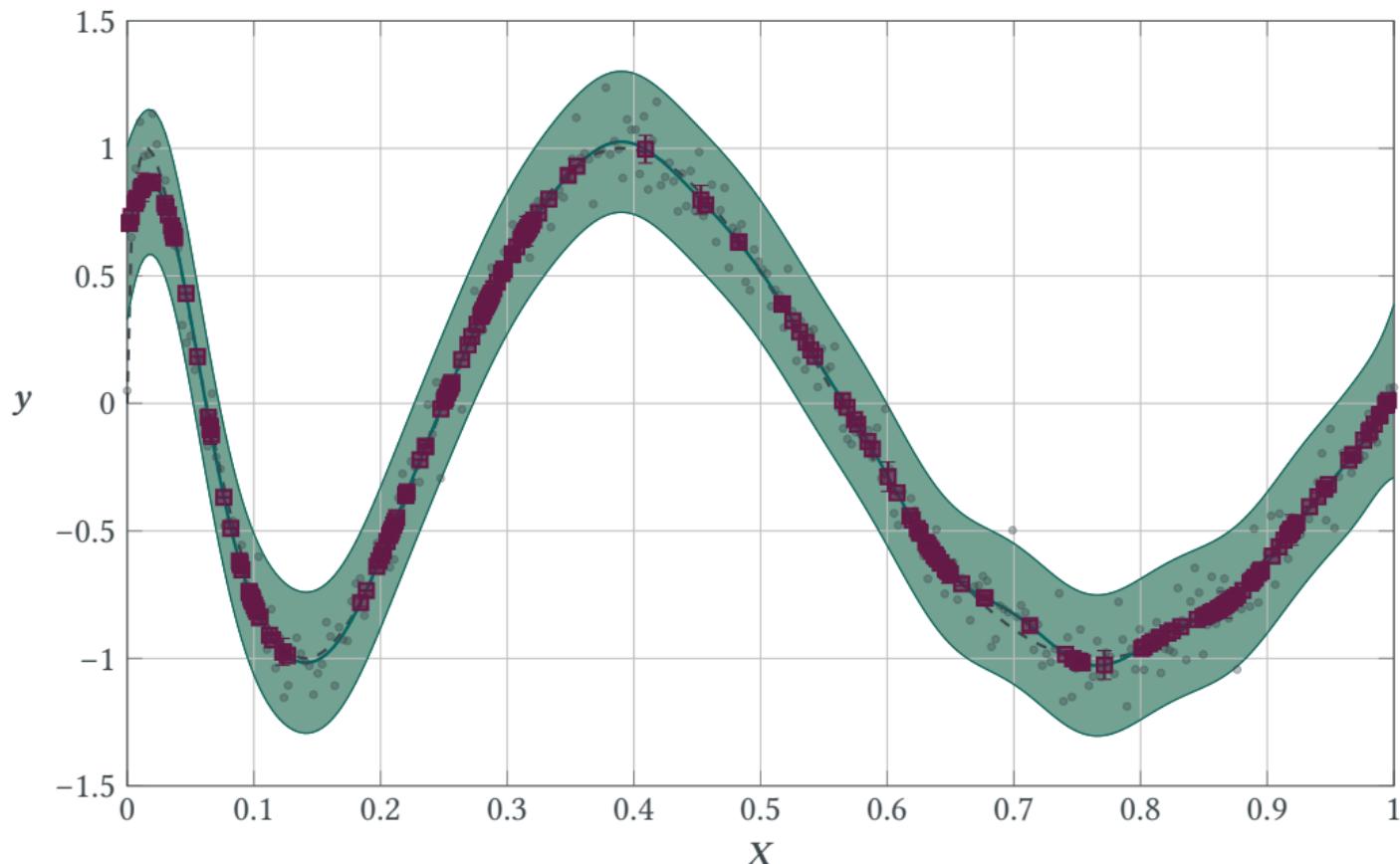
# SVGP Showcase



# SVGP Showcase



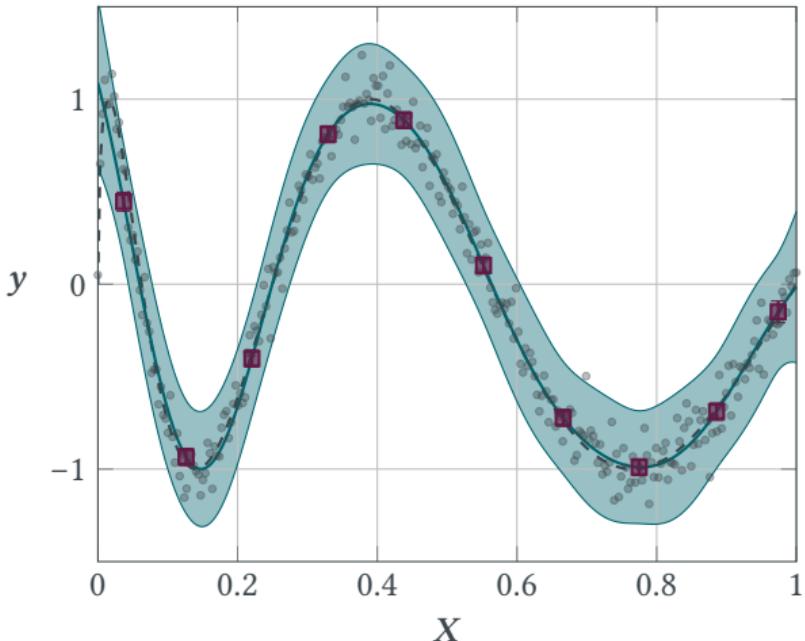
# SVGP Showcase



# Summary

We have made progress on seemingly fundamental restrictions of GPs:

- Scale to large data sets
- Scale to high dimensions (ish)
- Compact representations
- Acceptable speeds



Next time: Propagation of uncertainties?