



UNIVERSIDADE
FEDERAL DO CEARÁ

DEPARTAMENTO
DE COMPUTAÇÃO

The Art of Modeling with Gaussian Processes

César Lincoln Cavalcante Mattos

Department of Computer Science
Federal University of Ceará
Fortaleza, Ceará, Brazil
cesarlincoln@dc.ufc.br

2023

Table of Contents

① Part I

- Fundamentals of Bayesian inference
- Bayesian linear regression
- The Gaussian process

② Part II

- Learning the kernel and its hyperparameters
- Beyond Gaussian likelihood
- Sparse approximations

③ Part III

- GPs for Bayesian optimization
- Current trends on kernel design
- From GPLVM to Deep GPs
- Concluding remarks and additional topics

The Art of Gaussian Processes: Classical and Contemporary

Part of this course was taken from the NeurIPS 2021 tutorial jointly presented with Prof. Felipe Tobar and available at <https://neurips.cc/virtual/2021/tutorial/21890>.



César Lincoln C. Mattos is an Associate Professor at the Department of Computer Science, at Federal University of Ceará (UFC), Brazil. He is also an associate researcher at the Logics and Artificial Intelligence Group (LOGIA). He has research interests in (deep) probabilistic learning, approximate inference, system identification and anomaly detection.



Felipe Tobar is an Associate Professor at the Initiative for Data & Artificial Intelligence, Universidad de Chile, and the Coordinator of the Master of Data Science at the same University. He is also an Associate Researcher at the Center for Mathematical Modeling and the Advanced Centre for Electrical and Electronic Engineering. He teaches courses on Probability, Statistics and Machine Learning, and his research interests include time series, Bayesian inference, computational optimal transport and the societal impacts of machine learning.

Table of Contents

① Part I

Fundamentals of Bayesian inference

Bayesian linear regression

The Gaussian process

② Part II

Learning the kernel and its hyperparameters

Beyond Gaussian likelihood

Sparse approximations

③ Part III

GPs for Bayesian optimization

Current trends on kernel design

From GPLVM to Deep GPs

Concluding remarks and additional topics

Learning from data

- **Data:** partial observations of the world.
- **Model:** links observations and latent (unobserved) objects.
- **Uncertainty:** data is incomplete and noisy.
- **Learning:** estimate unknowns from known quantities.



The real
world is
complex Data is
fragmentary The
Model

(taken with permission from Thoughtfulnz on Twitter)

Probability

General < ----- > Particular

Population < ----- > Sample

Model < ----- > Data

Statistics

The statistical model

Notation

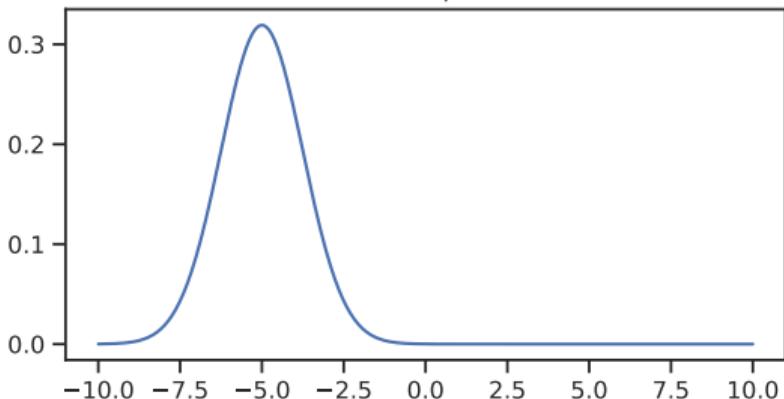
- **Sample space \mathcal{X} :** where the data lives.
- **Statistical model \mathcal{M} :** the set of probability distributions on \mathcal{X} indexed by a set \mathcal{T} :

$$\mathcal{M} = \{P_\theta \mid \theta \in \mathcal{T}\}. \quad (1)$$

- **Parameter θ and parameter space \mathcal{T} .**

Examples:

Gaussian with mean μ and variance σ^2



The statistical model

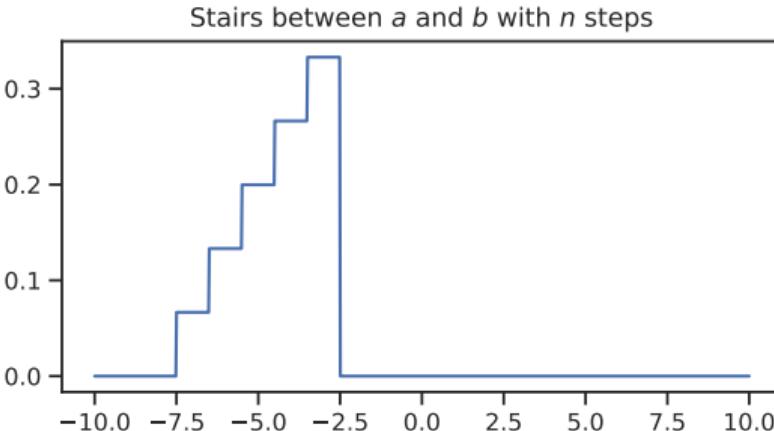
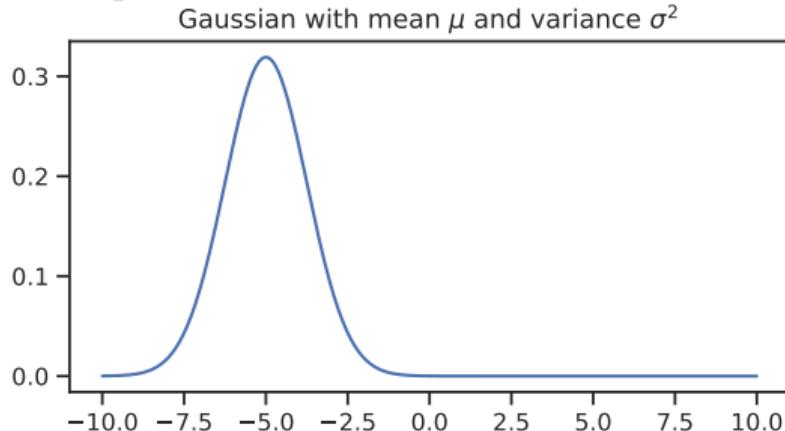
Notation

- **Sample space \mathcal{X} :** where the data lives.
- **Statistical model \mathcal{M} :** the set of probability distributions on \mathcal{X} indexed by a set \mathcal{T} :

$$\mathcal{M} = \{P_\theta \mid \theta \in \mathcal{T}\}. \quad (1)$$

- **Parameter θ and parameter space \mathcal{T} .**

Examples:



The statistical model

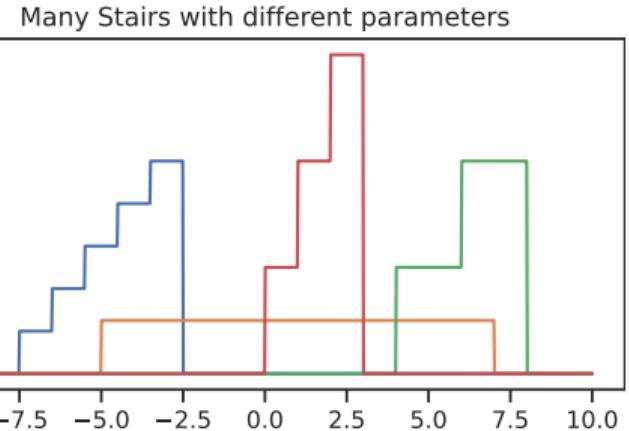
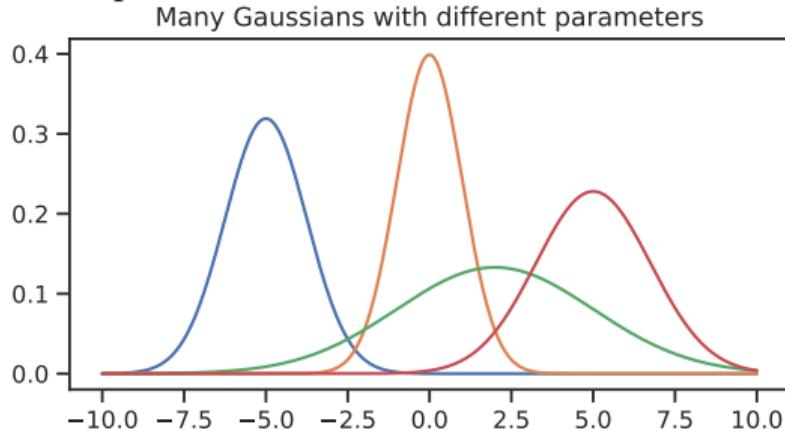
Notation

- **Sample space \mathcal{X} :** where the data lives.
- **Statistical model \mathcal{M} :** the set of probability distributions on \mathcal{X} indexed by a set \mathcal{T} :

$$\mathcal{M} = \{P_\theta \mid \theta \in \mathcal{T}\}. \quad (1)$$

- **Parameter θ and parameter space \mathcal{T} .**

Examples:

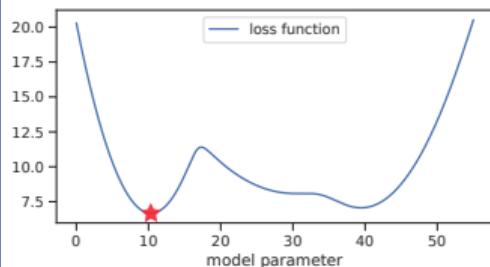


Finding the parameters

(and thus the model)

Alternative 1: best model

- i) Define the model space
- ii) Gather data
- iii) Give each model a (fit) score:

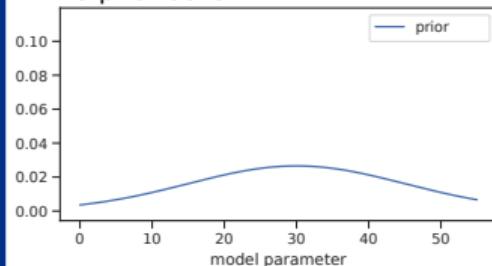


- iv) Choose the model with the best score ★

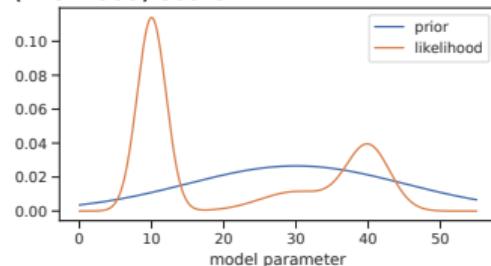
Predictions: "plug in" the chosen model into the predictor

Alternative 2: posterior belief

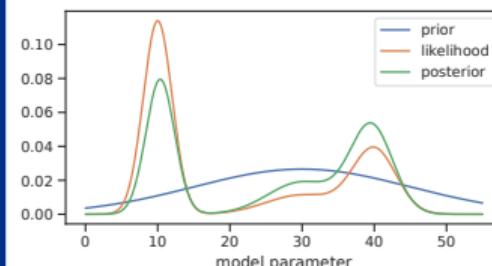
- i) Define the model space and equip it with a prior belief



- ii) Gather data and give each model a (likelihood) score



- iii) Update your prior belief into a posterior one



Predictions: predict with each model and average predictions wrt to posterior belief

Bayesian inference: a one-slide cheat sheet

Model specification: Define some parameters as random variables (θ) and some as deterministic values (λ):

$$\lambda \in \Lambda, \quad \theta \sim p_\lambda(\theta) \quad \Rightarrow \quad X|\theta, \lambda \sim p_\lambda(x|\theta). \quad (2)$$

After observing some data (\mathcal{D}):

- **Set λ (model comparison):** e.g., by maximizing the marginal likelihood

$$p_\lambda(\mathcal{D}) = \int p_\lambda(\mathcal{D}|\theta)p_\lambda(\theta)d\theta. \quad (3)$$

- **Bayesian update:** compute the posterior over θ

$$p_\lambda(\theta|\mathcal{D}) = \frac{p_\lambda(\mathcal{D}|\theta)p_\lambda(\theta)}{p_\lambda(\mathcal{D})}. \quad (4)$$

- **Prediction:** performed via posterior averages

$$p_\lambda(x|\mathcal{D}) = \int p_\lambda(x|\theta)p_\lambda(\theta|\mathcal{D})d\theta. \quad (5)$$

Important: The data is *summarized* in the posterior $p_\lambda(\theta|\mathcal{D})$ and is used for prediction.

Table of Contents

① Part I

Fundamentals of Bayesian inference

Bayesian linear regression

The Gaussian process

② Part II

Learning the kernel and its hyperparameters

Beyond Gaussian likelihood

Sparse approximations

③ Part III

GPs for Bayesian optimization

Current trends on kernel design

From GPLVM to Deep GPs

Concluding remarks and additional topics

Bayesian linear regression

- Consider a linear regression model with Gaussian noise:

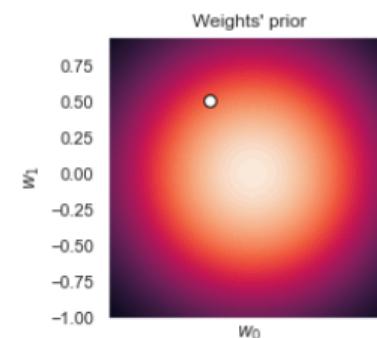
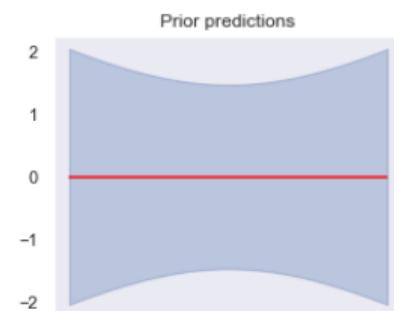
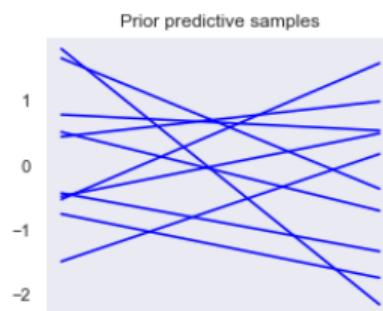
$$y = \mathbf{w}^\top \mathbf{x} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2).$$

- The likelihood function is given by:

$$p(y|\mathbf{x}, \mathbf{w}, \sigma^2) = \mathcal{N}(y|\mathbf{w}^\top \mathbf{x}, \sigma^2).$$

- The unknown parameter vector \mathbf{w} has a prior distribution, e.g., a multivariate Gaussian:

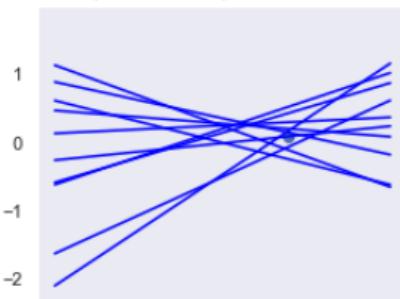
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0).$$



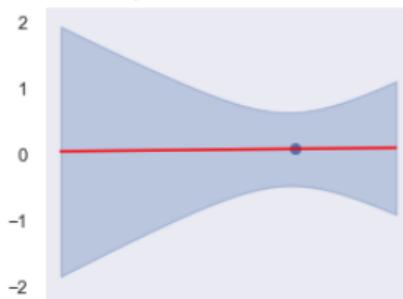
Bayesian linear regression

- After observing $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, we calculate the posterior $p(\mathbf{w}|\mathcal{D})$ and perform prediction.

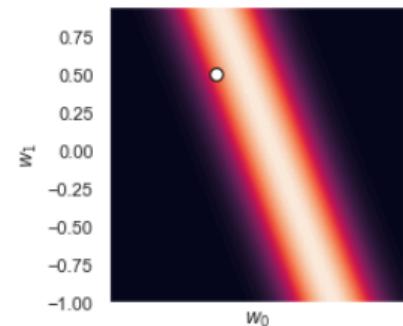
Posterior predictive samples after 1 observations



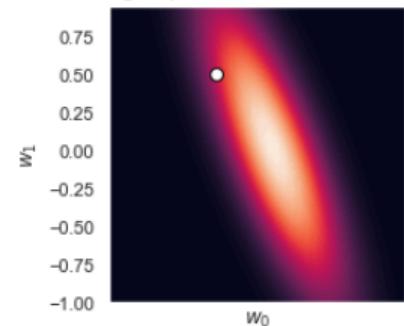
Posterior predictions after 1 observations



Likelihood after 1 observations



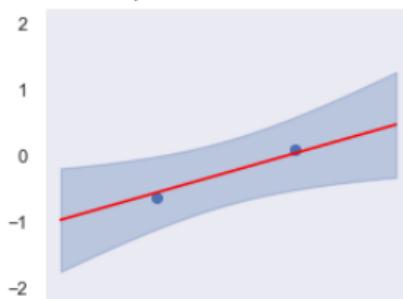
Weights' posterior after 1 observations



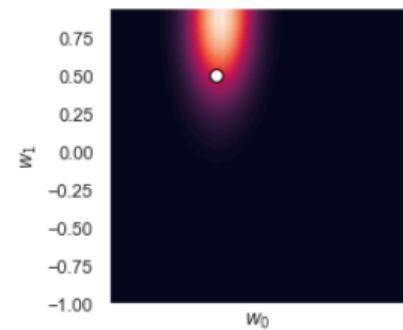
Posterior predictive samples after 2 observations



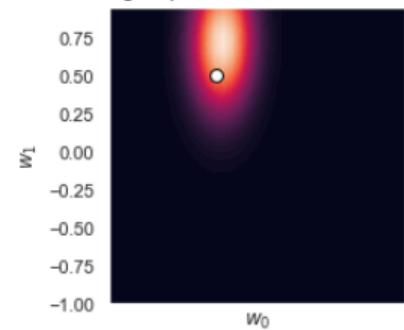
Posterior predictions after 2 observations



Likelihood after 2 observations



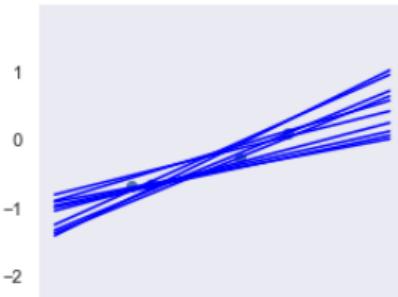
Weights' posterior after 2 observations



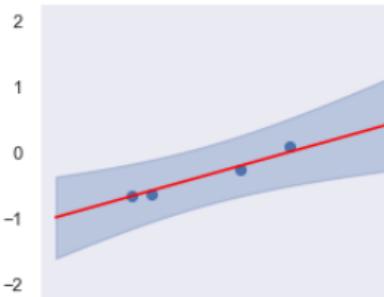
Bayesian linear regression

- After observing $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, we calculate the posterior $p(\mathbf{w}|\mathcal{D})$ and perform prediction.

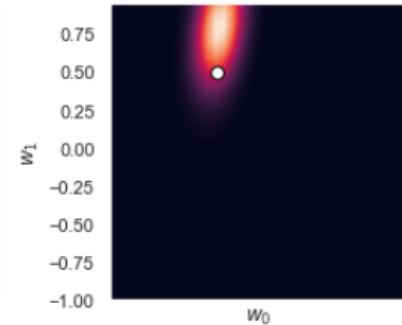
Posterior predictive samples after 4 observations



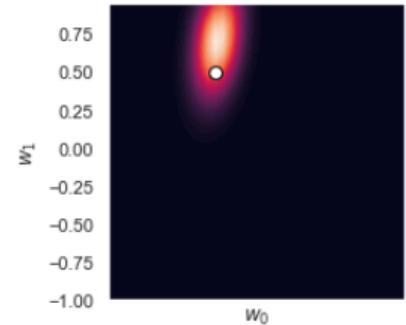
Posterior predictions after 4 observations



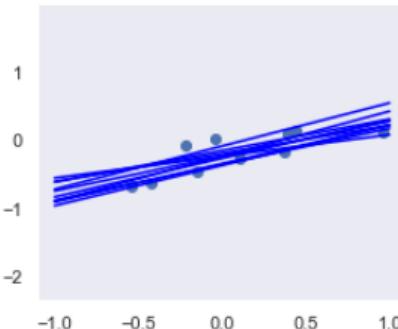
Likelihood after 4 observations



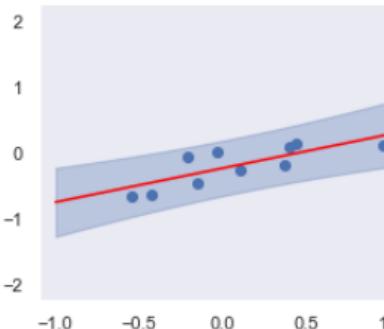
Weights' posterior after 4 observations



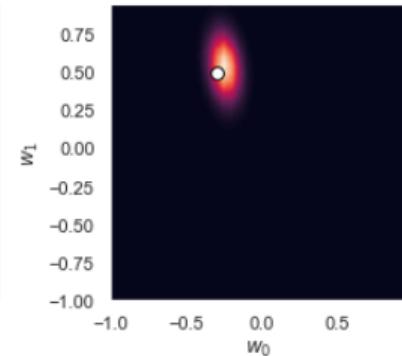
Posterior predictive samples after 10 observations



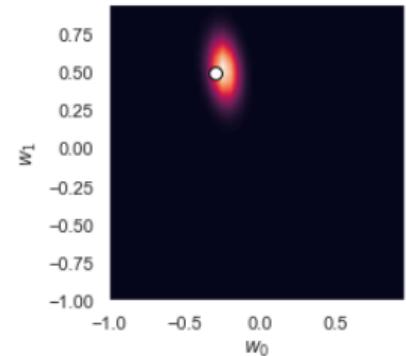
Posterior predictions after 10 observations



Likelihood after 10 observations



Weights' posterior after 10 observations

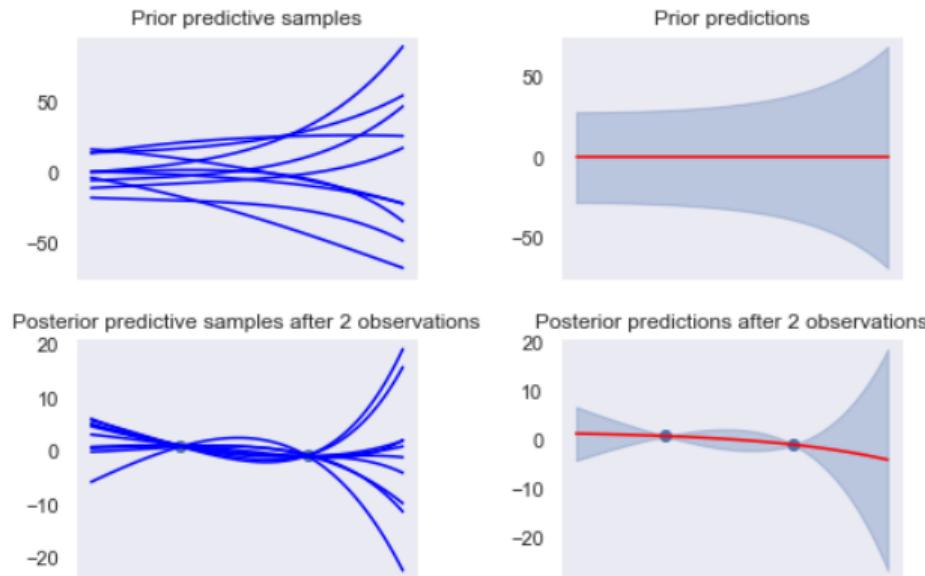


Bayesian polynomial regression

- Consider a 5-order polynomial regression model with Gaussian noise:

$$y = \mathbf{w}^\top \phi(x) + \epsilon, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0), \quad \epsilon \sim \mathcal{N}(0, \sigma^2),$$

where $\phi(x) = [1, x, x^2, x^3, x^4, x^5]^\top$.



Bayesian polynomial regression

- Consider a 5-order polynomial regression model with Gaussian noise:

$$y = \mathbf{w}^\top \phi(x) + \epsilon, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0), \quad \epsilon \sim \mathcal{N}(0, \sigma^2),$$

where $\phi(x) = [1, x, x^2, x^3, x^4, x^5]^\top$.

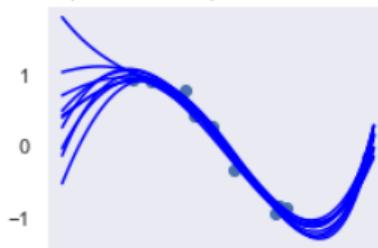
Posterior predictive samples after 6 observations



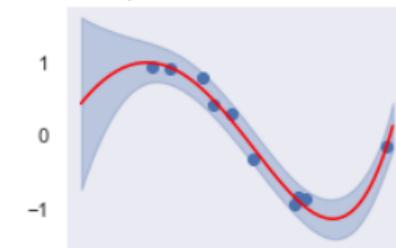
Posterior predictions after 6 observations



Posterior predictive samples after 10 observations



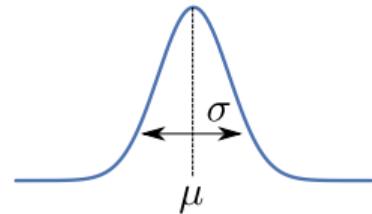
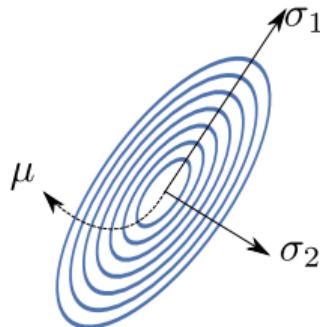
Posterior predictions after 10 observations



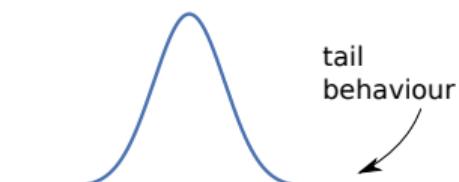
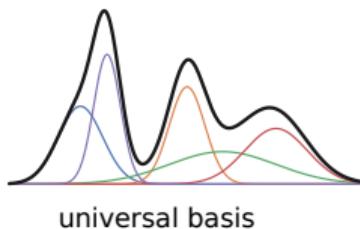
Why Gaussians?

Many useful properties:

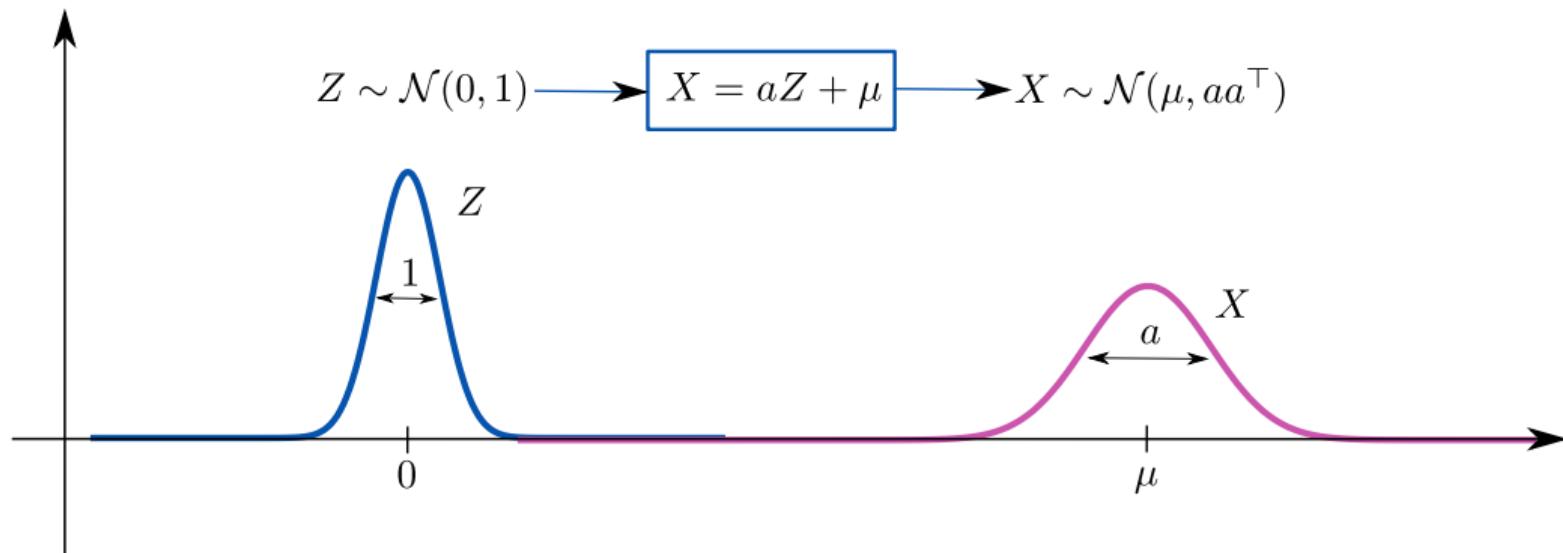
- central limit theorem;
- conjugacy;
- infinite support;
- infinitely differentiable;
- maximum entropy (fixed variance);
- convex loss (linear regression);
- conditional expectation is linear;
- equivalence to least squares;
- uncorrelatedness \Leftrightarrow independence;
- universal basis.



$$X \sim \mathcal{N}(\mu, \sigma^2) \Leftrightarrow p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

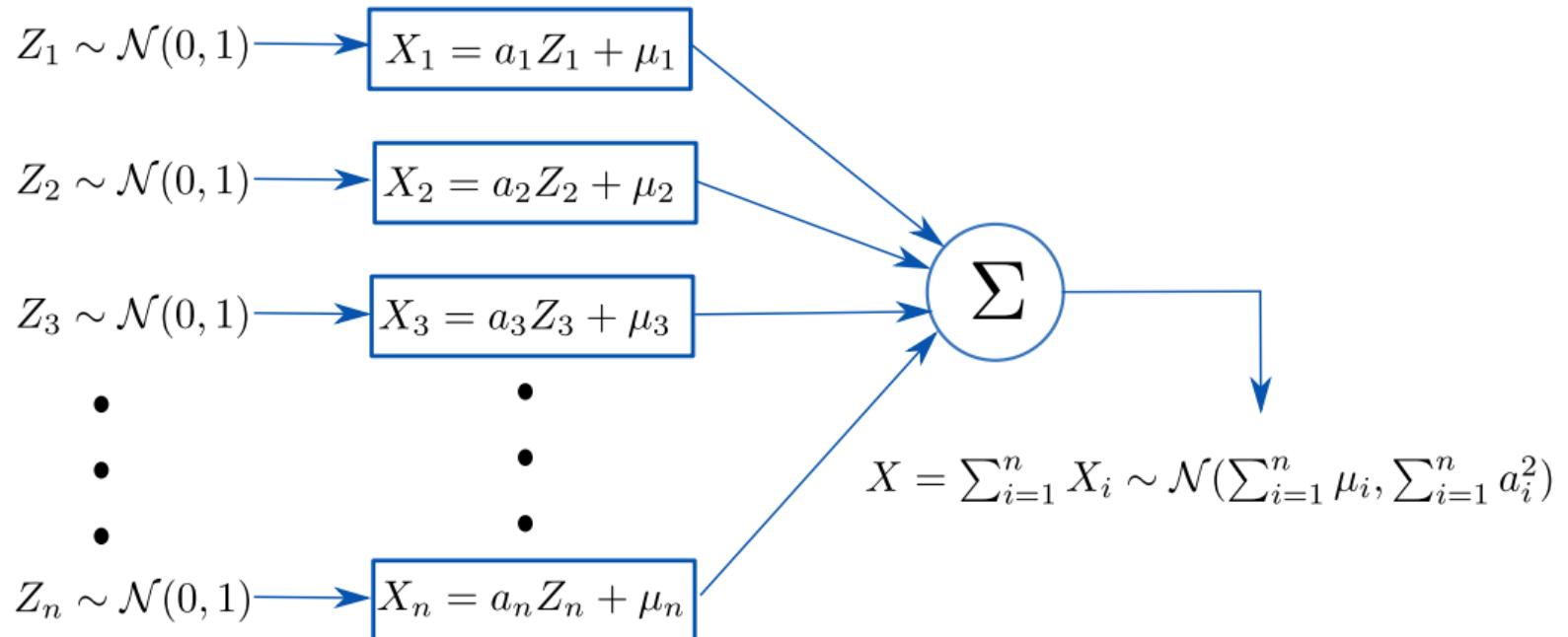


Linear transformation of Gaussians



- Closed under linear transformations.
- Any Gaussian RV can be produced from a linear transformation of $\mathcal{N}(0, 1)$.

Linear combinations of Gaussians



(we have assumed the RVs to be independent but the result holds for the correlated case)

Summary of convenient Gaussian properties

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}.$$

Marginalization

The observation of a larger collection of variables does not affect the distribution of smaller subsets:

$$x_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \text{ and } x_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}).$$

Conditioning

Gaussian conditioning results in another Gaussian:

$$p(\mathbf{x}_1 | \mathbf{x}_2 = \mathbf{z}) = \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{z} - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

Linearity

Linear combination of Gaussians results in another Gaussian:

$$p(a\mathbf{x}_1 + b\mathbf{x}_2) = \mathcal{N}(a\boldsymbol{\mu}_1 + b\boldsymbol{\mu}_2, a^2\boldsymbol{\Sigma}_{11} + b^2\boldsymbol{\Sigma}_{22}).$$

Bayesian linear regression

Math details

- Considering inputs $\mathbf{X} \in \mathbb{R}^{N \times D}$ and observed outputs $\mathbf{y} \in \mathbb{R}^N$:

$$\begin{aligned}y_i &= \mathbf{w}^\top \mathbf{x}_i + \epsilon_i, \quad p(\epsilon_i) = \mathcal{N}(\epsilon_i | 0, \sigma^2), \quad i = 1, \dots, N, \\p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) &= \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}), \\p(\mathbf{w}) &= \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0).\end{aligned}$$

- The marginal likelihood is tractable:

$$\begin{aligned}p(\mathbf{y} | \mathbf{X}) &= \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) p(\mathbf{w}) d\mathbf{w} \\&= \int \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) d\mathbf{w} \\&= \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2 \mathbf{I}).\end{aligned}$$

- The covariance between \mathbf{w} and \mathbf{y} is calculated by

$$\begin{aligned}\text{cov}[\mathbf{w}, \mathbf{y}] &= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}] \\&= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w}] = \mathbf{X}\text{cov}[\mathbf{w}, \mathbf{w}] = \mathbf{X}\mathbf{S}_0.\end{aligned}$$

Bayesian linear regression

Math details

- Considering inputs $\mathbf{X} \in \mathbb{R}^{N \times D}$ and observed outputs $\mathbf{y} \in \mathbb{R}^N$:

$$\begin{aligned}y_i &= \mathbf{w}^\top \mathbf{x}_i + \epsilon_i, \quad p(\epsilon_i) = \mathcal{N}(\epsilon_i | 0, \sigma^2), \quad i = 1, \dots, N, \\p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) &= \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}), \\p(\mathbf{w}) &= \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0).\end{aligned}$$

- The marginal likelihood is tractable:

$$\begin{aligned}p(\mathbf{y} | \mathbf{X}) &= \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) p(\mathbf{w}) d\mathbf{w} \\&= \int \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) d\mathbf{w} \\&= \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2 \mathbf{I}).\end{aligned}$$

- The covariance between \mathbf{w} and \mathbf{y} is calculated by

$$\begin{aligned}\text{cov}[\mathbf{w}, \mathbf{y}] &= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}] \\&= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w}] = \mathbf{X}\text{cov}[\mathbf{w}, \mathbf{w}] = \mathbf{X}\mathbf{S}_0.\end{aligned}$$

Bayesian linear regression

Math details

- Considering inputs $\mathbf{X} \in \mathbb{R}^{N \times D}$ and observed outputs $\mathbf{y} \in \mathbb{R}^N$:

$$\begin{aligned}y_i &= \mathbf{w}^\top \mathbf{x}_i + \epsilon_i, \quad p(\epsilon_i) = \mathcal{N}(\epsilon_i | 0, \sigma^2), \quad i = 1, \dots, N, \\p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) &= \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}), \\p(\mathbf{w}) &= \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0).\end{aligned}$$

- The marginal likelihood is tractable:

$$\begin{aligned}p(\mathbf{y} | \mathbf{X}) &= \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) p(\mathbf{w}) d\mathbf{w} \\&= \int \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) d\mathbf{w} \\&= \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2 \mathbf{I}).\end{aligned}$$

- The covariance between \mathbf{w} and \mathbf{y} is calculated by

$$\begin{aligned}\text{cov}[\mathbf{w}, \mathbf{y}] &= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}] \\&= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w}] = \mathbf{X}\text{cov}[\mathbf{w}, \mathbf{w}] = \mathbf{X}\mathbf{S}_0.\end{aligned}$$

Bayesian linear regression

Math details

- The joint distribution $p(\mathbf{w}, \mathbf{y})$ can be readily built:

$$p\left(\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \middle| \mathbf{X}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_0 \\ \mathbf{Xm}_0 \end{bmatrix}, \begin{bmatrix} \mathbf{S}_0 & \mathbf{S}_0 \mathbf{X}^\top \\ \mathbf{X} \mathbf{S}_0 & \mathbf{X} \mathbf{S}_0 \mathbf{X}^\top + \sigma^2 \mathbf{I} \end{bmatrix}\right)$$

- Using the Gaussian conditioning property, the posterior over \mathbf{w} after observing $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ is

$$\begin{aligned} p(\mathbf{w} | \mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}, \boldsymbol{\Sigma}), \\ \boldsymbol{\mu} &= \mathbf{m}_0 + (\mathbf{S}_0 \mathbf{X}^\top \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{S}_0 \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \mathbf{m}_0), \\ \boldsymbol{\Sigma} &= \mathbf{S}_0 - (\mathbf{S}_0 \mathbf{X}^\top \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{S}_0 \mathbf{X}^\top \mathbf{X} \mathbf{S}_0. \end{aligned}$$

- Predictions for new inputs \mathbf{X}_* are given by

$$\begin{aligned} p(\mathbf{y}_* | \mathbf{X}_*, \mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \int p(\mathbf{y}_* | \mathbf{w}, \sigma^2) p(\mathbf{w} | \mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) d\mathbf{w} \\ &= \int \mathcal{N}(\mathbf{y}_* | \mathbf{X}_* \mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}_* | \mathbf{X}_* \boldsymbol{\mu}, \mathbf{X}_* \boldsymbol{\Sigma} \mathbf{X}_*^\top + \sigma^2 \mathbf{I}). \end{aligned}$$

Bayesian linear regression

Math details

- The joint distribution $p(\mathbf{w}, \mathbf{y})$ can be readily built:

$$p\left(\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \middle| \mathbf{X}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_0 \\ \mathbf{X}\mathbf{m}_0 \end{bmatrix}, \begin{bmatrix} \mathbf{S}_0 & \mathbf{S}_0\mathbf{X}^\top \\ \mathbf{X}\mathbf{S}_0 & \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2\mathbf{I} \end{bmatrix}\right)$$

- Using the Gaussian conditioning property, the posterior over \mathbf{w} after observing $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ is

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \\ \boldsymbol{\mu} &= \mathbf{m}_0 + (\mathbf{S}_0\mathbf{X}^\top\mathbf{X} + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\mathbf{m}_0), \\ \boldsymbol{\Sigma} &= \mathbf{S}_0 - (\mathbf{S}_0\mathbf{X}^\top\mathbf{X} + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\mathbf{X}^\top\mathbf{X}\mathbf{S}_0. \end{aligned}$$

- Predictions for new inputs \mathbf{X}_* are given by

$$\begin{aligned} p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \int p(\mathbf{y}_*|\mathbf{w}, \sigma^2)p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2)d\mathbf{w} \\ &= \int \mathcal{N}(\mathbf{y}_*|\mathbf{X}_*\mathbf{w}, \sigma^2\mathbf{I})\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}_*|\mathbf{X}_*\boldsymbol{\mu}, \mathbf{X}_*\boldsymbol{\Sigma}\mathbf{X}_*^\top + \sigma^2\mathbf{I}). \end{aligned}$$

Bayesian linear regression

Math details

- The joint distribution $p(\mathbf{w}, \mathbf{y})$ can be readily built:

$$p\left(\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \middle| \mathbf{X}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_0 \\ \mathbf{X}\mathbf{m}_0 \end{bmatrix}, \begin{bmatrix} \mathbf{S}_0 & \mathbf{S}_0\mathbf{X}^\top \\ \mathbf{X}\mathbf{S}_0 & \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2\mathbf{I} \end{bmatrix}\right)$$

- Using the Gaussian conditioning property, the posterior over \mathbf{w} after observing $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ is

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \\ \boldsymbol{\mu} &= \mathbf{m}_0 + (\mathbf{S}_0\mathbf{X}^\top\mathbf{X} + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\mathbf{m}_0), \\ \boldsymbol{\Sigma} &= \mathbf{S}_0 - (\mathbf{S}_0\mathbf{X}^\top\mathbf{X} + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\mathbf{X}^\top\mathbf{X}\mathbf{S}_0. \end{aligned}$$

- Predictions for new inputs \mathbf{X}_* are given by

$$\begin{aligned} p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \int p(\mathbf{y}_*|\mathbf{w}, \sigma^2)p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2)d\mathbf{w} \\ &= \int \mathcal{N}(\mathbf{y}_*|\mathbf{X}_*\mathbf{w}, \sigma^2\mathbf{I})\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}_*|\mathbf{X}_*\boldsymbol{\mu}, \mathbf{X}_*\boldsymbol{\Sigma}\mathbf{X}_*^\top + \sigma^2\mathbf{I}). \end{aligned}$$

Table of Contents

① Part I

Fundamentals of Bayesian inference

Bayesian linear regression

The Gaussian process

② Part II

Learning the kernel and its hyperparameters

Beyond Gaussian likelihood

Sparse approximations

③ Part III

GPs for Bayesian optimization

Current trends on kernel design

From GPLVM to Deep GPs

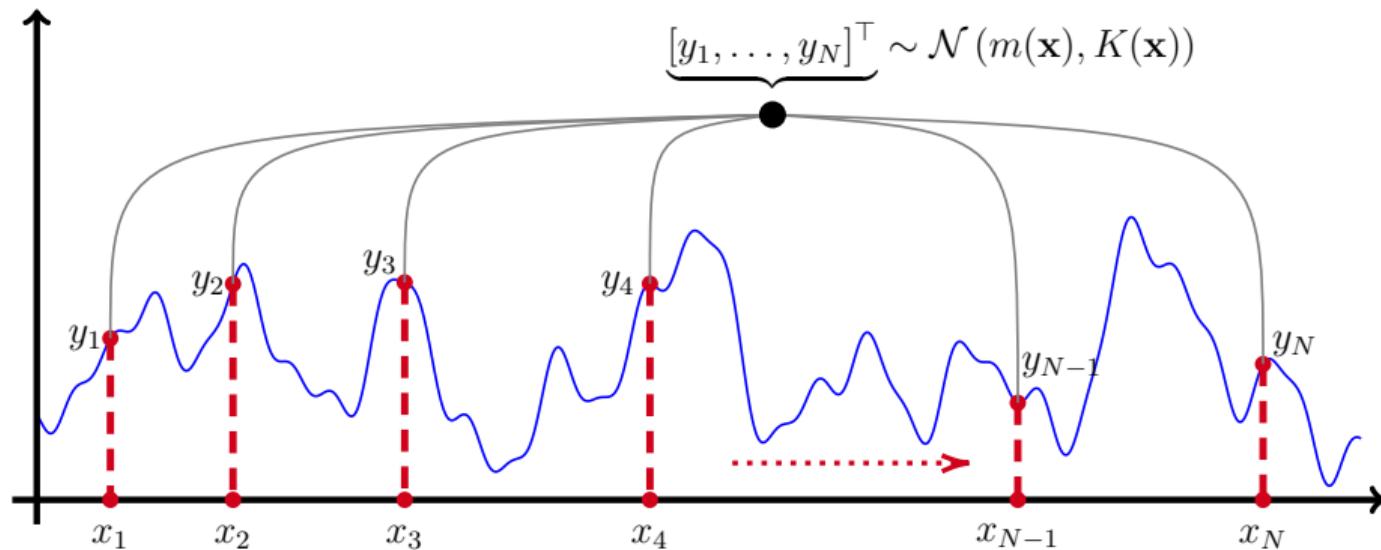
Concluding remarks and additional topics

What do we seek?

- Avoid a finite parametric formulation that acts as a **bottleneck** between data and predictions.
- Proper **uncertainty handling** and predictive distributions.
- Balance between **data fit** and **regularization**.
- (Few) hyperparameters that can be **directly** adjusted.
- Solid theoretical framework with **useful properties** from Gaussians.
- **Gaussian process**: a Bayesian nonparametric modeling framework.

The Gaussian process

Definition: A GP is a stochastic process such that any finite collection of values follows a multivariate Gaussian distribution.



The Gaussian process

- Notation for N observations:

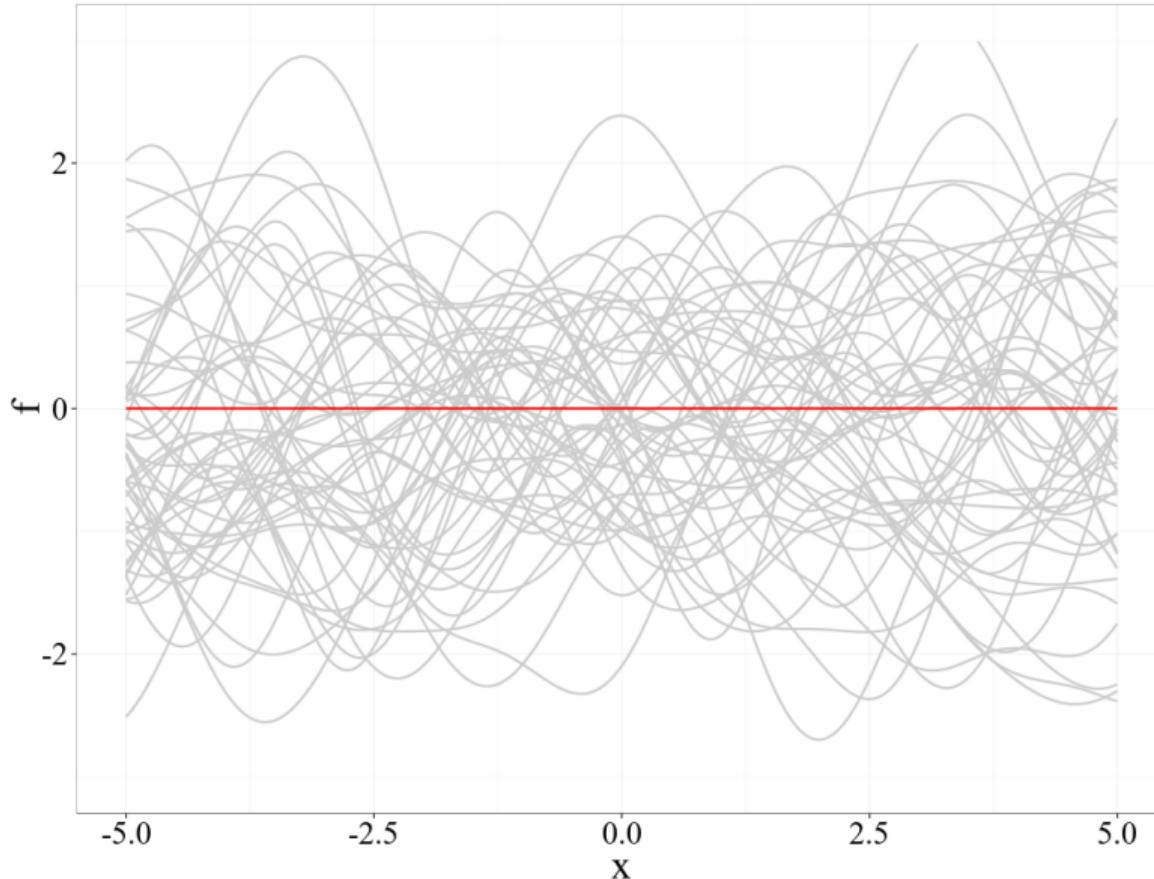
$$f \sim \mathcal{GP}(\boldsymbol{\mu}, \mathbf{K}) \iff f(\mathbf{x}) \sim \mathcal{N}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')),$$

$$p(\mathbf{f}|\boldsymbol{\mu}, \mathbf{K}) = \frac{1}{(2\pi)^{\frac{N}{2}} |\mathbf{K}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{f} - \boldsymbol{\mu})^\top \mathbf{K}^{-1}(\mathbf{f} - \boldsymbol{\mu})\right).$$

Important

- The GP is a distribution over **functions**.
- The likelihood $p(\mathbf{y}|\mathbf{f})$ relates the noisy observations \mathbf{y} and the latent function values \mathbf{f} .
- The GP is entirely defined by its mean ($m(\mathbf{x})$) and covariance ($K(\mathbf{x}, \mathbf{x}')$) functions.
- We usually choose a zero mean prior ($\boldsymbol{\mu} = \mathbf{0}$).
- The GP allows us to analyze functions (**infinite objects**) from finite distributions (**multivariate Gaussians**).
- The function f is a single (infinite dimensional) GP sample, and, accordingly, the vector \mathbf{f} is a single (finite dimensional) multivariate Gaussian sample.

Samples from the GP prior



Nonlinear regression

- Consider a dataset

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) |_{i=1}^N\} = (\mathbf{X}, \mathbf{y}),$$

where $\mathbf{x}_i \in \mathbb{R}^D$ are the inputs, $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $\mathbf{y} \in \mathbb{R}^N$ are the observed outputs.

- A general nonlinear regression task is expressed as

$$y_i = f(\mathbf{x}_i) + \epsilon_i,$$

where $\epsilon_i \in \mathbb{R}$ is the observation noise.

- The values $f_i = f(\mathbf{x}_i)$ are not directly observed (they are *latent*) and the true $f(\cdot)$ is unknown.

Nonlinear regression with GPs

- Choose a GP prior for the unknown function:

$$f \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}), \\ \mathbf{f} = f(\mathbf{X}) \sim \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K}),$$

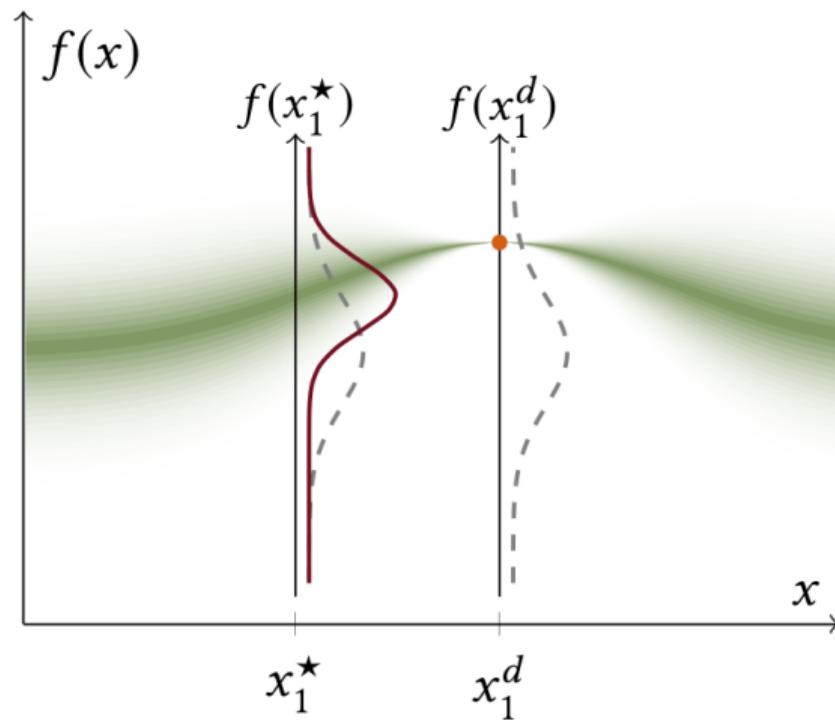
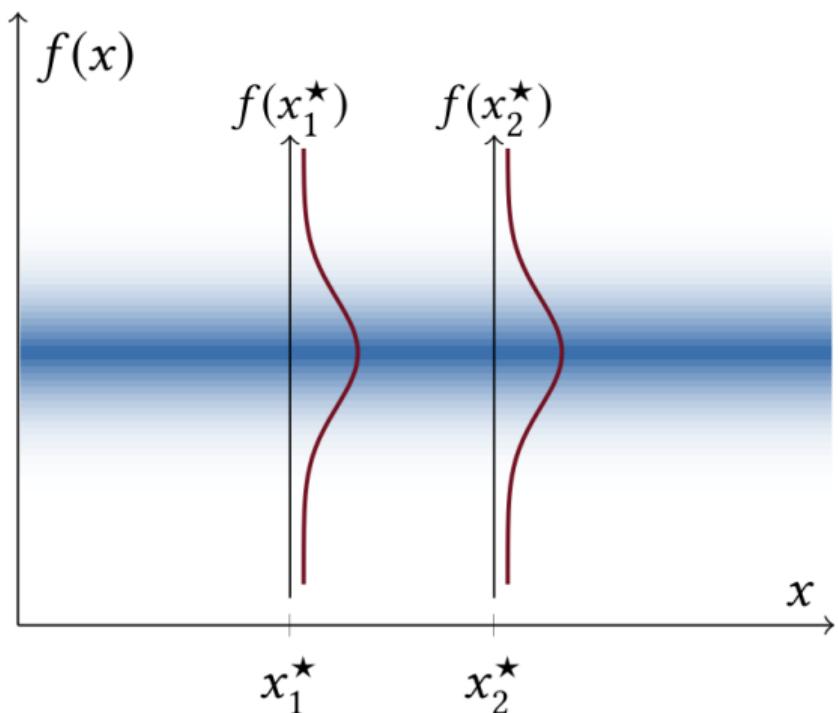
where $\mathbf{K} \in \mathbb{R}^{N \times N}$, $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, is the covariance matrix, obtained from the covariance, or *kernel* function, $k(\cdot, \cdot)$.

- Note that the output behavior is entirely explained by the covariance matrix \mathbf{K} , obtained from relations between the inputs.
- A common choice is the *squared exponential* (or *Radial Basis Function* - RBF):

$$K(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_{d=1}^D \frac{(x_{id} - x_{jd})^2}{l_d^2}\right).$$

- The vector of hyperparameters $\boldsymbol{\theta} = [\sigma_f^2, l_1^2, \dots, l_D^2]^\top$ characterizes the model covariance.

Nonlinear regression with GPs



Nonlinear regression with GPs

- If the observation noise is Gaussian, i.e., $\epsilon_i \sim \mathcal{N}(0, \sigma_y^2)$, all the expressions of interest are analytical.

Likelihood

Since the noise is Gaussian, the likelihood is also Gaussian:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}),$$

where $\mathbf{I} \in \mathbb{R}^{N \times N}$ is the N -th order identity matrix.

- Important: Different noise considerations lead to distinct likelihood functions and break the tractability (more on that later!).

Nonlinear regression with GPs

- If the observation noise is Gaussian, i.e., $\epsilon_i \sim \mathcal{N}(0, \sigma_y^2)$, all the expressions of interest are analytical.

Likelihood

Since the noise is Gaussian, the likelihood is also Gaussian:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}),$$

where $\mathbf{I} \in \mathbb{R}^{N \times N}$ is the N -th order identity matrix.

- **Important:** Different noise considerations lead to distinct likelihood functions and break the tractability (more on that later!).

Nonlinear regression with GPs

Marginal likelihood

The marginal likelihood is calculated by integrating \mathbf{f} :

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X})d\mathbf{f} = \int \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I})\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})d\mathbf{f}, \\ &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K} + \sigma_y^2 \mathbf{I}). \end{aligned}$$

Posteriori

The posterior over \mathbf{f} , after observing \mathbf{y} , is also Gaussian:

$$\begin{aligned} \underbrace{p(\mathbf{f}|\mathbf{y}, \mathbf{X})}_{\text{posterior}} &= \frac{\overbrace{p(\mathbf{y}|\mathbf{f})}^{\text{likelihood}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood}}} = \frac{\mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I})\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})}{\mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K} + \sigma_y^2 \mathbf{I})}, \\ &= \mathcal{N}(\mathbf{f}|\mathbf{K}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1}\mathbf{y}, \mathbf{K} - \mathbf{K}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1}\mathbf{K}). \end{aligned}$$

Nonlinear regression with GPs

Marginal likelihood

The marginal likelihood is calculated by integrating \mathbf{f} :

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X})d\mathbf{f} = \int \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I})\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})d\mathbf{f}, \\ &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K} + \sigma_y^2 \mathbf{I}). \end{aligned}$$

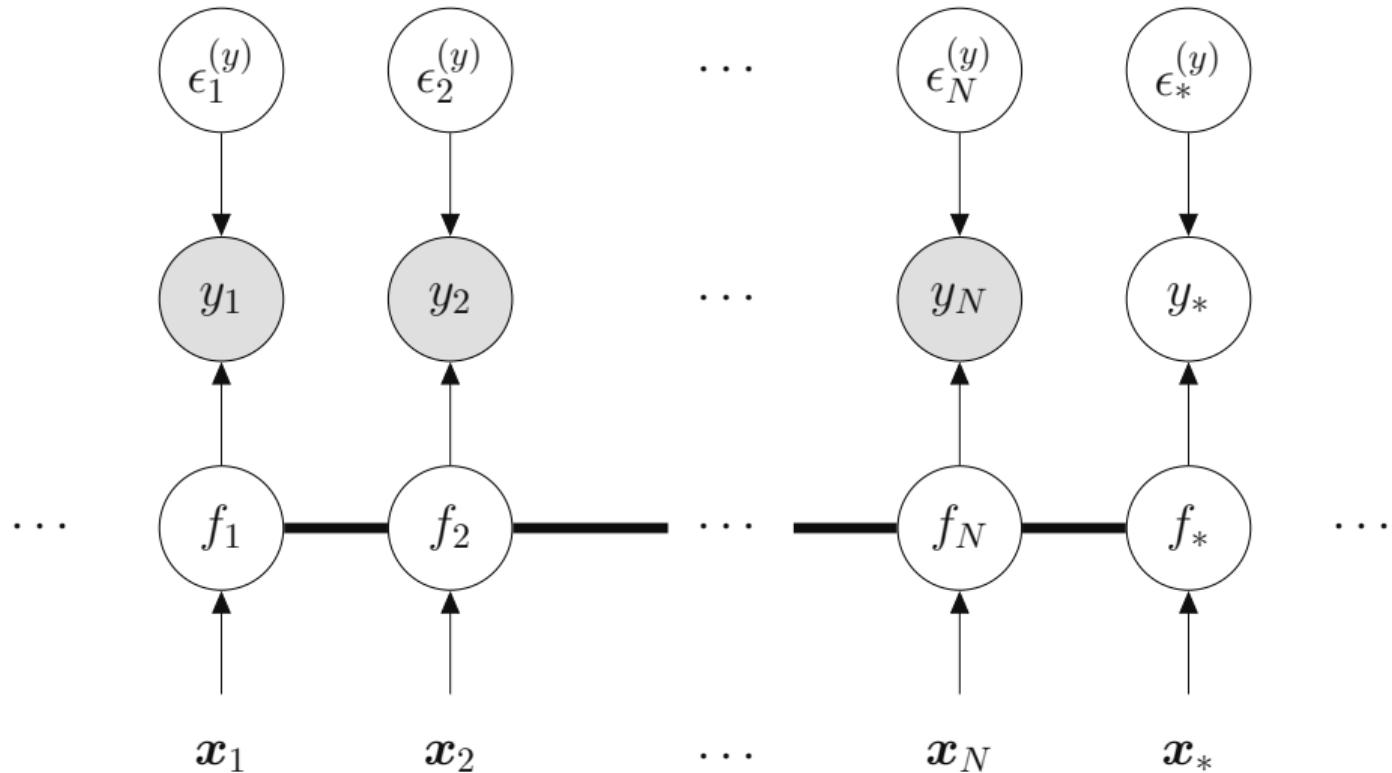
Posteriori

The posterior over \mathbf{f} , after observing \mathbf{y} , is also Gaussian:

$$\begin{aligned} \underbrace{p(\mathbf{f}|\mathbf{y}, \mathbf{X})}_{\text{posterior}} &= \frac{\overbrace{p(\mathbf{y}|\mathbf{f})}^{\text{likelihood}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood}}} = \frac{\mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I})\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})}{\mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K} + \sigma_y^2 \mathbf{I})}, \\ &= \mathcal{N}(\mathbf{f}|\mathbf{K}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1}\mathbf{y}, \mathbf{K} - \mathbf{K}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1}\mathbf{K}). \end{aligned}$$

Nonlinear regression with GPs

Graphical model



Predictions with GPs

- Inference for f_* (latent space) given \mathbf{x}_* is obtained by using the posterior $p(\mathbf{f}|\mathbf{y}, \mathbf{X})$:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f} | \mathbf{y}, \mathbf{X}) d\mathbf{f}, \\ &= \int \mathcal{N}(f_* | \mathbf{k}_{*f} \mathbf{K}_f^{-1} \mathbf{f}, k_{**} - \mathbf{k}_{*f} \mathbf{K}_f^{-1} \mathbf{k}_{f*}) p(\mathbf{f} | \mathbf{y}, \mathbf{X}) d\mathbf{f}, \\ &= \mathcal{N}(f_* | \mu_*, \sigma_*^2), \\ \mu_* &= \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \\ \sigma_*^2 &= k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}, \\ \text{where } \mathbf{k}_{f*} &= [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)]^\top, \quad k_{**} = k(\mathbf{x}_*, \mathbf{x}_*). \end{aligned}$$

- Inference for y_* (observation space) is obtained by marginalizing f_* :

$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* | f_*) p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) df \\ &= \mathcal{N}(y_* | \mu_*, \sigma_*^2 + \sigma_y^2). \end{aligned}$$

Predictions with GPs

- Inference for f_* (latent space) given \mathbf{x}_* is obtained by using the posterior $p(\mathbf{f}|\mathbf{y}, \mathbf{X})$:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f} | \mathbf{y}, \mathbf{X}) d\mathbf{f}, \\ &= \int \mathcal{N}(f_* | \mathbf{k}_{*f} \mathbf{K}_f^{-1} \mathbf{f}, k_{**} - \mathbf{k}_{*f} \mathbf{K}_f^{-1} \mathbf{k}_{f*}) p(\mathbf{f} | \mathbf{y}, \mathbf{X}) d\mathbf{f}, \\ &= \mathcal{N}(f_* | \mu_*, \sigma_*^2), \\ \mu_* &= \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \\ \sigma_*^2 &= k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}, \\ \text{where } \mathbf{k}_{f*} &= [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)]^\top, \quad k_{**} = k(\mathbf{x}_*, \mathbf{x}_*). \end{aligned}$$

- Inference for y_* (observation space) is obtained by marginalizing f_* :

$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* | f_*) p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) d\mathbf{f} \\ &= \mathcal{N}(y_* | \mu_*, \sigma_*^2 + \sigma_y^2). \end{aligned}$$

Predictions with GPs

Alternative derivation

- Inference for f_* given \mathbf{x}_* may also be obtained from the joint distribution:

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_y^2 \mathbf{I} & \mathbf{k}_{f*} \\ \mathbf{k}_{*f} & k_{**} \end{bmatrix}\right),$$

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(f_* | \mu_*, \sigma_*^2),$$

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y_* | \mu_*, \sigma_*^2 + \sigma_y^2),$$

$$\mu_* = \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\sigma_*^2 = k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*},$$

where $\mathbf{k}_{f*} = [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)]^\top$, $k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$.

Stable implementation via Cholesky decomposition

$$\mathbf{K} + \sigma_y^2 \mathbf{I} = \mathbf{L} \mathbf{L}^\top, \quad (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} = \mathbf{L}^{-\top} \mathbf{L}^{-1},$$

$$\mu_* = (\mathbf{L}^{-1} \mathbf{k}_{f*})^\top \mathbf{L}^{-1} \mathbf{y},$$

$$\sigma_*^2 = k_{**} - (\mathbf{L}^{-1} \mathbf{k}_{f*})^\top \mathbf{L}^{-1} \mathbf{k}_{f*}.$$

Predictions with GPs

Alternative derivation

- Inference for f_* given \mathbf{x}_* may also be obtained from the joint distribution:

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_y^2 \mathbf{I} & \mathbf{k}_{f*} \\ \mathbf{k}_{*f} & k_{**} \end{bmatrix}\right),$$

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(f_* | \mu_*, \sigma_*^2),$$

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y_* | \mu_*, \sigma_*^2 + \sigma_y^2),$$

$$\mu_* = \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\sigma_*^2 = k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*},$$

where $\mathbf{k}_{f*} = [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)]^\top$, $k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$.

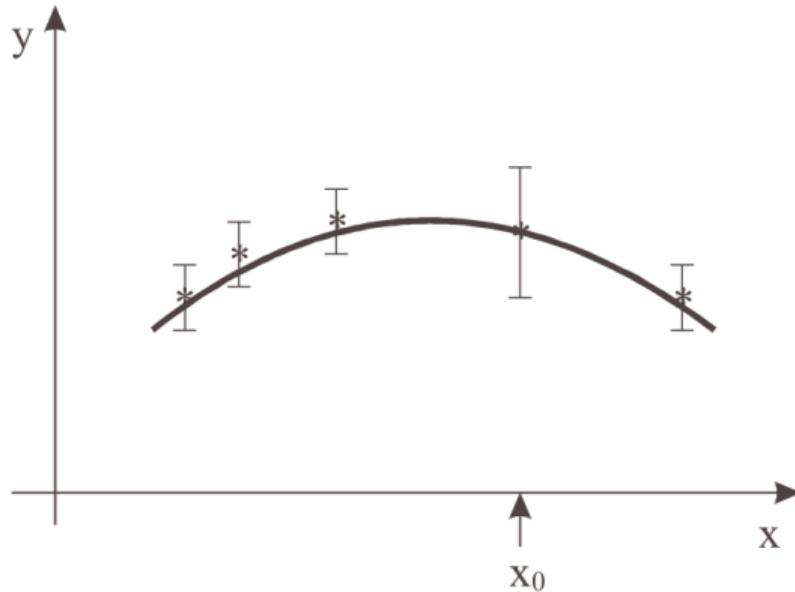
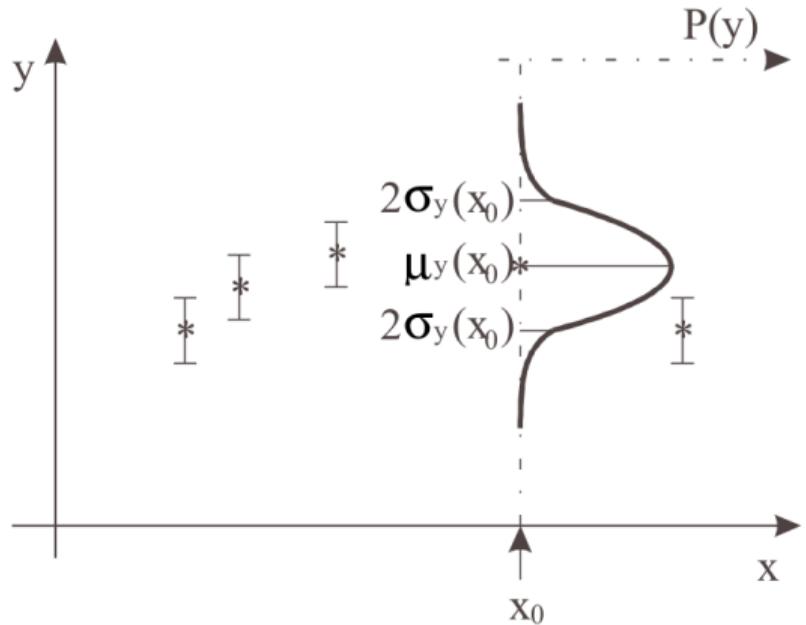
Stable implementation via Cholesky decomposition

$$\mathbf{K} + \sigma_y^2 \mathbf{I} = \mathbf{L} \mathbf{L}^\top, \quad (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} = \mathbf{L}^{-\top} \mathbf{L}^{-1},$$

$$\mu_* = (\mathbf{L}^{-1} \mathbf{k}_{f*})^\top \mathbf{L}^{-1} \mathbf{y},$$

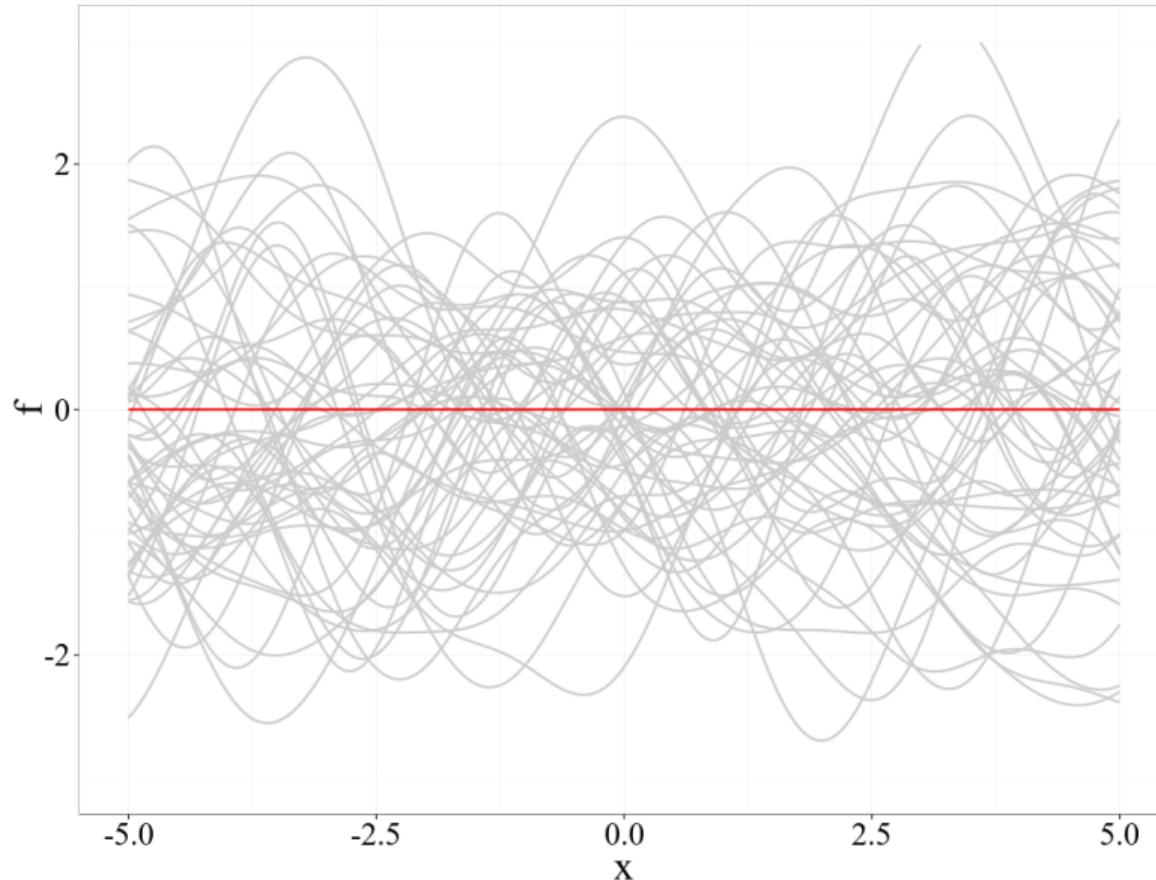
$$\sigma_*^2 = k_{**} - (\mathbf{L}^{-1} \mathbf{k}_{f*})^\top \mathbf{L}^{-1} \mathbf{k}_{f*}.$$

GP predictive distribution

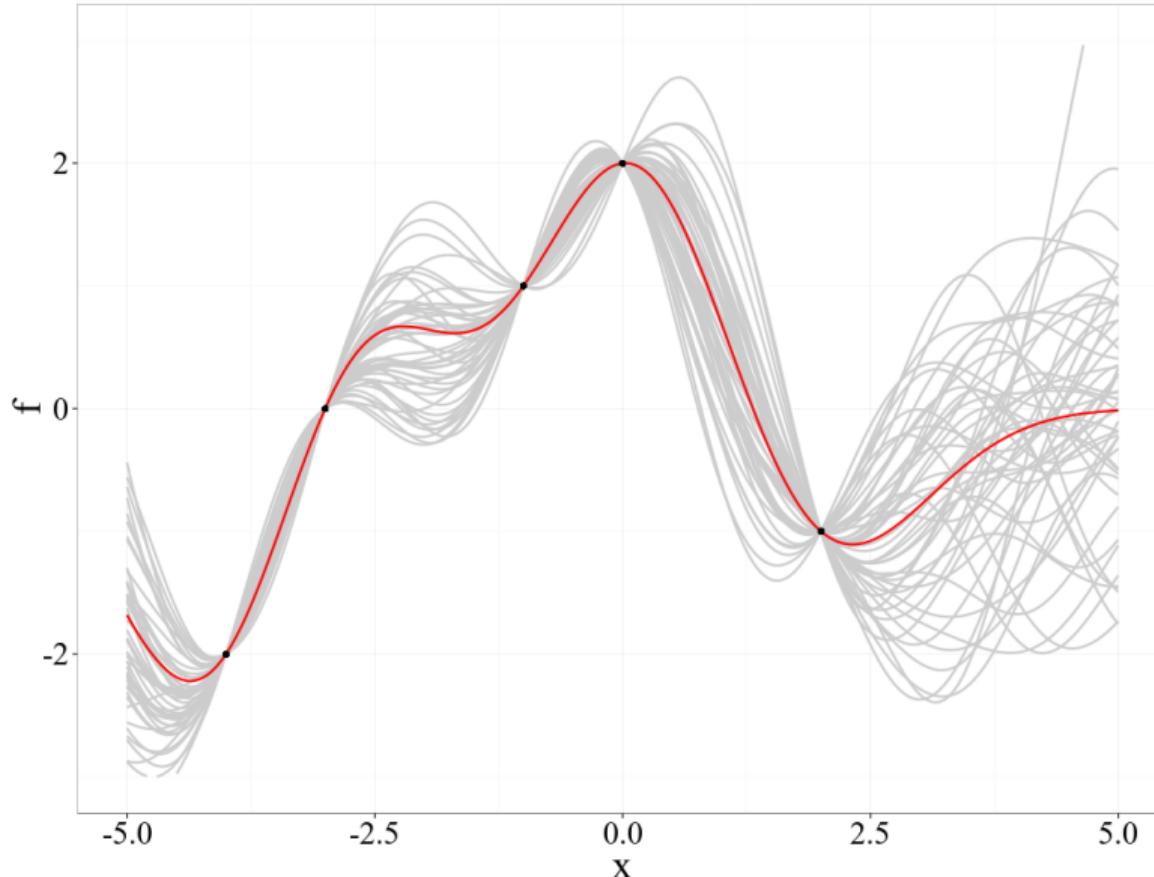


Note that each prediction is a fully defined Gaussian distribution.

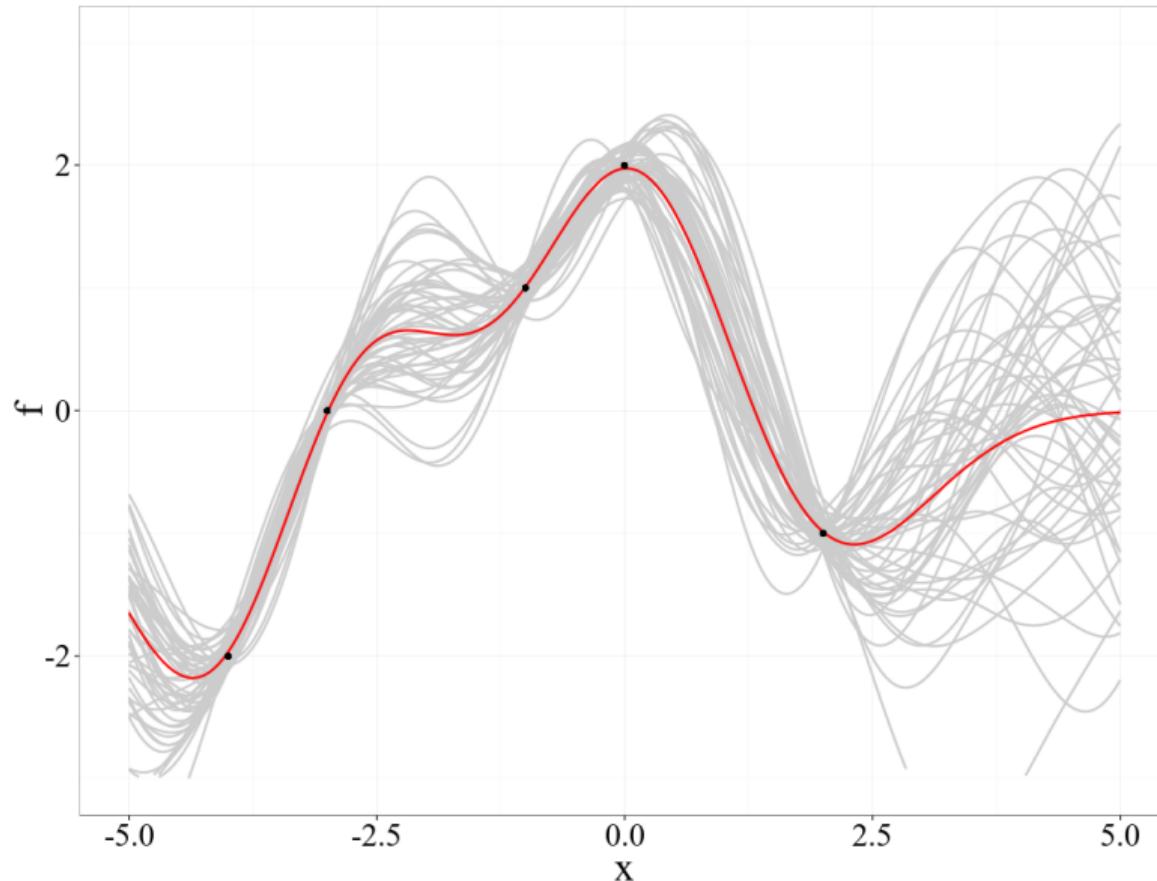
Samples from the GP posterior



Samples from the GP posterior



Samples from the GP posterior



From feature spaces to GPs

Math details

- Let $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^Q$ be a mapping function and $\mathbf{w}_i \in \mathbb{R}^Q$ be a parameter vector:

$$y_i = \mathbf{w}^\top \phi(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma_y^2).$$

- If we consider the prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_w)$ and the matrix $\Phi = \phi(\mathbf{X})$, where each row is given by $\phi(\mathbf{x}_i)$, the posterior is

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \mathbf{w}, \mathbf{X}) p(\mathbf{w})}{p(\mathbf{y} | \mathbf{X})} = \mathcal{N}\left(\mathbf{w} \middle| \frac{1}{\sigma_y^2} \mathbf{A}^{-1} \Phi \mathbf{y}, \mathbf{A}^{-1}\right),$$

where $\mathbf{A} \in \mathbb{R}^{Q \times Q} = \frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1}$.

- Predictions are made by marginalizing the parameters using the posterior:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(f_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{y}, \mathbf{X}) d\mathbf{w} \\ &= \mathcal{N}\left(f_* \middle| \frac{1}{\sigma_y^2} \phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \phi(\mathbf{x}_*)\right). \end{aligned}$$

From feature spaces to GPs

Math details

- Let $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^Q$ be a mapping function and $\mathbf{w}_i \in \mathbb{R}^Q$ be a parameter vector:

$$y_i = \mathbf{w}^\top \phi(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma_y^2).$$

- If we consider the prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_w)$ and the matrix $\Phi = \phi(\mathbf{X})$, where each row is given by $\phi(\mathbf{x}_i)$, the posterior is

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \mathbf{w}, \mathbf{X}) p(\mathbf{w})}{p(\mathbf{y} | \mathbf{X})} = \mathcal{N}\left(\mathbf{w} \middle| \frac{1}{\sigma_y^2} \mathbf{A}^{-1} \Phi \mathbf{y}, \mathbf{A}^{-1}\right),$$

where $\mathbf{A} \in \mathbb{R}^{Q \times Q} = \frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1}$.

- Predictions are made by marginalizing the parameters using the posterior:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(f_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{y}, \mathbf{X}) d\mathbf{w} \\ &= \mathcal{N}\left(f_* \middle| \frac{1}{\sigma_y^2} \phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \phi(\mathbf{x}_*)\right). \end{aligned}$$

From feature spaces to GPs

Math details

- Let $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^Q$ be a mapping function and $\mathbf{w}_i \in \mathbb{R}^Q$ be a parameter vector:

$$y_i = \mathbf{w}^\top \phi(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma_y^2).$$

- If we consider the prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_w)$ and the matrix $\Phi = \phi(\mathbf{X})$, where each row is given by $\phi(\mathbf{x}_i)$, the posterior is

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \mathbf{w}, \mathbf{X}) p(\mathbf{w})}{p(\mathbf{y} | \mathbf{X})} = \mathcal{N}\left(\mathbf{w} \middle| \frac{1}{\sigma_y^2} \mathbf{A}^{-1} \Phi \mathbf{y}, \mathbf{A}^{-1}\right),$$

where $\mathbf{A} \in \mathbb{R}^{Q \times Q} = \frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1}$.

- Predictions are made by marginalizing the parameters using the posterior:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(f_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{y}, \mathbf{X}) d\mathbf{w} \\ &= \mathcal{N}\left(f_* \middle| \frac{1}{\sigma_y^2} \phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \phi(\mathbf{x}_*)\right). \end{aligned}$$

From feature spaces to GPs

Math details

- We replace the definition of \mathbf{A} and use the *matrix inversion lemma* to rewrite the expression:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}\left(f_* \left| \frac{1}{\sigma_y^2} \phi(\mathbf{x}_*)^\top \left[\frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1} \right]^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^\top \left[\frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1} \right]^{-1} \phi(\mathbf{x}_*) \right. \right) \\ &= \mathcal{N}\left(f_* \left| \phi(\mathbf{x}_*)^\top \Sigma_w \Phi (\Phi^\top \Sigma_w \Phi + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \right. \right. \\ &\quad \left. \left. \phi(\mathbf{x}_*)^\top \Sigma_w \phi(\mathbf{x}_*) - \phi(\mathbf{x}_*)^\top \Sigma_w \Phi (\Phi^\top \Sigma_w \Phi + \sigma_y^2 \mathbf{I})^{-1} \Phi^\top \Sigma_w \phi(\mathbf{x}_*) \right) \right). \end{aligned}$$

- The kernel trick can now be applied:

$$\begin{aligned} \Phi^\top \Sigma_w \Phi &= \Phi^\top \Sigma_w^{1/2} \Sigma_w^{1/2} \Phi = \Psi^\top \Psi = k(\mathbf{X}, \mathbf{X}) = \mathbf{K}, \\ \phi(\mathbf{x}_*)^\top \Sigma_w \Phi &= k(\mathbf{x}_*, \mathbf{X}) = \mathbf{k}_{*f}, \\ \Phi^\top \Sigma_w \phi(\mathbf{x}_*) &= k(\mathbf{X}, \mathbf{x}_*) = \mathbf{k}_{f*}, \\ \phi(\mathbf{x}_*)^\top \Sigma_w \phi(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) = k_{**}. \end{aligned}$$

- Finally, the standard GP predictive distribution is obtained:

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{k}_{*f}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \mathbf{k}_{**} - \mathbf{k}_{*f}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}).$$

From feature spaces to GPs

Math details

- We replace the definition of \mathbf{A} and use the *matrix inversion lemma* to rewrite the expression:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}\left(f_* \middle| \frac{1}{\sigma_y^2} \phi(\mathbf{x}_*)^\top \left[\frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1} \right]^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^\top \left[\frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1} \right]^{-1} \phi(\mathbf{x}_*)\right) \\ &= \mathcal{N}\left(f_* \middle| \phi(\mathbf{x}_*)^\top \Sigma_w \Phi (\Phi^\top \Sigma_w \Phi + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \right. \\ &\quad \left. \phi(\mathbf{x}_*)^\top \Sigma_w \phi(\mathbf{x}_*) - \phi(\mathbf{x}_*)^\top \Sigma_w \Phi (\Phi^\top \Sigma_w \Phi + \sigma_y^2 \mathbf{I})^{-1} \Phi^\top \Sigma_w \phi(\mathbf{x}_*)\right). \end{aligned}$$

- The **kernel trick** can now be applied:

$$\begin{aligned} \Phi^\top \Sigma_w \Phi &= \Phi^\top \Sigma_w^{1/2} \Sigma_w^{1/2} \Phi = \Psi^\top \Psi = k(\mathbf{X}, \mathbf{X}) = \mathbf{K}, \\ \phi(\mathbf{x}_*)^\top \Sigma_w \Phi &= k(\mathbf{x}_*, \mathbf{X}) = \mathbf{k}_{*f}, \\ \Phi^\top \Sigma_w \phi(\mathbf{x}_*) &= k(\mathbf{X}, \mathbf{x}_*) = \mathbf{k}_{f*}, \\ \phi(\mathbf{x}_*)^\top \Sigma_w \phi(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) = k_{**}. \end{aligned}$$

- Finally, the standard GP predictive distribution is obtained:

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(k_{*f}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - k_{*f}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}).$$

From feature spaces to GPs

Math details

- We replace the definition of \mathbf{A} and use the *matrix inversion lemma* to rewrite the expression:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}\left(f_* \left| \frac{1}{\sigma_y^2} \phi(\mathbf{x}_*)^\top \left[\frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1} \right]^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^\top \left[\frac{1}{\sigma_y^2} \Phi \Phi^\top + \Sigma_w^{-1} \right]^{-1} \phi(\mathbf{x}_*) \right. \right) \\ &= \mathcal{N}\left(f_* \left| \phi(\mathbf{x}_*)^\top \Sigma_w \Phi (\Phi^\top \Sigma_w \Phi + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \right. \right. \\ &\quad \left. \phi(\mathbf{x}_*)^\top \Sigma_w \phi(\mathbf{x}_*) - \phi(\mathbf{x}_*)^\top \Sigma_w \Phi (\Phi^\top \Sigma_w \Phi + \sigma_y^2 \mathbf{I})^{-1} \Phi^\top \Sigma_w \phi(\mathbf{x}_*) \right). \end{aligned}$$

- The **kernel trick** can now be applied:

$$\begin{aligned} \Phi^\top \Sigma_w \Phi &= \Phi^\top \Sigma_w^{1/2} \Sigma_w^{1/2} \Phi = \Psi^\top \Psi = k(\mathbf{X}, \mathbf{X}) = \mathbf{K}, \\ \phi(\mathbf{x}_*)^\top \Sigma_w \Phi &= k(\mathbf{x}_*, \mathbf{X}) = \mathbf{k}_{*f}, \\ \Phi^\top \Sigma_w \phi(\mathbf{x}_*) &= k(\mathbf{X}, \mathbf{x}_*) = \mathbf{k}_{f*}, \\ \phi(\mathbf{x}_*)^\top \Sigma_w \phi(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) = k_{**}. \end{aligned}$$

- Finally, the standard GP predictive distribution is obtained:

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{k}_{*f}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{*f}(\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}).$$

Table of Contents

① Part I

- Fundamentals of Bayesian inference
- Bayesian linear regression
- The Gaussian process

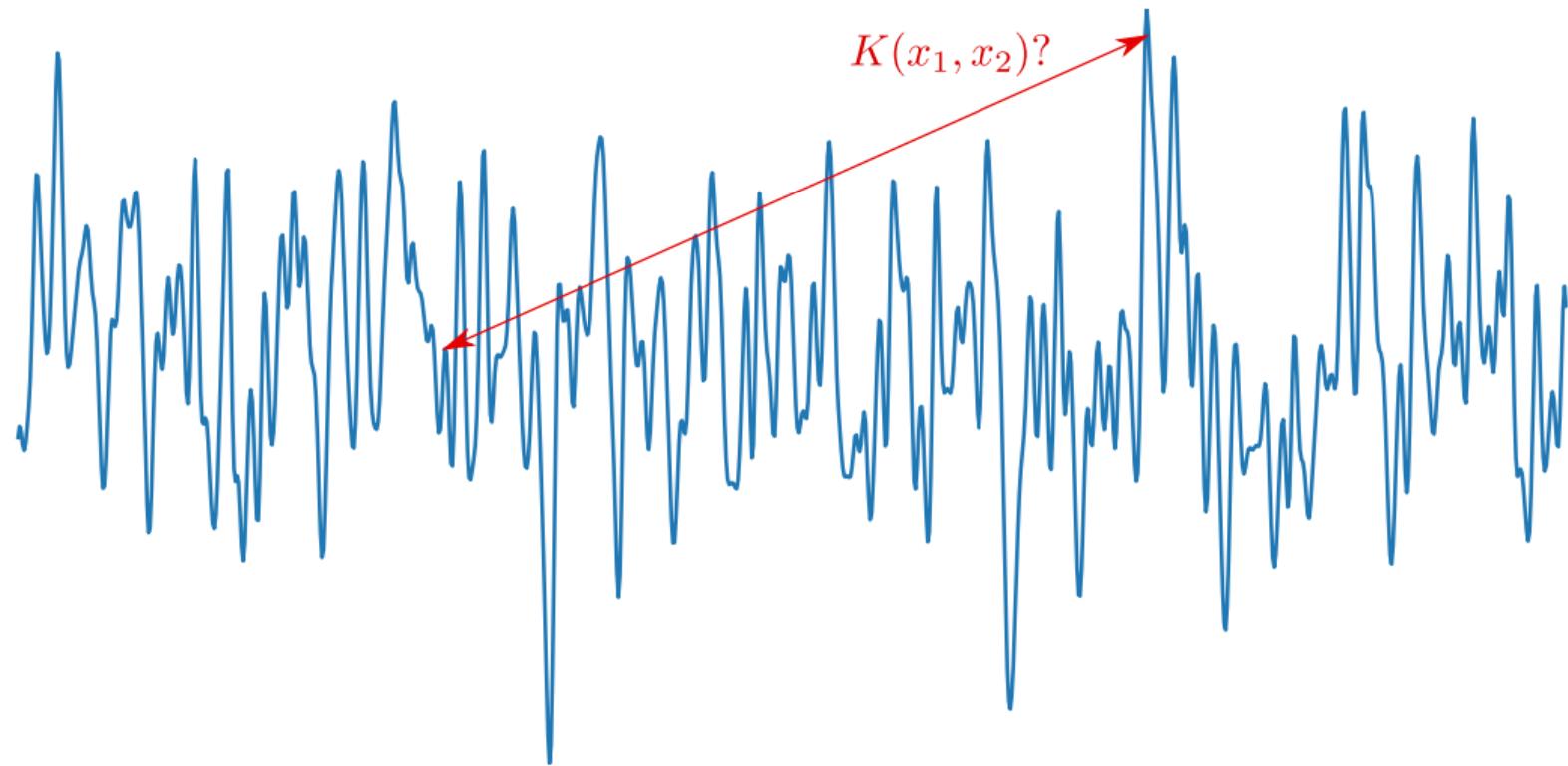
② Part II

- Learning the kernel and its hyperparameters
- Beyond Gaussian likelihood
- Sparse approximations

③ Part III

- GPs for Bayesian optimization
- Current trends on kernel design
- From GPLVM to Deep GPs
- Concluding remarks and additional topics

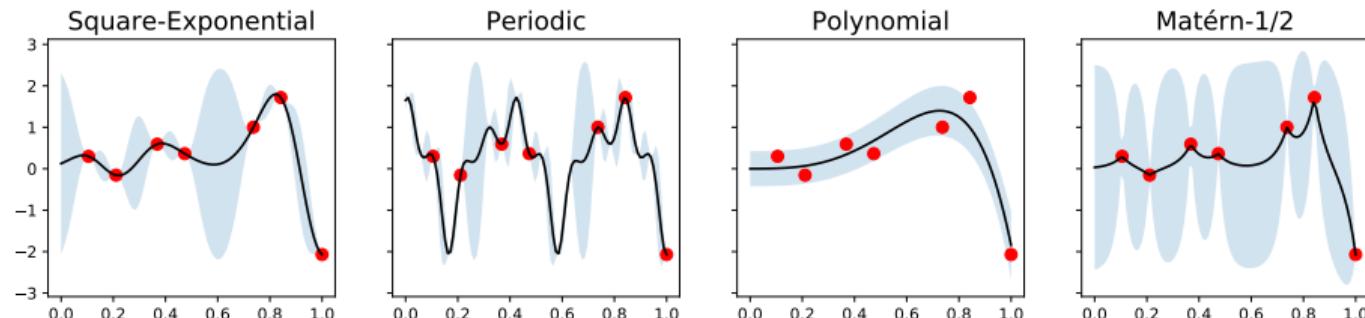
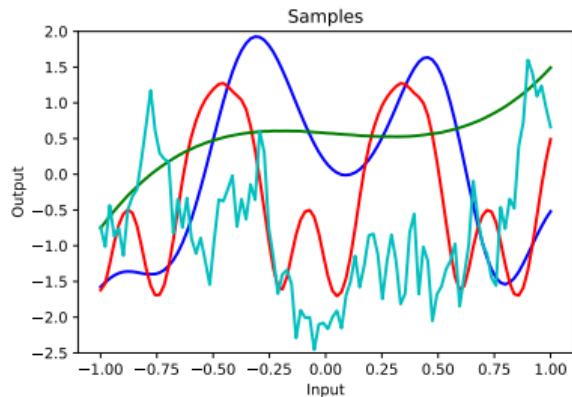
How to parameterize the kernel?



Different kernels, different predictions

- The choice of the covariance/kernel function directly affects the model behavior.
- Some of the most common options:

- Square-Exponential: $k_{\text{SE}}(x, x') = \sigma_f^2 \exp\left(-\frac{(x-x')^2}{2l^2}\right)$
- Periodic: $k_{\text{Per}}(x, x') = \sigma_f^2 \exp\left(-\frac{2 \sin(\pi|x-x'|/p)}{l^2}\right)$
- Polynomial: $k_{\text{Pol}}(x, x') = (x^\top x' + c)^d$
- Matérn-1/2: $k_{\text{M}}(x, x') = \sigma_f^2 \exp\left(-\frac{|x-x'|}{2l^2}\right)$



Creating kernel functions

- Any function that generates a **positive semidefinite** covariance matrix (i.e. $\mathbf{a}^\top \mathbf{K} \mathbf{a} \geq 0, \forall \mathbf{a}$) is a valid kernel function.
- New kernel functions may be created by summing or multiplying multiple valid kernels, with optional scalings, which further improves model expressivity:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_m A_m k_m(\mathbf{x}_i, \mathbf{x}_j) + \prod_n B_n k_n(\mathbf{x}_i, \mathbf{x}_j),$$

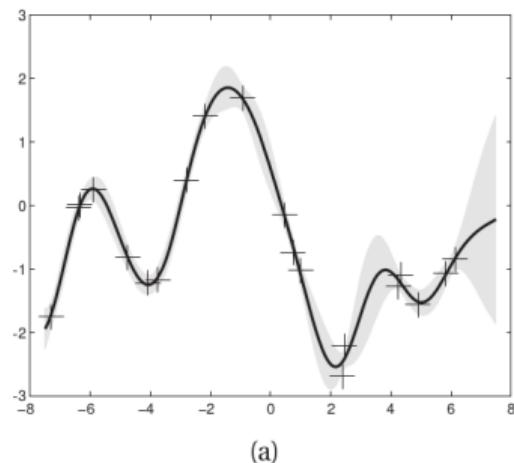
where A_m and B_n are positive real constants.

- Mapping the inputs via an arbitrary nonlinear function $g(\cdot)$ (even if it changes the input dimension) before feeding them to the kernel function also results in valid kernels:

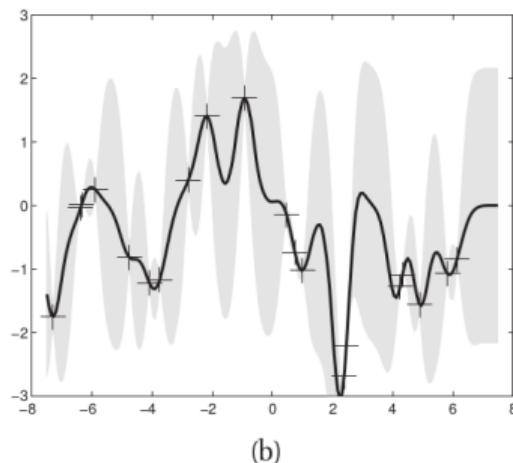
$$k_2(\mathbf{x}_i, \mathbf{x}_j) = k(g(\mathbf{x}_i), g(\mathbf{x}_j)).$$

Samples from the GP posterior

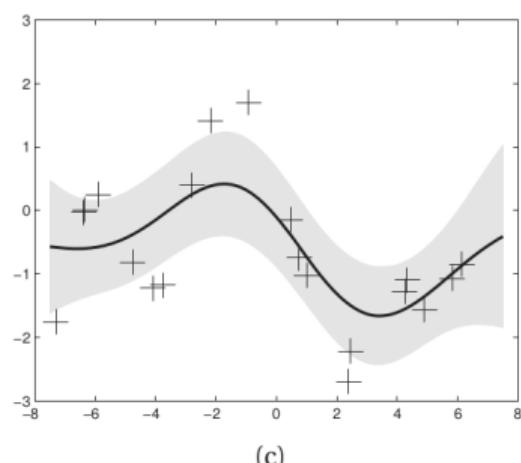
- The same kernel function and data can generate very different samples when distinct hyperparameters are used.



(a)



(b)



(c)

GP samples with larger lengthscales l (a), with smaller lengthscales l (b), and larger noise variance σ_y^2 (c).

Model selection for GPs

Hyperparameter optimization

The hyperparameter vector $\boldsymbol{\theta} = [\sigma_f^2, l_1^2, \dots, l_D^2, \sigma_y^2]^\top$ (we have included the noise variance σ_y^2) can be optimized via maximization of the marginal log-likelihood, the so-called *model evidence*:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &= \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) \\ &= -\frac{1}{2} \underbrace{\log |\mathbf{K} + \sigma_y^2 \mathbf{I}|}_{\text{model capacity}} - \frac{1}{2} \underbrace{\mathbf{y}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}}_{\text{data fit}} - \frac{N}{2} \log(2\pi).\end{aligned}$$

Stable implementation via Cholesky decomposition

$$\begin{aligned}\mathbf{K} + \sigma_y^2 \mathbf{I} &= \mathbf{L} \mathbf{L}^\top, \quad \boldsymbol{\alpha} = \mathbf{L}^{-1} \mathbf{y}, \\ \mathcal{L}(\boldsymbol{\theta}) &= - \sum_i \log L_{ii} - \frac{1}{2} \boldsymbol{\alpha}^\top \boldsymbol{\alpha} - \frac{N}{2} \log(2\pi).\end{aligned}$$

Model selection for GPs

Hyperparameter optimization

The hyperparameter vector $\theta = [\sigma_f^2, l_1^2, \dots, l_D^2, \sigma_y^2]^\top$ (we have included the noise variance σ_y^2) can be optimized via maximization of the marginal log-likelihood, the so-called *model evidence*:

$$\begin{aligned}\mathcal{L}(\theta) &= \log p(\mathbf{y}|\mathbf{X}, \theta) \\ &= -\frac{1}{2} \underbrace{\log |\mathbf{K} + \sigma_y^2 \mathbf{I}|}_{\text{model capacity}} - \frac{1}{2} \underbrace{\mathbf{y}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}}_{\text{data fit}} - \frac{N}{2} \log(2\pi).\end{aligned}$$

Stable implementation via Cholesky decomposition

$$\begin{aligned}\mathbf{K} + \sigma_y^2 \mathbf{I} &= \mathbf{L} \mathbf{L}^\top, \quad \boldsymbol{\alpha} = \mathbf{L}^{-1} \mathbf{y}, \\ \mathcal{L}(\theta) &= - \sum_i \log L_{ii} - \frac{1}{2} \boldsymbol{\alpha}^\top \boldsymbol{\alpha} - \frac{N}{2} \log(2\pi).\end{aligned}$$

Model selection for GPs

Hyperparameter optimization

- The optimization follows analytical gradients:

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \theta_i} = -\frac{1}{2} \text{Tr} \left((\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \frac{\partial (\mathbf{K} + \sigma_y^2 \mathbf{I})}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{y}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \frac{\partial (\mathbf{K} + \sigma_y^2 \mathbf{I})}{\partial \theta_i} (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}.$$

- It is not necessary to perform grid-search with cross-validations in the model selection phase!

Automatic Relevance Determination (ARD)

The optimized lengthscale hyperparameters l_1^2, \dots, l_D^2 indicate the relevance of each input dimension, where more relevant dimensions correspond to lower values of l_d^2 .

Model selection for GPs

Hyperparameter optimization

- The optimization follows analytical gradients:

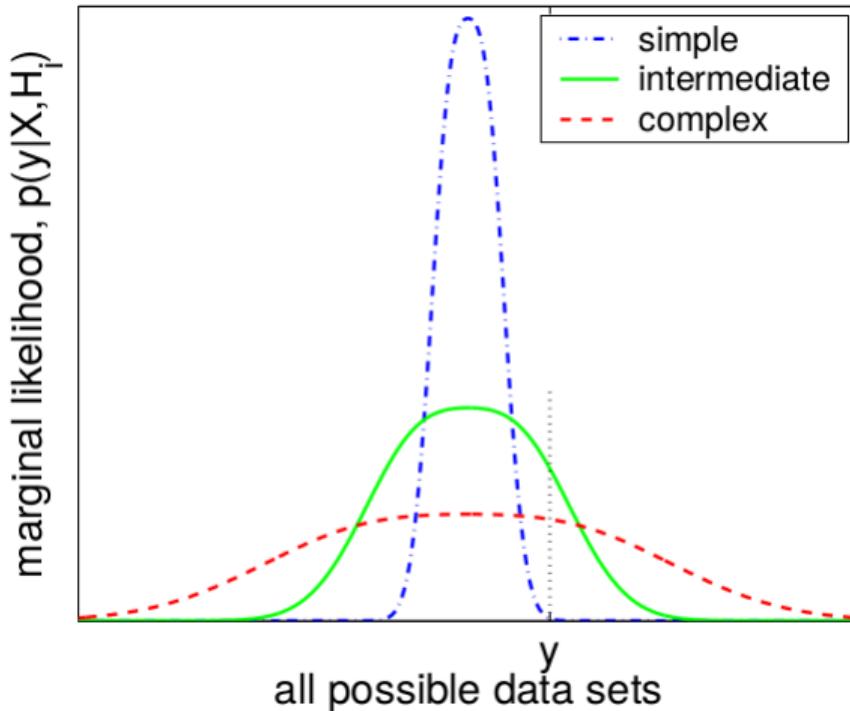
$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \theta_i} = -\frac{1}{2} \text{Tr} \left((\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \frac{\partial (\mathbf{K} + \sigma_y^2 \mathbf{I})}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{y}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \frac{\partial (\mathbf{K} + \sigma_y^2 \mathbf{I})}{\partial \theta_i} (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}.$$

- It is not necessary to perform grid-search with cross-validations in the model selection phase!

Automatic Relevance Determination (ARD)

The optimized lengthscale hyperparameters l_1^2, \dots, l_D^2 indicate the relevance of each input dimension, where more relevant dimensions correspond to lower values of l_d^2 .

Bayesian model selection



The evidence framework privileges models with intermediate complexity.

GP for nonlinear regression

Summary of the algorithm

- Estimation step (model selection)

① Initialize the hyperparameters $\boldsymbol{\theta} = [\sigma_f^2, l_1^2, \dots, l_D^2, \sigma_y^2]^\top$;

→ Example: $\sigma_f^2 = \mathbb{V}[y]$, $l_d^2 = \mathbb{V}[X_{:d}]$, $\sigma_y^2 = 0.01\sigma_f^2$.

② Repeat until convergence or for a maximum number of iterations:

① Compute the model evidence $\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$;

② Compute the evidence analytical gradients $\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$;

③ Update $\boldsymbol{\theta}$ with the computed gradients (e.g. via BFGS);

④ Return the optimized hyperparameters $\hat{\boldsymbol{\theta}}$.

- Prediction step

① Given a new input \mathbf{x}_* , return the predictive distribution

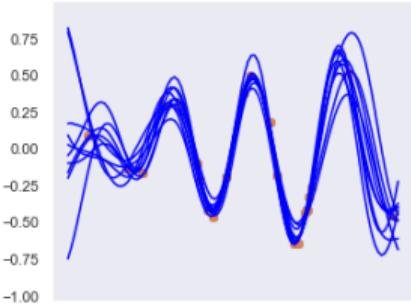
$$p(y_*|\mathbf{x}_*, \mathbf{y}, \mathbf{X}, \hat{\boldsymbol{\theta}}) = \mathcal{N}(y_*|\mu_*, \sigma_*^2 + \sigma_y^2),$$

$$\mu_* = \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y},$$

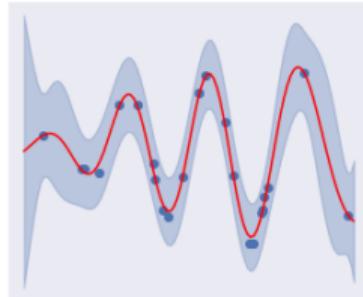
$$\sigma_*^2 = k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}.$$

Samples from non-optimized and optimized GPs

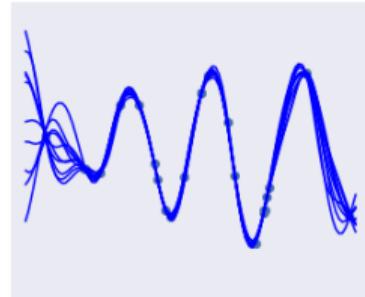
Predictive samples after 25 observations
Before hyperparameter optimization



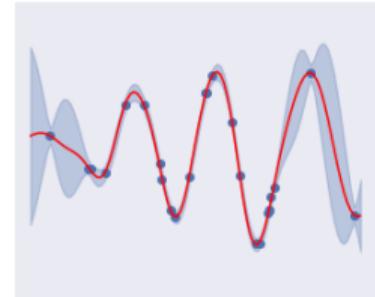
Posterior predictions after 25 observations
Before hyperparameter optimization
Evidence = 5.66e+00



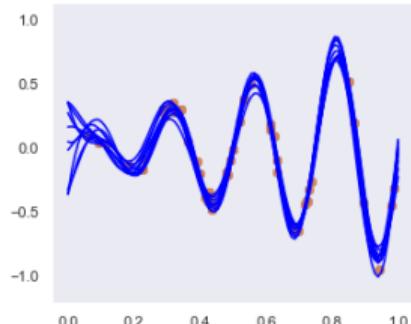
Posterior predictive samples after 25 observations
After hyperparameter optimization



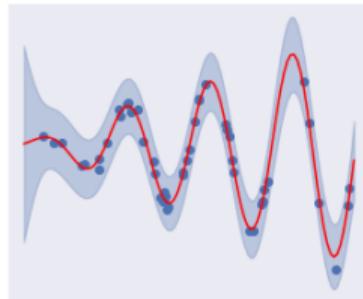
Posterior predictions after 25 observations
After hyperparameter optimization
Evidence = 2.21e+01



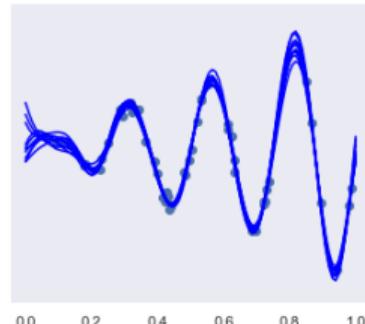
Predictive samples after 50 observations
Before hyperparameter optimization



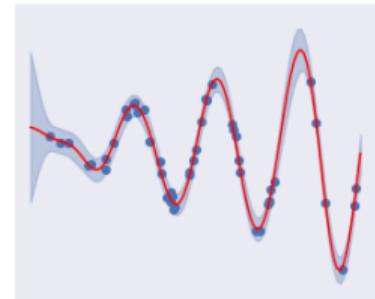
Posterior predictions after 50 observations
Before hyperparameter optimization
Evidence = 2.76e+01



Posterior predictive samples after 50 observations
After hyperparameter optimization



Posterior predictions after 50 observations
After hyperparameter optimization
Evidence = 4.87e+01



Samples from non-optimized and optimized GPs

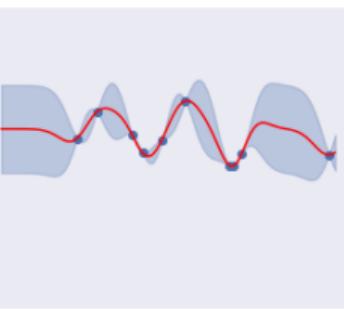
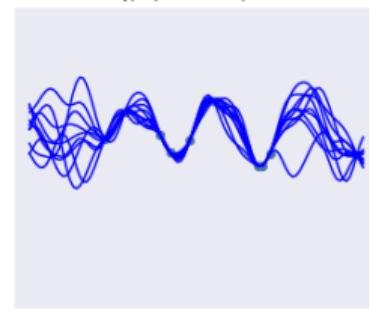
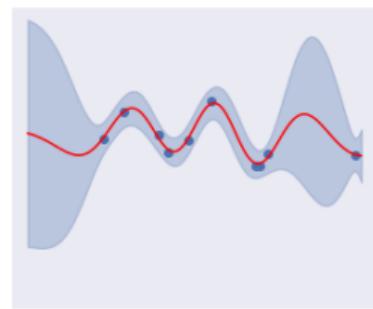
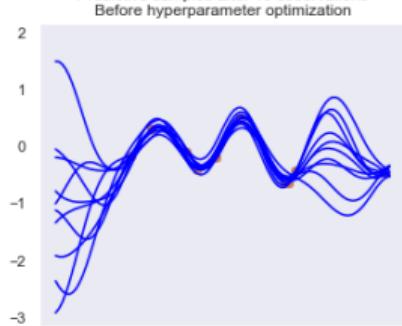
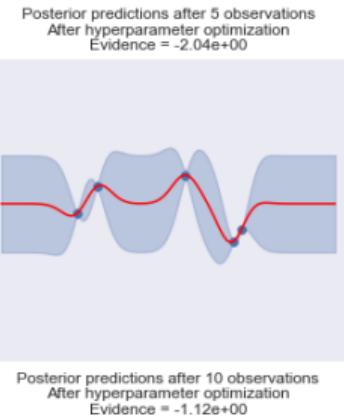
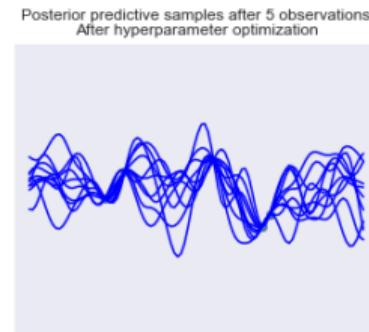
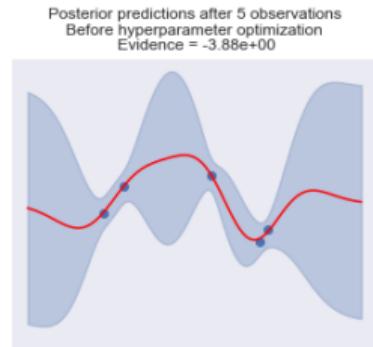
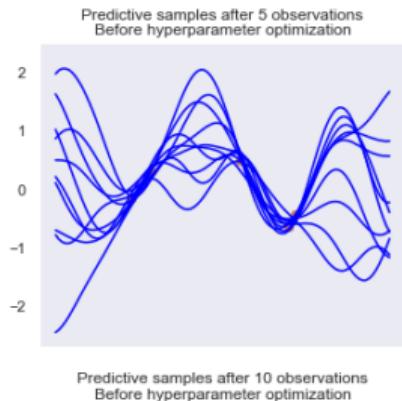


Table of Contents

① Part I

- Fundamentals of Bayesian inference
- Bayesian linear regression
- The Gaussian process

② Part II

- Learning the kernel and its hyperparameters
- Beyond Gaussian likelihood
- Sparse approximations

③ Part III

- GPs for Bayesian optimization
- Current trends on kernel design
- From GPLVM to Deep GPs
- Concluding remarks and additional topics

Beyond Gaussian likelihood

- Gaussian likelihood is suitable only when Gaussian observations are expected:

$$\begin{aligned}\mathbf{y} &= \mathbf{f} + \boldsymbol{\epsilon}, \\ p(\mathbf{f}|\mathbf{X}) &\sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}), \\ p(\mathbf{y}|\mathbf{f}) &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}).\end{aligned}$$

- For non-Gaussian data, we choose other likelihood $p(\mathbf{y}|\mathbf{f}, \phi)$ parameterized by ϕ .
 - Heavy-tailed regression: e.g. Student-t or Laplace likelihoods;
 - Non-negative regression: e.g. gamma or exponential likelihoods;
 - Discrete regression: e.g. Poisson likelihood;
 - Binary classification: e.g. Bernoulli likelihood;
 - Multiclass classification: e.g. categorical likelihood.
- **Problem:** The tractable integrals of GP models require a Gaussian likelihood.

Beyond Gaussian likelihood

- Gaussian likelihood is suitable only when Gaussian observations are expected:

$$\begin{aligned}\mathbf{y} &= \mathbf{f} + \boldsymbol{\epsilon}, \\ p(\mathbf{f}|\mathbf{X}) &\sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}), \\ p(\mathbf{y}|\mathbf{f}) &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}).\end{aligned}$$

- For non-Gaussian data, we choose other likelihood $p(\mathbf{y}|\mathbf{f}, \phi)$ parameterized by ϕ .
 - Heavy-tailed regression: e.g. **Student-t** or **Laplace** likelihoods;
 - Non-negative regression: e.g. **gamma** or **exponential** likelihoods;
 - Discrete regression: e.g. **Poisson** likelihood;
 - Binary classification: e.g. **Bernoulli** likelihood;
 - Multiclass classification: e.g. **categorical** likelihood.
- Problem: The tractable integrals of GP models require a Gaussian likelihood.

Beyond Gaussian likelihood

- Gaussian likelihood is suitable only when Gaussian observations are expected:

$$\begin{aligned}\mathbf{y} &= \mathbf{f} + \boldsymbol{\epsilon}, \\ p(\mathbf{f}|\mathbf{X}) &\sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}), \\ p(\mathbf{y}|\mathbf{f}) &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}).\end{aligned}$$

- For non-Gaussian data, we choose other likelihood $p(\mathbf{y}|\mathbf{f}, \phi)$ parameterized by ϕ .
 - Heavy-tailed regression: e.g. **Student-t** or **Laplace** likelihoods;
 - Non-negative regression: e.g. **gamma** or **exponential** likelihoods;
 - Discrete regression: e.g. **Poisson** likelihood;
 - Binary classification: e.g. **Bernoulli** likelihood;
 - Multiclass classification: e.g. **categorical** likelihood.
- **Problem:** The tractable integrals of GP models require a Gaussian likelihood.

Beyond Gaussian likelihood

- Consider a **binary classification task** with Bernoulli observations:

$$p(y) = \mathcal{B}(y|a) = a^y(1 - a^{1-y}), \quad a = p(y = 1).$$

- A Bernoulli likelihood $p(y|f)$ is obtained with $a = \sigma(f) = p(y = 1|f)$, where $\sigma(f) \in [0, 1], \forall f \in \mathbb{R}$, is a **sigmoid function**:

$$\begin{aligned} p(y|f) &= \sigma(f)^y(1 - \sigma(f))^{1-y}, \\ p(\mathbf{y}|\mathbf{f}) &= \prod_{i=1}^N \sigma(f_i)^{y_i}(1 - \sigma(f_i))^{1-y_i}. \end{aligned}$$



- $\sigma(f) = a$ is the **response (or transformation) function**.
- $\sigma^{-1}(a) = f$ is the **link function**, which must be monotonic.
- $f = \theta^\top x$ results in **logistic regression**, an example of a generalized linear model.

Beyond Gaussian likelihood

- Consider a **binary classification task** with Bernoulli observations:

$$p(y) = \mathcal{B}(y|a) = a^y(1 - a^{1-y}), \quad a = p(y=1).$$

- A Bernoulli likelihood $p(y|f)$ is obtained with $a = \sigma(f) = p(y=1|f)$, where $\sigma(f) \in [0, 1], \forall f \in \mathbb{R}$, is a **sigmoid function**:

$$\begin{aligned} p(y|f) &= \sigma(f)^y(1 - \sigma(f))^{1-y}, \\ p(\mathbf{y}|\mathbf{f}) &= \prod_{i=1}^N \sigma(f_i)^{y_i}(1 - \sigma(f_i))^{1-y_i}. \end{aligned}$$



- $\sigma(f) = a$ is the **response** (or transformation) function.
- $\sigma^{-1}(a) = f$ is the **link function**, which must be monotonic.
- $f = \theta^\top x$ results in logistic regression, an example of a generalized linear model.

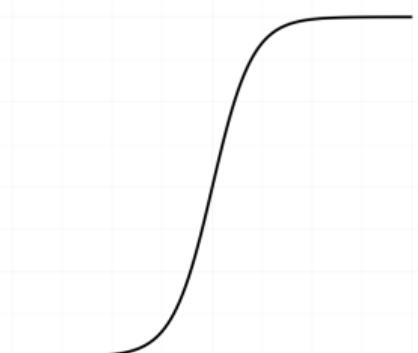
Beyond Gaussian likelihood

- Consider a **binary classification task** with Bernoulli observations:

$$p(y) = \mathcal{B}(y|a) = a^y(1 - a^{1-y}), \quad a = p(y=1).$$

- A Bernoulli likelihood $p(y|f)$ is obtained with $a = \sigma(f) = p(y=1|f)$, where $\sigma(f) \in [0, 1], \forall f \in \mathbb{R}$, is a **sigmoid function**:

$$\begin{aligned} p(y|f) &= \sigma(f)^y(1 - \sigma(f))^{1-y}, \\ p(\mathbf{y}|\mathbf{f}) &= \prod_{i=1}^N \sigma(f_i)^{y_i}(1 - \sigma(f_i))^{1-y_i}. \end{aligned}$$



- $\sigma(f) = a$ is the **response (or transformation) function**.
- $\sigma^{-1}(a) = f$ is the **link function**, which must be monotonic.
- $f = \theta^\top x$ results in **logistic regression**, an example of a **generalized linear model**.

Beyond Gaussian likelihood

- GP models for binary classification can be similarly obtained:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i} (1 - \sigma(f_i)^{1-y_i}), \quad \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f).$$

- \mathbf{f} is a **nuisance variable** and we are interested in predictions in the observed space:

$$\underbrace{p(f_*|x_*, \mathbf{X}, \mathbf{y})}_{\text{latent space}} = \int p(f_*|\mathbf{f}, x_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f},$$

$$\underbrace{p(y_*|x_*, \mathbf{X}, \mathbf{y})}_{\text{observed space}} = \int p(y_*|f_*) p(f_*|x_*, \mathbf{X}, \mathbf{y}) df_*.$$

- For a binary classification task with Bernoulli likelihood:

$$p(y_*|x_*, \mathbf{X}, \mathbf{y}) = \int \overbrace{\sigma(f_*)^{y_*} (1 - \sigma(f_*)^{1-y_*})}^{p(y_*|f_*)} \overbrace{p(f_*|\mathbf{f}, x_*, \mathbf{X})}^{\text{cond. Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}^{\text{posterior}} d\mathbf{f} df_*,$$

$$p(y_* = 1|x_*, \mathbf{X}, \mathbf{y}) = \int \sigma(f_*) p(f_*|\mathbf{f}, x_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f} df_*.$$

Beyond Gaussian likelihood

- GP models for binary classification can be similarly obtained:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i} (1 - \sigma(f_i)^{1-y_i}), \quad \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f).$$

- \mathbf{f} is a **nuisance variable** and we are interested in predictions in the observed space:

$$\underbrace{p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{latent space}} = \int p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f} | \mathbf{X}, \mathbf{y}) d\mathbf{f},$$

$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{observed space}} = \int p(y_* | f_*) p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) df_*.$$

- For a binary classification task with Bernoulli likelihood:

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \overbrace{\sigma(f_*)^{y_*} (1 - \sigma(f_*)^{1-y_*})}^{p(y_* | f_*)} \overbrace{p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X})}^{\text{cond. Gaussian}} \overbrace{p(\mathbf{f} | \mathbf{X}, \mathbf{y})}^{\text{posterior}} df_* df_*,$$

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \sigma(f_*) p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f} | \mathbf{X}, \mathbf{y}) df_* df_*,$$

Beyond Gaussian likelihood

- GP models for binary classification can be similarly obtained:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i} (1 - \sigma(f_i)^{1-y_i}), \quad \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f).$$

- \mathbf{f} is a **nuisance variable** and we are interested in predictions in the observed space:

$$\underbrace{p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{latent space}} = \int p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f},$$

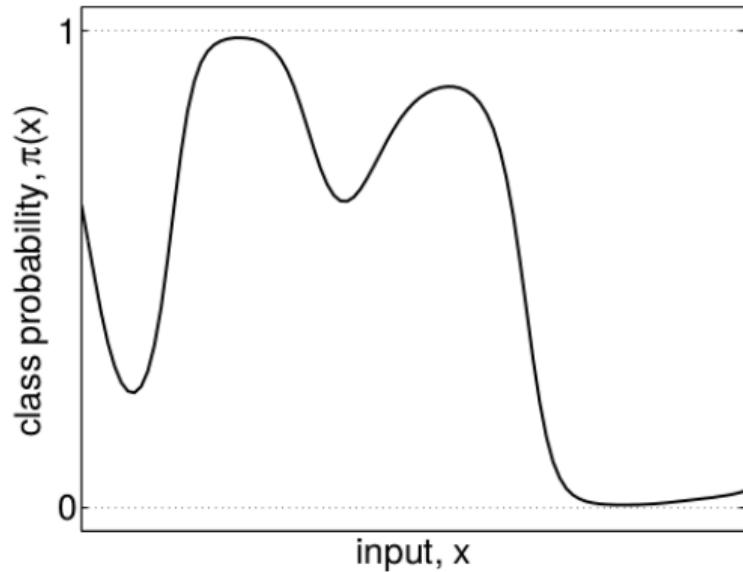
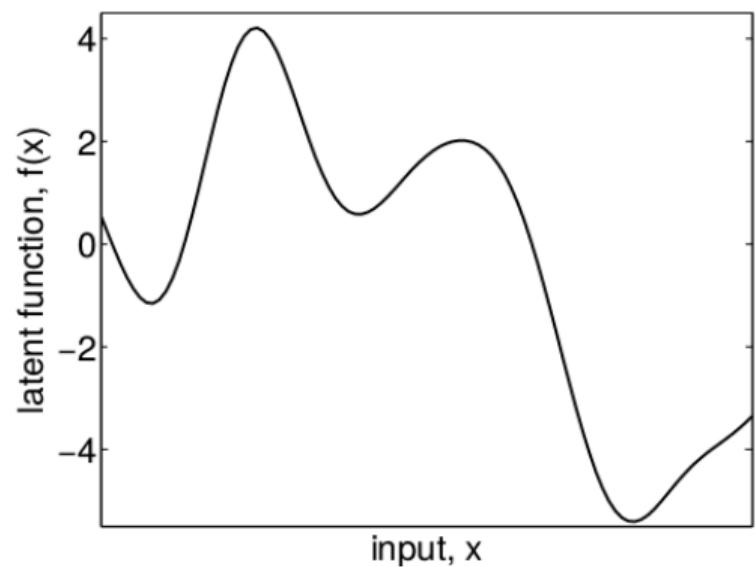
$$\underbrace{p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{observed space}} = \int p(y_*|f_*) p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) df_*.$$

- For a binary classification task with Bernoulli likelihood:

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \overbrace{\sigma(f_*)^{y_*} (1 - \sigma(f_*)^{1-y_*})}^{p(y_*|f_*)} \overbrace{p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X})}^{\text{cond. Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}^{\text{posterior}} d\mathbf{f} df_*,$$

$$p(y_* = 1|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \sigma(f_*) p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f} df_*.$$

Beyond Gaussian likelihood



(Rasmussen and Williams, 2006)

Beyond Gaussian likelihood

- **Problem:** The posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is intractable for non-Gaussian likelihood:

$$\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}^{\text{non-Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{Gaussian}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{non-Gaussian}}}.$$

- **Problem:** The marginal likelihood is also intractable for a non-Gaussian likelihood:

$$\begin{aligned}\widehat{p(\mathbf{y}|\mathbf{X})} &= \widehat{\int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f}} \\ &= \int \underbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}_{\text{non-Gaussian}} \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}} d\mathbf{f}.\end{aligned}$$

- We need approximations for both the posterior and the marginal likelihood.

Beyond Gaussian likelihood

- **Problem:** The posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is intractable for non-Gaussian likelihood:

$$\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}^{\text{non-Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{Gaussian}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{non-Gaussian}}}.$$

- **Problem:** The marginal likelihood is also intractable for a non-Gaussian likelihood:

$$\begin{aligned}\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marg. likel.}} &= \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f} \\ &= \int \underbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}_{\text{non-Gaussian}} \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}} d\mathbf{f}.\end{aligned}$$

- We need approximations for both the posterior and the marginal likelihood.

Beyond Gaussian likelihood

- **Problem:** The posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is intractable for non-Gaussian likelihood:

$$\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}^{\text{non-Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{Gaussian}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{non-Gaussian}}}.$$

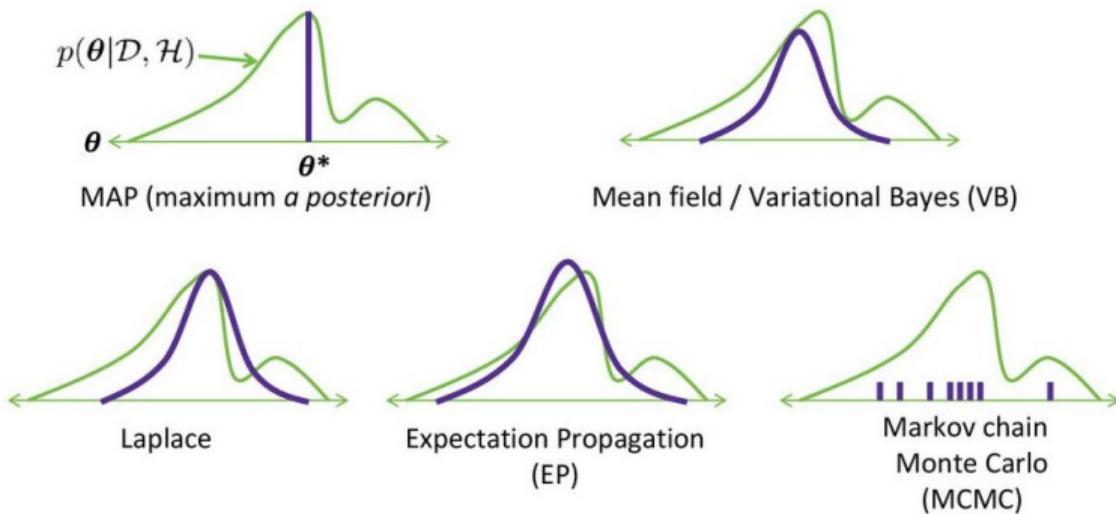
- **Problem:** The marginal likelihood is also intractable for a non-Gaussian likelihood:

$$\begin{aligned}\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marg. likel.}} &= \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f} \\ &= \int \underbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}_{\text{non-Gaussian}} \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}} d\mathbf{f}.\end{aligned}$$

- We need **approximations** for both the posterior and the marginal likelihood.

Beyond Gaussian likelihood

- The approximations are usually performed by
 - Laplace approximation;
 - Variational inference;
 - Expectation Propagation;
 - Sampling methods (e.g. MCMC).



(Ben-Elazar, *et al.*, 2017)

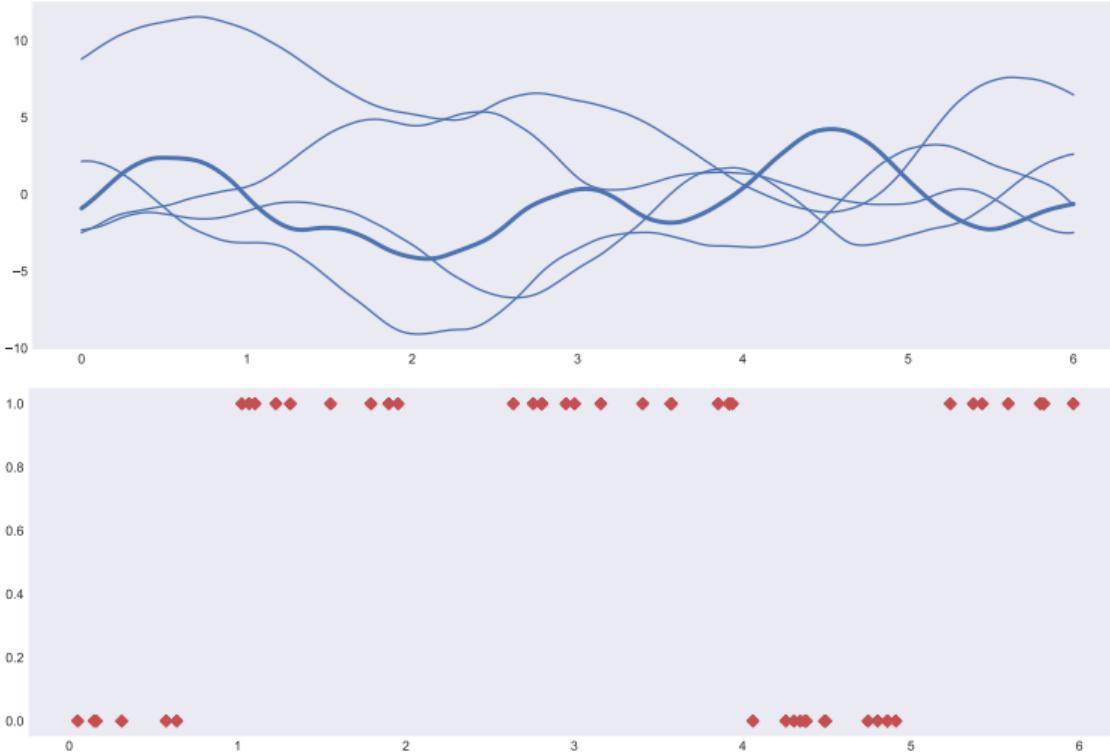
Beyond Gaussian likelihood

- For a **Gaussian approximate posterior** $q(\mathbf{f}|\mathbf{X}, \mathbf{y})$, the prediction becomes:

$$\begin{aligned} p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int \underbrace{p(y_*|f_*)}_{\text{non-Gaussian}} p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) \underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{intractable}} d\mathbf{f} df_* \\ &\approx \int p(y_*|f_*) \left[\int \underbrace{p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X})}_{\text{cond. Gaussian}} \underbrace{q(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{Gaussian approx.}} d\mathbf{f} \right] df_*. \end{aligned}$$

- The integral on \mathbf{f} is tractable, since it involves only Gaussians.
- The integral on f_* is usually intractable, but since it is one-dimensional, it can be cheaply approximated by Monte Carlo, Gauss-Hermite quadrature, etc.
- Importantly, the previous frameworks also approximate the marginal likelihood.

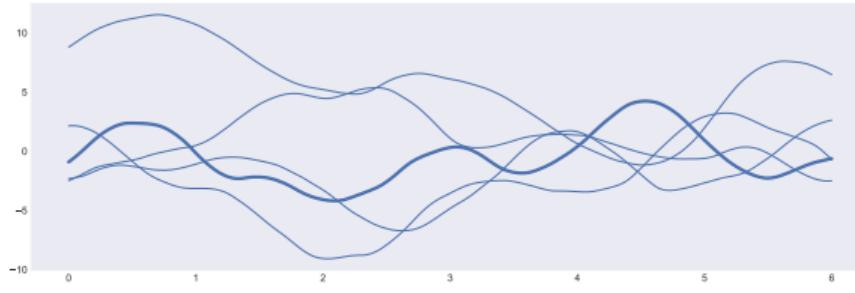
Beyond Gaussian likelihood



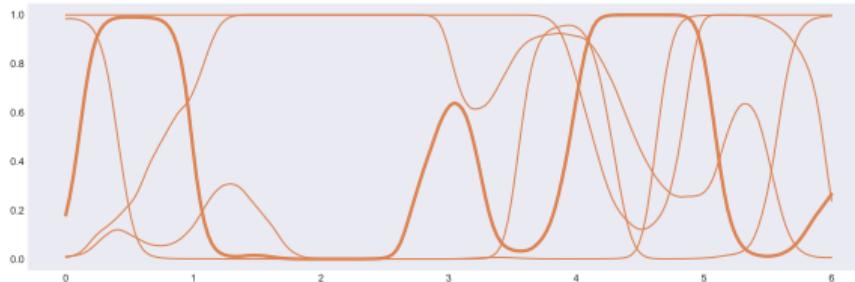
$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f)$$

$$(\mathbf{X}, \mathbf{y}), y_i \in \{0, 1\}, \forall i$$

Beyond Gaussian likelihood

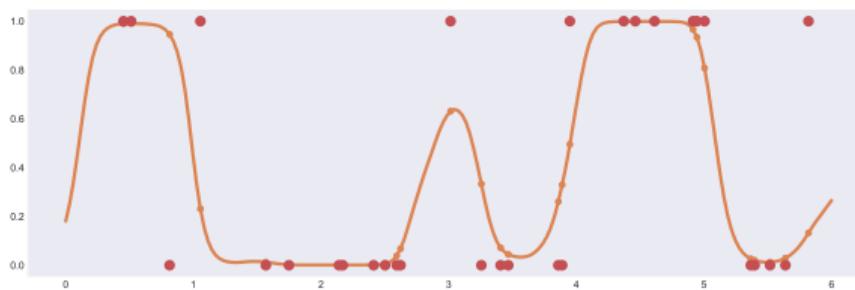


$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f)$$



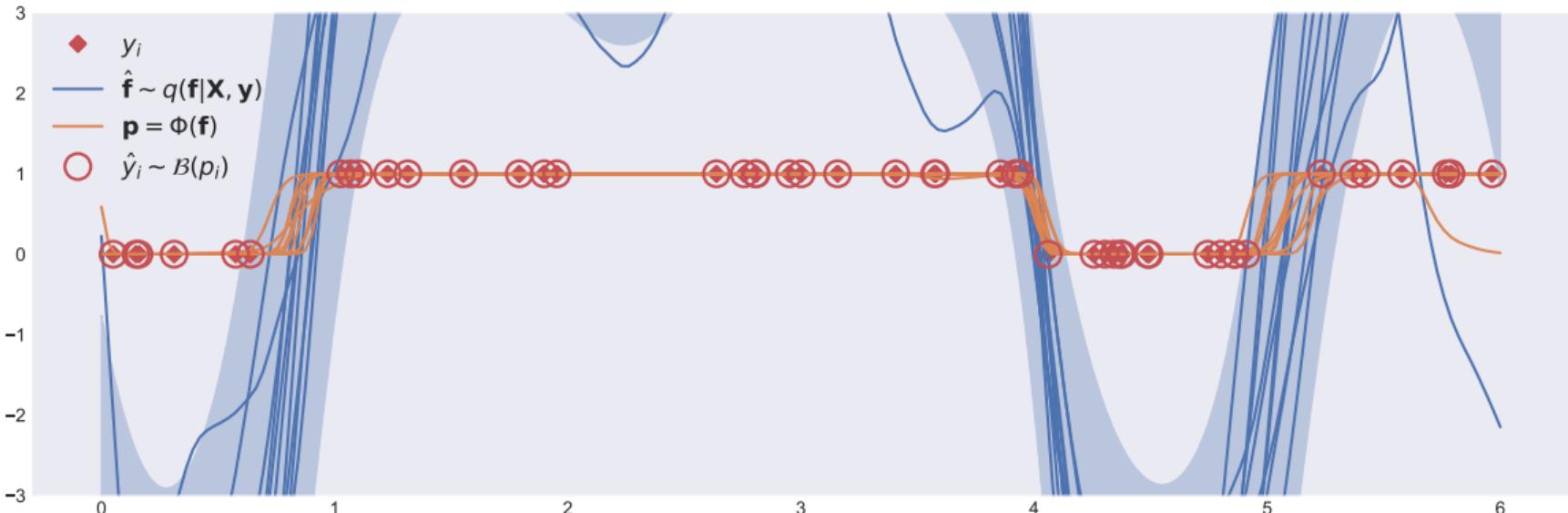
$$\mathbf{p} = \Phi(\mathbf{f}) \triangleq \int_{-\infty}^f \mathcal{N}(f|0, 1) df$$

(Normal CDF)



$$y_i \sim \mathcal{B}(p_i), \forall i$$

Beyond Gaussian likelihood



- Gaussian approximate posterior $q(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{\Sigma})$.
- Variational inference with **GPflow** (whose documentation this illustration was based on).
- Latent samples are continuous, but samples in the observed space remain discrete.

Warped GPs

- A **warping function** may be a useful preprocessing step for non-Gaussian observations.
 - For instance, if $y > 0$, $f = \log(y)$.
- We can also consider a parameterized warping function $w(\cdot|\phi)$:

$$f = w(y|\phi).$$

- $w(\cdot|\phi)$ must be monotonic, so that $y = w^{-1}(f|\phi)$ is unambiguous.
- If $p(f|X)$ is a GP, we obtain a **warped GP model** (Snelson *et al.*, 2004).
- A **Bayesian warped GP** is obtained by directly modeling the inverse warping function (Lázaro-Gredilla, 2012):

$$y = g(f(x)) + \epsilon,$$

where both $f(\cdot)$ and $g(\cdot)$ have GP priors and ϵ is an observation noise.

→ Analogous to a 2-layer Deep GP!

Warped GPs

- A **warping function** may be a useful preprocessing step for non-Gaussian observations.
 - For instance, if $y > 0$, $f = \log(y)$.
- We can also consider a parameterized warping function $w(\cdot|\phi)$:

$$f = w(y|\phi).$$

- $w(\cdot|\phi)$ must be monotonic, so that $y = w^{-1}(f|\phi)$ is unambiguous.
- If $p(f|X)$ is a GP, we obtain a **warped GP model** (Snelson *et al.*, 2004).
- A **Bayesian warped GP** is obtained by directly modeling the inverse warping function (Lázaro-Gredilla, 2012):

$$y = g(f(x)) + \epsilon,$$

where both $f(\cdot)$ and $g(\cdot)$ have GP priors and ϵ is an observation noise.

→ Analogous to a 2-layer Deep GP!

Warped GPs

- A **warping function** may be a useful preprocessing step for non-Gaussian observations.
 - For instance, if $y > 0$, $f = \log(y)$.
- We can also consider a parameterized warping function $w(\cdot|\phi)$:

$$f = w(y|\phi).$$

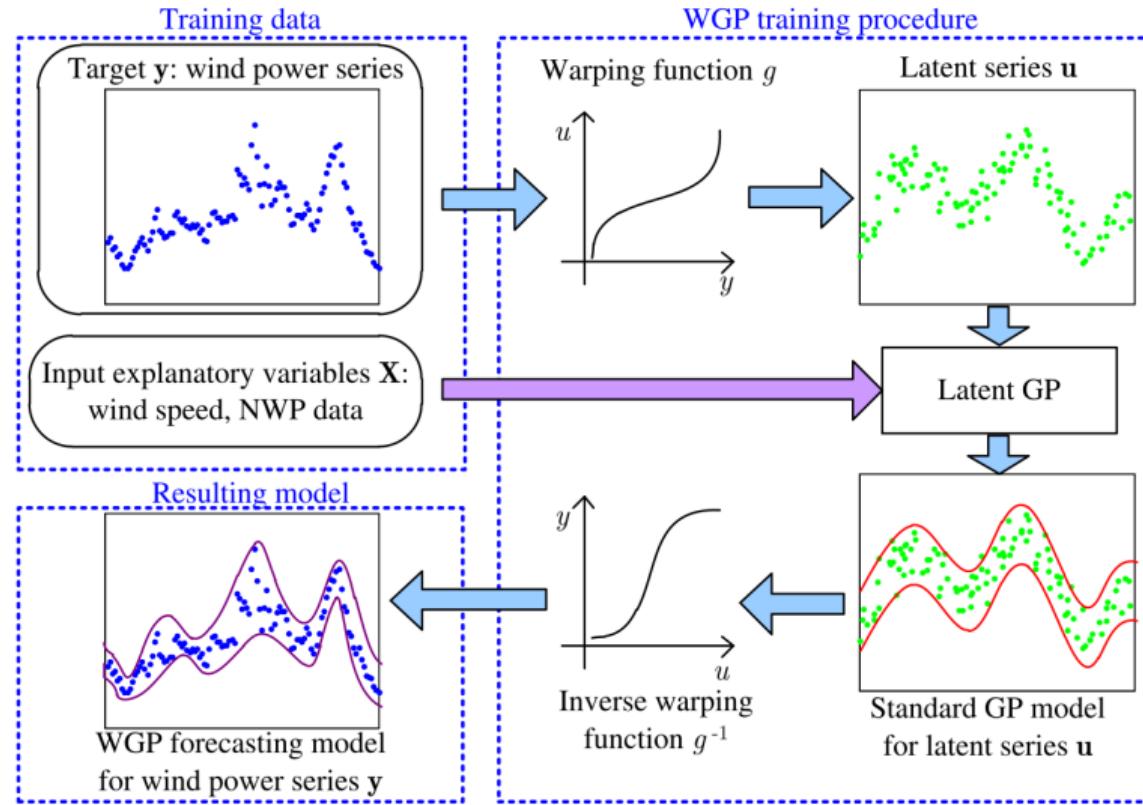
- $w(\cdot|\phi)$ must be monotonic, so that $y = w^{-1}(f|\phi)$ is unambiguous.
- If $p(f|X)$ is a GP, we obtain a **warped GP model** (Snelson *et al.*, 2004).
- A **Bayesian warped GP** is obtained by directly modeling the inverse warping function (Lázaro-Gredilla, 2012):

$$y = g(f(x)) + \epsilon,$$

where both $f(\cdot)$ and $g(\cdot)$ have GP priors and ϵ is an observation noise.

- Analogous to a 2-layer Deep GP!

Warped GPs



Kou, Peng et al. "Sparse online warped Gaussian process for wind power probabilistic forecasting". Applied Energy, 2013.

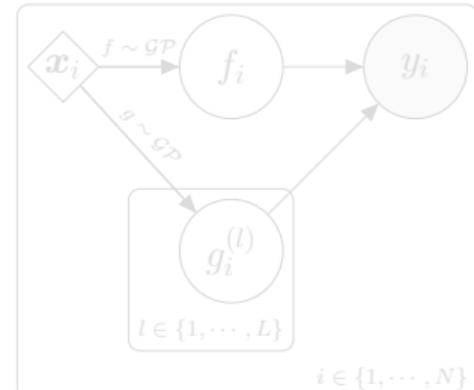
Beyond Gaussian likelihood

- Flexible models can be obtained by considering **multiple GP priors**.
- **Heteroscedastic Gaussian likelihood**, i.e. input-dependent noise variance (Goldberg *et al.*, 1998; Lázaro-Gredilla and Titsias, 2011):

$$y = f(\mathbf{x}) + \epsilon,$$
$$f|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f), \quad \epsilon|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \exp(g(\mathbf{x}))), \quad g|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_g, \mathbf{K}_g).$$

- Chained GPs with L GP priors for the likelihood parameters ϕ (Saul *et al*, 2016):

$$\phi = [\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}],$$
$$p(\mathbf{y}|\mathbf{f}, \phi) = p(\mathbf{y}|\mathbf{f}, \mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}),$$
$$\mathbf{f}|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$
$$\mathbf{g}^{(l)}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_{g^{(l)}}, \mathbf{K}_{g^{(l)}}), \quad 1 \leq l \leq L.$$



Beyond Gaussian likelihood

- Flexible models can be obtained by considering **multiple GP priors**.
- **Heteroscedastic Gaussian likelihood**, i.e. input-dependent noise variance (Goldberg *et al.*, 1998; Lázaro-Gredilla and Titsias, 2011):

$$y = f(\mathbf{x}) + \epsilon,$$

$$f|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f), \quad \epsilon|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \exp(g(\mathbf{x}))), \quad g|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_g, \mathbf{K}_g).$$

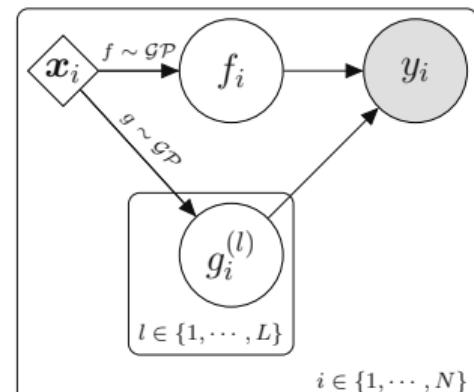
- **Chained GPs** with L GP priors for the likelihood parameters ϕ (Saul *et al.*, 2016):

$$\phi = [\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}],$$

$$p(\mathbf{y}|\mathbf{f}, \phi) = p(\mathbf{y}|\mathbf{f}, \mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}),$$

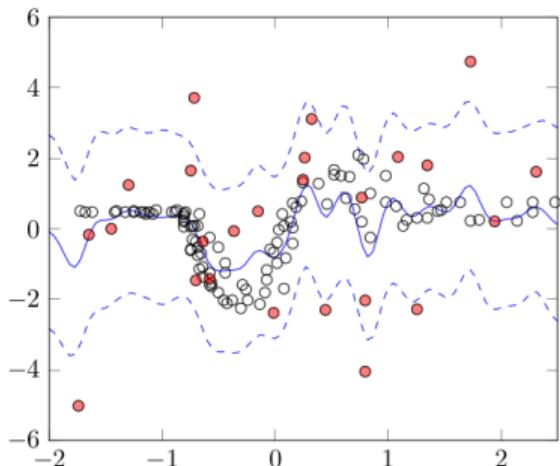
$$\mathbf{f}|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$

$$\mathbf{g}^{(l)}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_{g^{(l)}}, \mathbf{K}_{g^{(l)}}), \quad 1 \leq l \leq L.$$

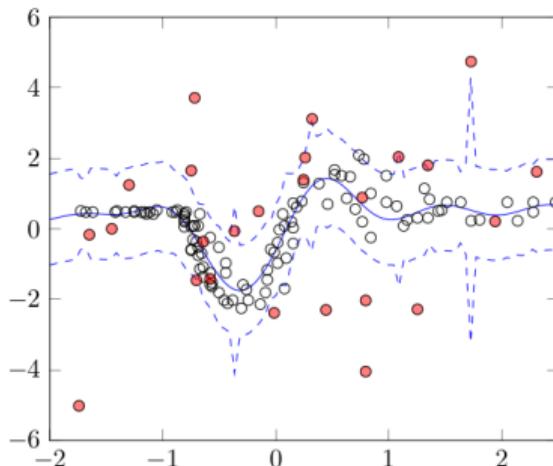


Chained GPs

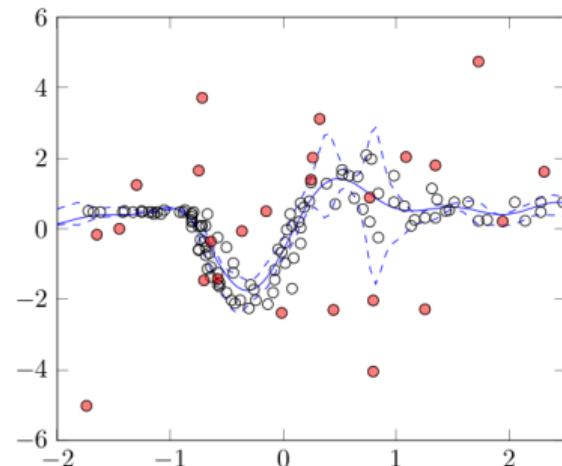
Standard Gaussian Process



Heteroscedastic Gaussian



Heteroscedastic Student-t



Saul, Alan, et al. “**Chained Gaussian processes**”. AISTATS, 2016.

$$y_i \sim \text{St}(\mu = f(\mathbf{x}_i), \sigma^2 = \exp(g(\mathbf{x}_i)), \nu),$$

$$\mathbf{f}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_f, \mathbf{K}_f),$$

$$\mathbf{g}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_g, \mathbf{K}_g).$$

Beyond Gaussian likelihood

Summary

- GP models are tractable only for Gaussian likelihood.
- Other noise models are equivalent to consider **non-Gaussian likelihoods**.
- We approximate the posterior $\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{predictions}}$ and the marginal likelihood $\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{model selection}}$.
 - Usual approaches: LA, VI, EP or MCMC.
- The observations can also be **warped**, either by a parametric function or by directly considering a **second GP prior** for the **inverse warping function**.
- By considering **multiple latent functions**, we can obtain flexible **heteroscedastic** likelihoods.

Table of Contents

① Part I

- Fundamentals of Bayesian inference
- Bayesian linear regression
- The Gaussian process

② Part II

- Learning the kernel and its hyperparameters
- Beyond Gaussian likelihood
- Sparse approximations

③ Part III

- GPs for Bayesian optimization
- Current trends on kernel design
- From GPLVM to Deep GPs
- Concluding remarks and additional topics

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

$$\underbrace{\log p(\mathbf{y} | \mathbf{X})}_{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$
$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N}(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2).$$

- Naive computation scales with $\mathcal{O}(N^3)$.
- Storage scales with $\mathcal{O}(N^2)$.
- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that summarizes the available data.

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

marg. likelihood

$$\overbrace{\log p(\mathbf{y} | \mathbf{X})}^{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$

$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N}(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2).$$

→ Naive computation scales with $\mathcal{O}(N^3)$.

→ Storage scales with $\mathcal{O}(N^2)$.

- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that summarizes the available data.

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

marg. likelihood

$$\overbrace{\log p(\mathbf{y} | \mathbf{X})}^{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$

$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N}(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2).$$

→ Naive computation scales with $\mathcal{O}(N^3)$.

→ Storage scales with $\mathcal{O}(N^2)$.

- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that summarizes the available data.

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

marg. likelihood

$$\overbrace{\log p(\mathbf{y} | \mathbf{X})}^{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$

$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N}(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2).$$

→ Naive computation scales with $\mathcal{O}(N^3)$.

→ Storage scales with $\mathcal{O}(N^2)$.

- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that **summarizes the available data**.

Sparse approximations

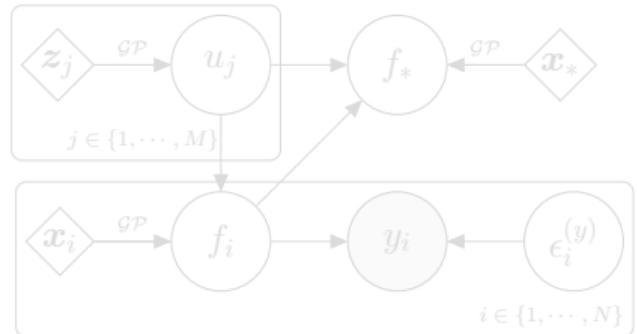
- Consider M **inducing variables** $\mathbf{u} \in \mathbb{R}^M$ from the same GP prior of \mathbf{f} .
→ \mathbf{u} corresponds to the evaluation of $f(\cdot)$ at M **pseudo-inputs** $\mathbf{z}_j|_{j=1}^M \in \mathbb{R}^D$.

marg. lik.

$$\overbrace{p(y|X)}^{\text{marg. lik.}} = \mathcal{N}(y|\mathbf{0}, K_f + \sigma_y^2 I),$$

$$\underbrace{p(f|X)}_{\text{prior}} = \mathcal{N}(f|\mathbf{0}, K_f), \quad [K_f]_{ii'} = k(x_i, x_{i'}),$$

$$\underbrace{p(u|Z)}_{\text{prior}} = \mathcal{N}(u|\mathbf{0}, K_u), \quad [K_u]_{jj'} = k(z_j, z_{j'}).$$



- If the knowledge from \mathbf{f} is concentrated in \mathbf{u} , the approximate posterior becomes

$$q(f_*) = \int \underbrace{p(f_*|u)}_{\approx p(f_*|u, f)} \underbrace{q(u)}_{\approx p(u|y)} du.$$

→ If $M = N$ and $z_i = x_i, \forall i$, we get $\mathbf{u} = \mathbf{f}$ and recover the original posterior exactly.

Sparse approximations

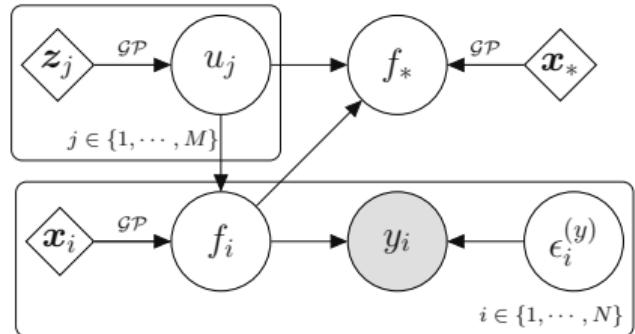
- Consider M inducing variables $\mathbf{u} \in \mathbb{R}^M$ from the same GP prior of \mathbf{f} .
→ \mathbf{u} corresponds to the evaluation of $f(\cdot)$ at M **pseudo-inputs** $\mathbf{z}_j|_{j=1}^M \in \mathbb{R}^D$.

marg. lik.

$$\overbrace{p(\mathbf{y}|X)}^{} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_f + \sigma_y^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{f}|X)}_{\text{prior}} = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad [\mathbf{K}_f]_{ii'} = k(\mathbf{x}_i, \mathbf{x}_{i'}),$$

$$\underbrace{p(\mathbf{u}|Z)}_{\text{prior}} = \mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{K}_u), \quad [\mathbf{K}_u]_{jj'} = k(\mathbf{z}_j, \mathbf{z}_{j'}).$$



- If the knowledge from \mathbf{f} is concentrated in \mathbf{u} , the approximate posterior becomes

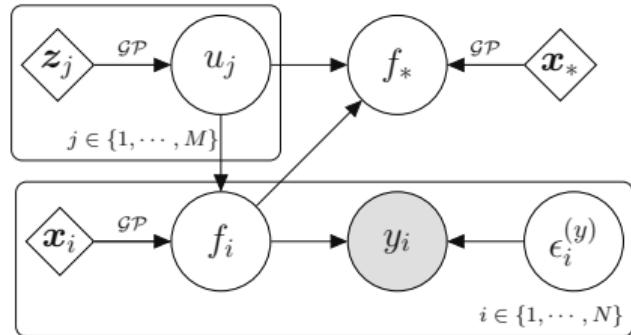
$$q(f_*) = \int \underbrace{p(f_*|u)}_{\approx p(f_*|u, f) \approx p(u|y)} \underbrace{q(u)}_{\approx p(u|y)} du.$$

→ If $M = N$ and $\mathbf{z}_i = \mathbf{x}_i, \forall i$, we get $\mathbf{u} = \mathbf{f}$ and recover the original posterior exactly.

Sparse approximations

- Consider M inducing variables $\mathbf{u} \in \mathbb{R}^M$ from the same GP prior of f .
 $\rightarrow \mathbf{u}$ corresponds to the evaluation of $f(\cdot)$ at M **pseudo-inputs** $\mathbf{z}_j |_{j=1}^M \in \mathbb{R}^D$.

$$\begin{aligned} \text{marg. lik.} \\ \underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{prior}} &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_f + \sigma_y^2 \mathbf{I}), \\ p(\mathbf{f}|\mathbf{X}) &= \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad [\mathbf{K}_f]_{ii'} = k(\mathbf{x}_i, \mathbf{x}_{i'}), \\ \underbrace{p(\mathbf{u}|\mathbf{Z})}_{\text{prior}} &= \mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{K}_u), \quad [\mathbf{K}_u]_{jj'} = k(\mathbf{z}_j, \mathbf{z}_{j'}). \end{aligned}$$



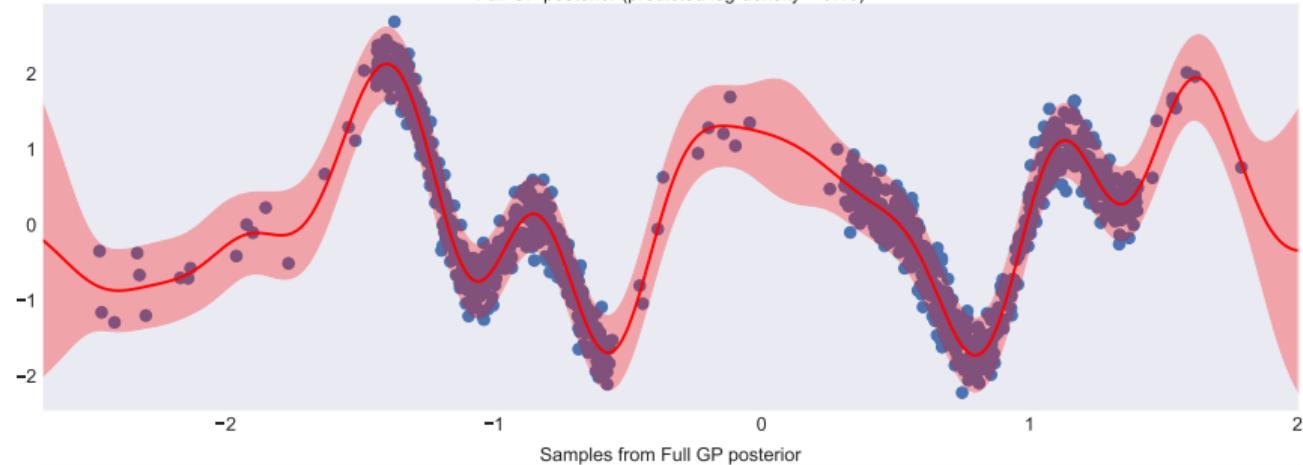
- If the knowledge from f is concentrated in u , the **approximate posterior** becomes

$$q(f_*) = \int \underbrace{p(f_* | \mathbf{u})}_{\approx p(f_* | \mathbf{u}, \mathbf{f})} \underbrace{q(\mathbf{u})}_{\approx p(\mathbf{u} | \mathbf{y})} d\mathbf{u}.$$

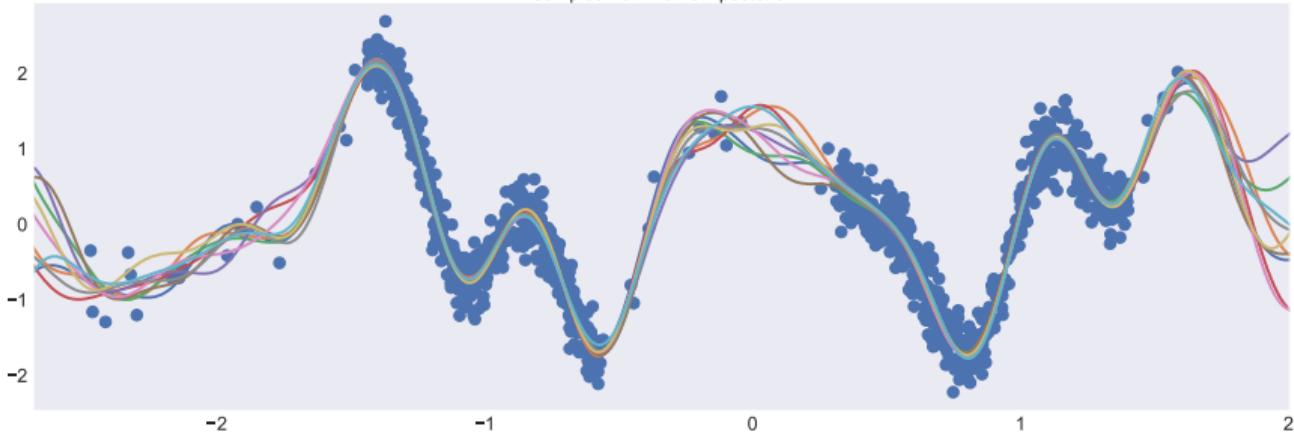
→ If $M = N$ and $\mathbf{z}_i = \mathbf{x}_i, \forall i$, we get $\mathbf{u} = \mathbf{f}$ and recover the original posterior exactly.

Full GP posterior (1083 observations)

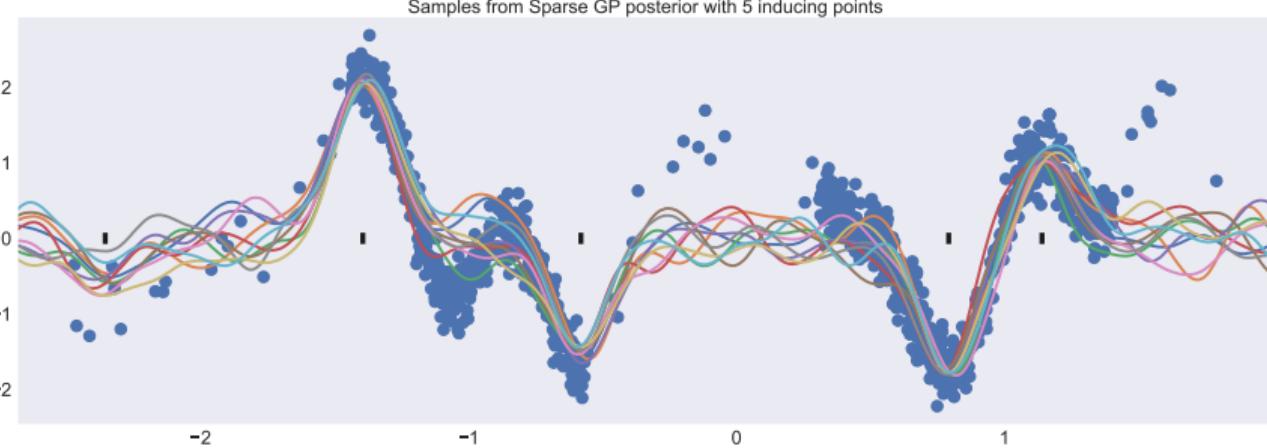
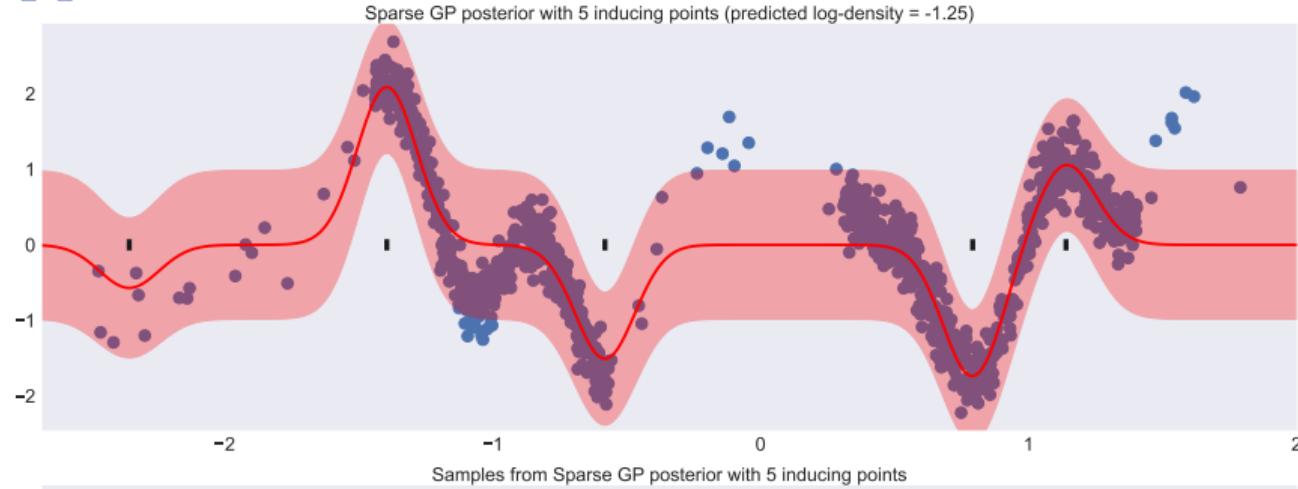
Full GP posterior (predicted log-density = 0.16)



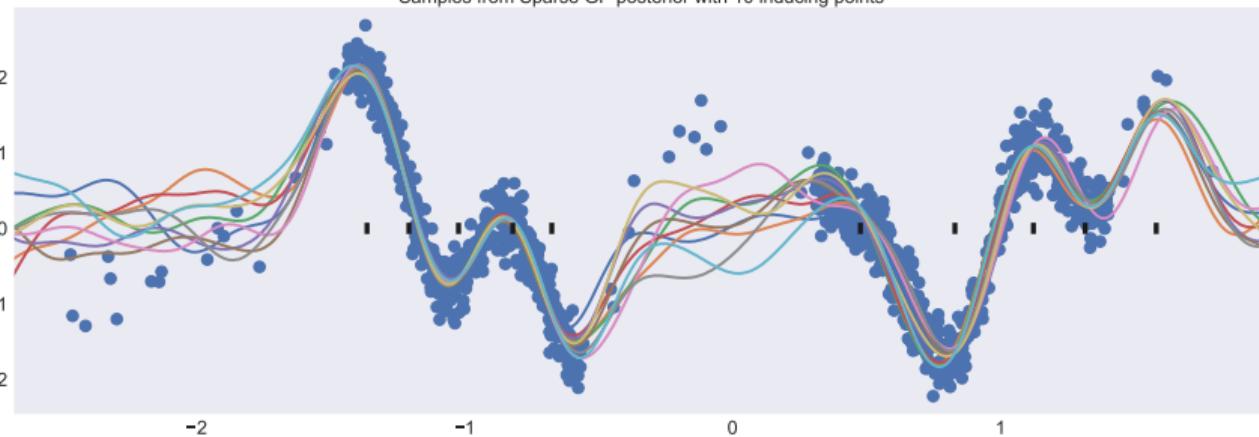
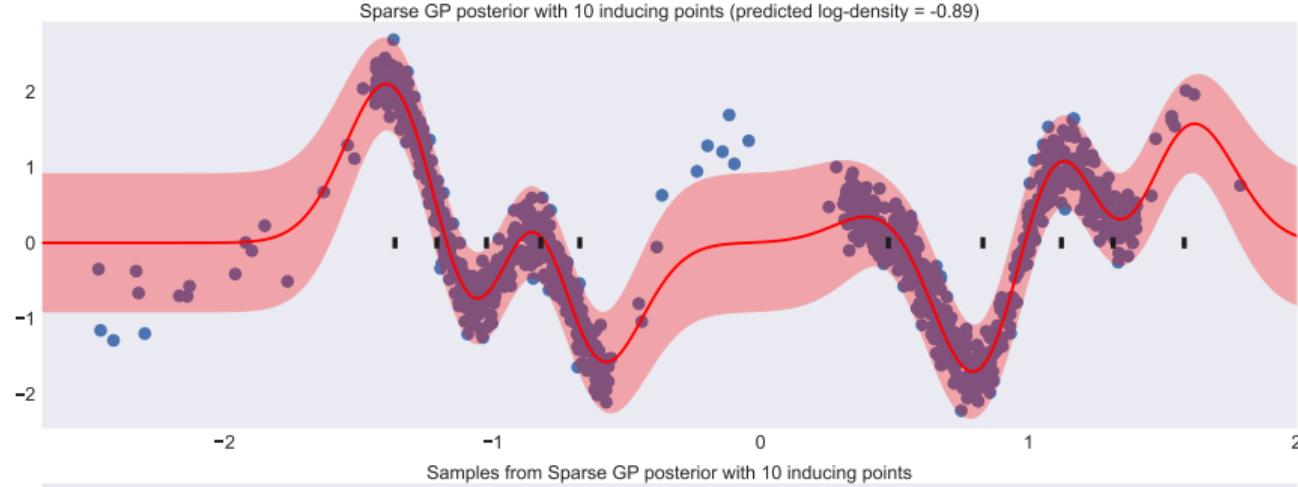
Samples from Full GP posterior



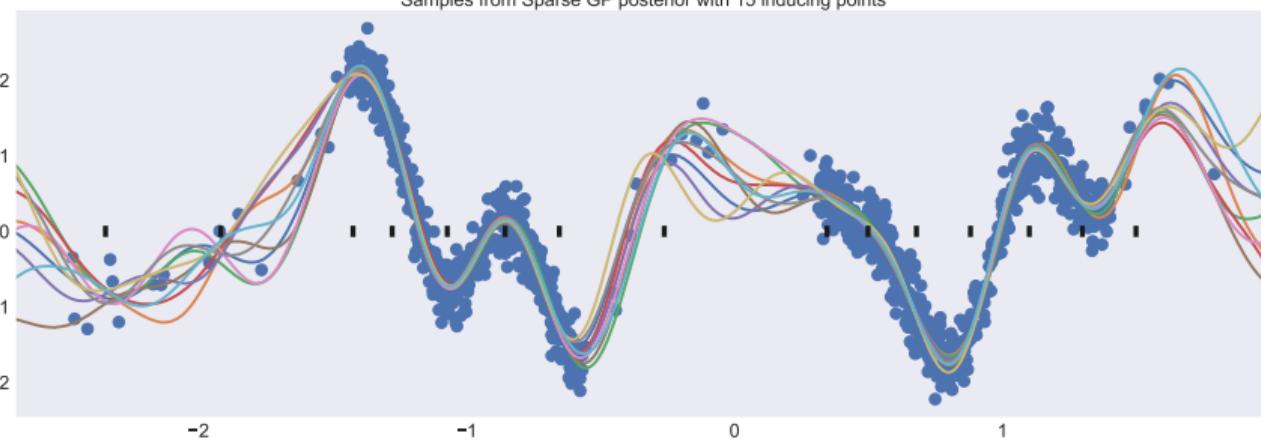
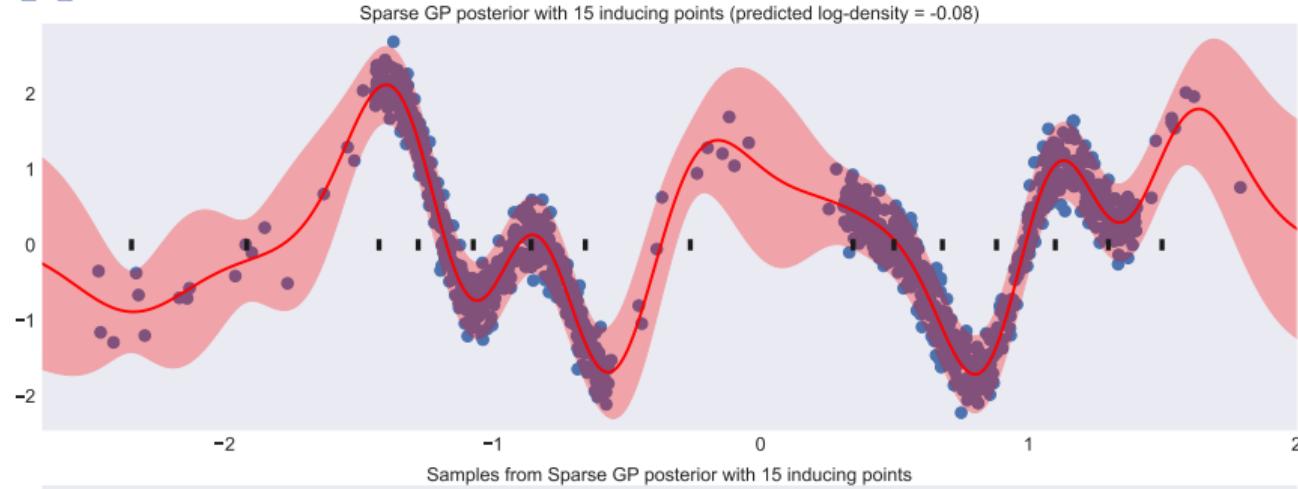
Sparse approximations (1083 observations, 5 inducing points)



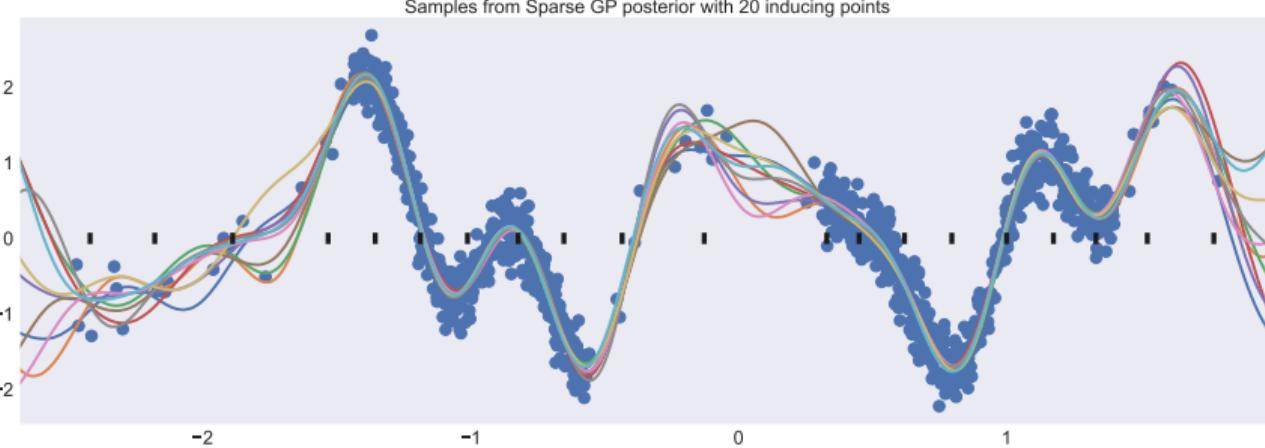
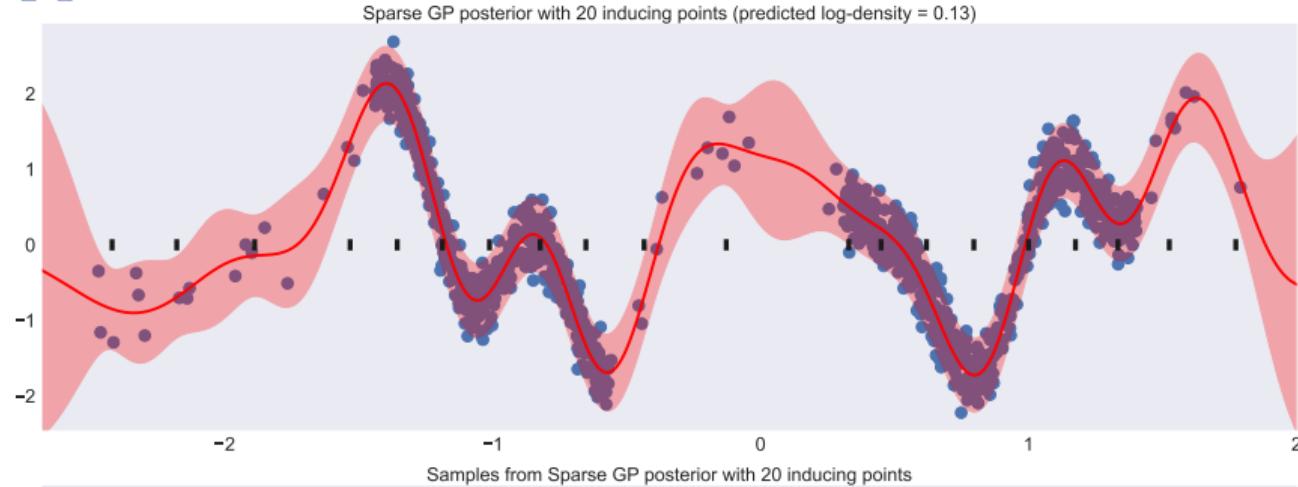
Sparse approximations (1083 observations, 10 inducing points)



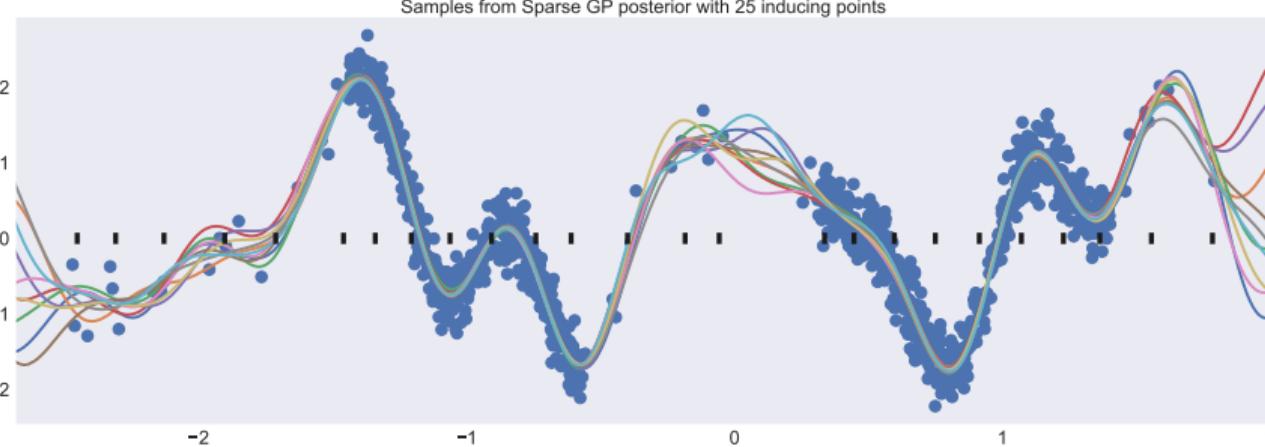
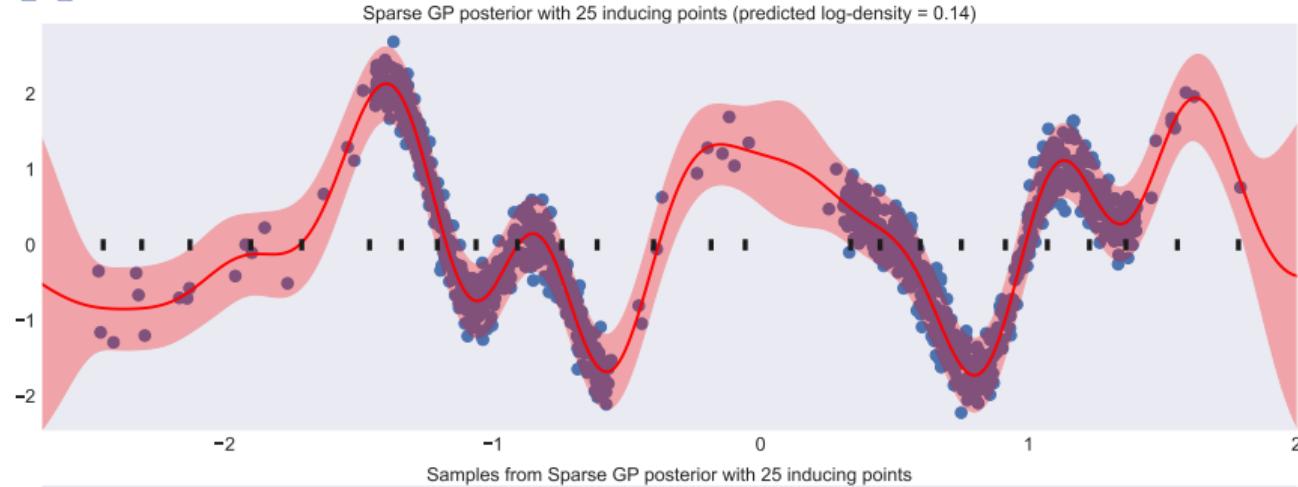
Sparse approximations (1083 observations, 15 inducing points)



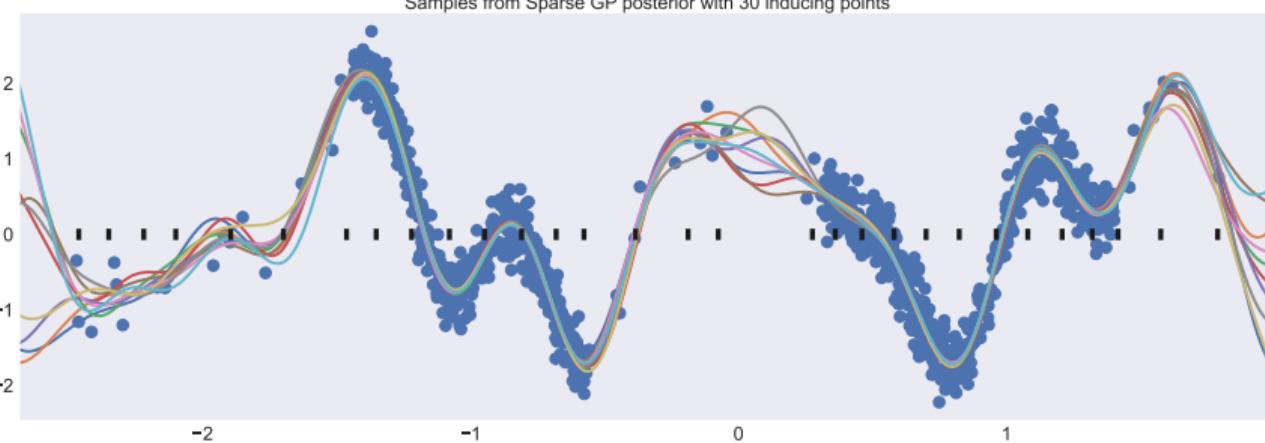
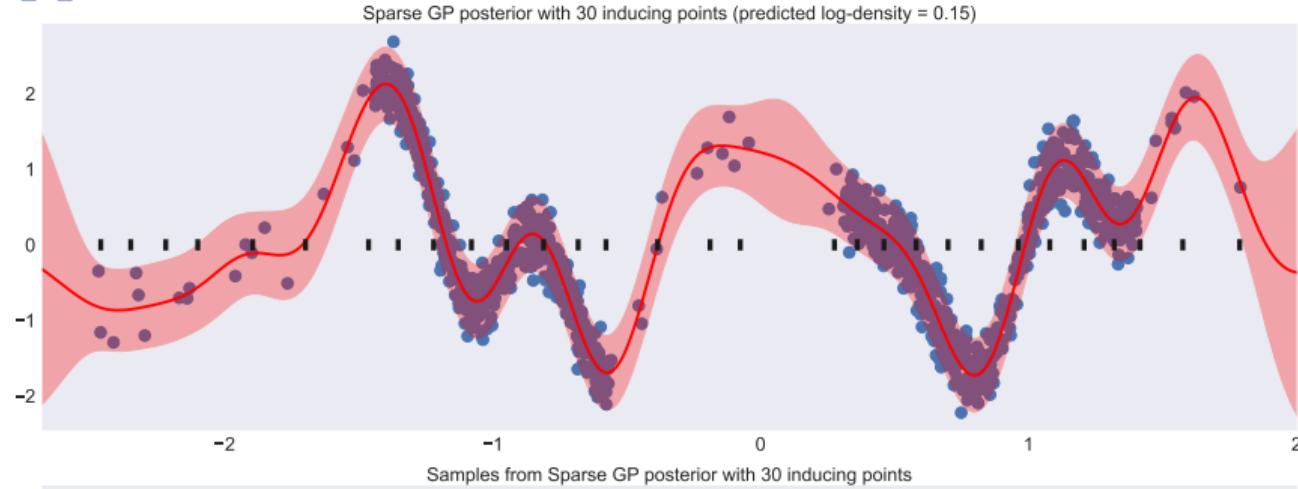
Sparse approximations (1083 observations, 20 inducing points)



Sparse approximations (1083 observations, 25 inducing points)

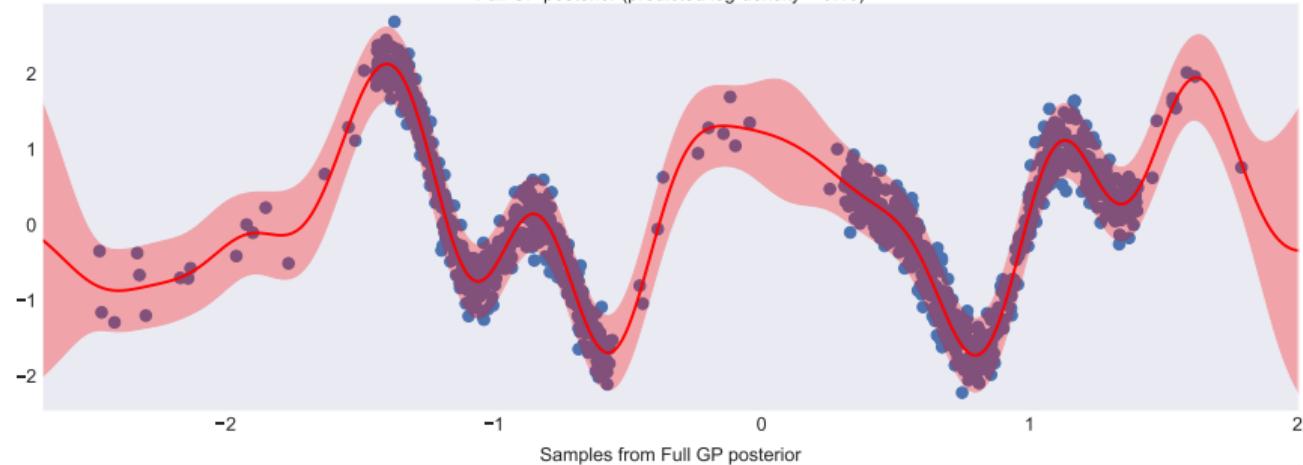


Sparse approximations (1083 observations, 30 inducing points)

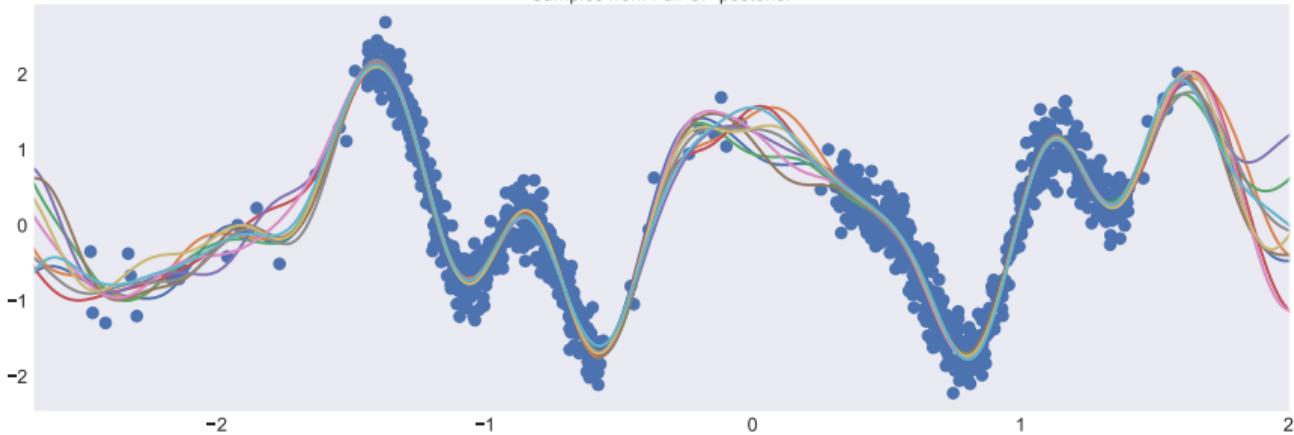


Full GP posterior (1083 observations)

Full GP posterior (predicted log-density = 0.16)



Samples from Full GP posterior



Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the Kullback-Leibler divergence:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- We replace $q(\mathbf{f}) = \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u}$ and (conveniently) choose $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{u}) \underbrace{p(\mathbf{f}|\mathbf{u})}_{\text{cond. prior}}$:

$$\begin{aligned}\underbrace{\log p(\mathbf{y})}_{\text{evidence}} &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})\cancel{p(\mathbf{f}|\mathbf{u})}p(\mathbf{u})}{q(\mathbf{u})\cancel{p(\mathbf{f}|\mathbf{u})}} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u})} \left[\log \frac{p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) = \underbrace{\mathcal{L}}_{\text{ELBO}}.\end{aligned}$$

- The kernel hyperparameters and the pseudo-inputs $\mathbf{z}_j|_1^M$ can be optimized by maximizing the **ELBO (evidence lower bound)** \mathcal{L} , improving the approximation.

Sparse approximations

- We replace $q(\mathbf{f}) = \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u}$ and (conveniently) choose $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{u}) \underbrace{p(\mathbf{f}|\mathbf{u})}_{\text{cond. prior}}$:

$$\begin{aligned}\underbrace{\log p(\mathbf{y})}_{\text{evidence}} &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u})} \left[\log \frac{p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) = \underbrace{\mathcal{L}}_{\text{ELBO}}.\end{aligned}$$

- The kernel hyperparameters and the pseudo-inputs $\mathbf{z}_j|_1^M$ can be optimized by maximizing the **ELBO (evidence lower bound)** \mathcal{L} , improving the approximation.

Sparse approximations

- We replace $q(\mathbf{f}) = \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u}$ and (conveniently) choose $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{u}) \underbrace{p(\mathbf{f}|\mathbf{u})}_{\text{cond. prior}}$:

$$\begin{aligned}\underbrace{\log p(\mathbf{y})}_{\text{evidence}} &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u})} \left[\log \frac{p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) = \underbrace{\mathcal{L}}_{\text{ELBO}}.\end{aligned}$$

- The kernel hyperparameters and the pseudo-inputs $\mathbf{z}_j|_1^M$ can be optimized by maximizing the **ELBO (evidence lower bound)** \mathcal{L} , improving the approximation.

Sparse approximations

- What we need to compute the ELBO \mathcal{L} ?

$$\log p(\mathbf{y}) \geq \mathcal{L} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})}[\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})).$$

- Since the joint $p(\mathbf{f}, \mathbf{u})$ is Gaussian, the conditional $p(\mathbf{f}|\mathbf{u})$ is also Gaussian and given by (with omitted dependence on \mathbf{X} and \mathbf{z}):

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_f & \mathbf{K}_{fu} \\ \mathbf{K}_{fu}^\top & \mathbf{K}_u \end{bmatrix}\right), \quad [\mathbf{K}_{fu}]_{ij} = k(\mathbf{x}_i, \mathbf{z}_j),$$
$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{u}, \mathbf{K}_f - \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{K}_{fu}^\top).$$

- How to define $q(\mathbf{u})$?

Sparse approximations

- What we need to compute the ELBO \mathcal{L} ?

$$\log p(\mathbf{y}) \geq \mathcal{L} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})}[\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})).$$

- Since the joint $p(\mathbf{f}, \mathbf{u})$ is Gaussian, the conditional $p(\mathbf{f}|\mathbf{u})$ is also Gaussian and given by (with omitted dependence on \mathbf{X} and \mathbf{z}):

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_f & \mathbf{K}_{fu} \\ \mathbf{K}_{fu}^\top & \mathbf{K}_u \end{bmatrix}\right), \quad [\mathbf{K}_{fu}]_{ij} = k(\mathbf{x}_i, \mathbf{z}_j),$$

$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{u}, \mathbf{K}_f - \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{K}_{fu}^\top).$$

- How to define $q(\mathbf{u})$?

Sparse approximations

- What we need to compute the ELBO \mathcal{L} ?

$$\log p(\mathbf{y}) \geq \mathcal{L} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})}[\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})).$$

- Since the joint $p(\mathbf{f}, \mathbf{u})$ is Gaussian, the conditional $p(\mathbf{f}|\mathbf{u})$ is also Gaussian and given by (with omitted dependence on \mathbf{X} and \mathbf{z}):

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_f & \mathbf{K}_{fu} \\ \mathbf{K}_{fu}^\top & \mathbf{K}_u \end{bmatrix}\right), \quad [\mathbf{K}_{fu}]_{ij} = k(\mathbf{x}_i, \mathbf{z}_j),$$

$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{u}, \mathbf{K}_f - \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{K}_{fu}^\top).$$

- How to define $q(\mathbf{u})$?

Sparse Variational GP

- The original sparse variational approach (Titsias, 2009) analytically optimizes $q(\mathbf{u})$:

$$q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \hat{\mathbf{m}}, \hat{\mathbf{S}}),$$

$$\hat{\mathbf{m}} = \frac{1}{\sigma_y^2} \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_{fu}^\top \mathbf{y},$$

$$\hat{\mathbf{S}} = \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_u.$$

- The result is a **collapsed bound** (considering a Gaussian likelihood):

$$\mathcal{L}_{\text{collapsed}} = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma_y^2 \mathbf{I} + \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top) - \frac{1}{2\sigma_y^2} \text{Tr}(\mathbf{K}_f - \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top).$$

→ Computation and storage scale with $\mathcal{O}(NM^2)$ and $\mathcal{O}(NM)$.

Sparse Variational GP

- The original sparse variational approach (Titsias, 2009) analytically optimizes $q(\mathbf{u})$:

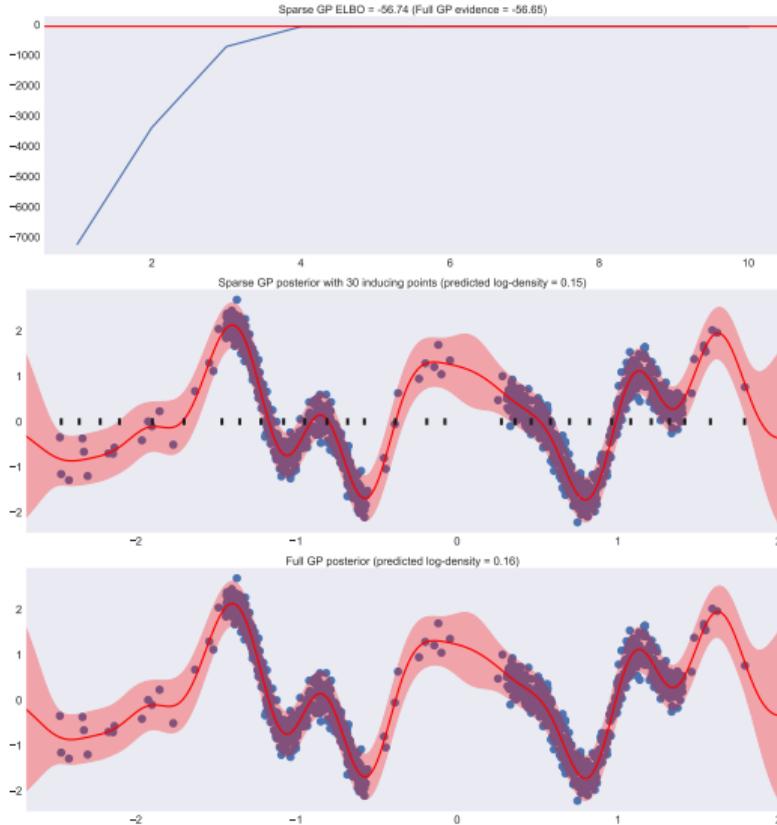
$$\begin{aligned} q(\mathbf{u}) &= \mathcal{N}(\mathbf{u} | \hat{\mathbf{m}}, \hat{\mathbf{S}}), \\ \hat{\mathbf{m}} &= \frac{1}{\sigma_y^2} \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_{fu}^\top \mathbf{y}, \\ \hat{\mathbf{S}} &= \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_u. \end{aligned}$$

- The result is a **collapsed bound** (considering a Gaussian likelihood):

$$\mathcal{L}_{\text{collapsed}} = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma_y^2 \mathbf{I} + \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top) - \frac{1}{2\sigma_y^2} \text{Tr}(\mathbf{K}_f - \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top).$$

→ Computation and storage scale with $\mathcal{O}(NM^2)$ and $\mathcal{O}(NM)$.

Sparse Variational GP



Observations

- The ELBO **monotonically increases**, which always **improves** the approximation.
- The gradients of the ELBO can be computed **analytically**.
- The exact gradients enable the use of **efficient optimization methods**, e.g. L-BFGS.

Stochastic Variational Inference for GPs

- SVI requires a set of **global parameters** and a **factorized ELBO**.
- SVGP explicitly parameterizes $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$ and obtains an **uncollapsed bound** that enables minibatch updates (Hensman *et al.*, 2013):

$$\mathcal{L}_{\text{uncollapsed}} = \sum_{i=1} \mathcal{L}_i - \text{KL}(q(\mathbf{u})||p(\mathbf{u})),$$

$$\text{where } \mathcal{L}_i = \log \mathcal{N}(y_i | \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{m}, \sigma_y^2) - \frac{1}{2\sigma_y^2} ([\mathbf{K}_f]_{ii} - \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{k}_i) - \frac{1}{2} \text{Tr} \left(\frac{1}{\sigma_y^2} \mathbf{S} \mathbf{K}_u^{-1} \mathbf{k}_i \mathbf{k}_i^\top \right)$$

$$\text{and } \mathbf{k}_i^\top = [\mathbf{K}_{fz}]_{i:}.$$

→ Given a minibatch size B , computation and storage scale with $\mathcal{O}(BM^2 + M^3)$ and $\mathcal{O}(BM + M^2)$.

Stochastic Variational Inference for GPs

- SVI requires a set of **global parameters** and a **factorized ELBO**.
- SVGP explicitly parameterizes $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$ and obtains an **uncollapsed bound** that enables minibatch updates (Hensman *et al.*, 2013):

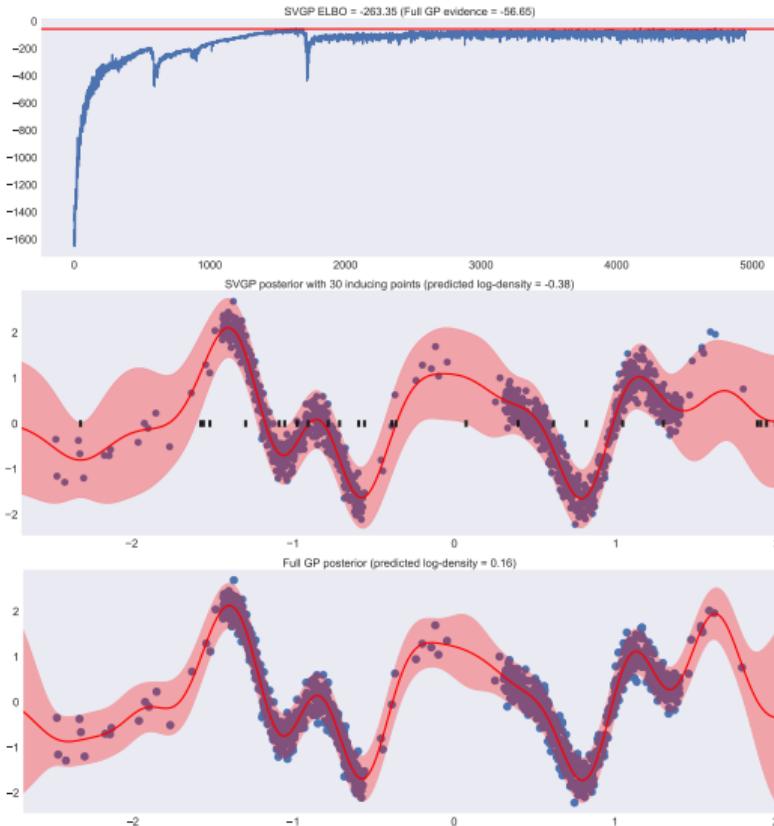
$$\mathcal{L}_{\text{uncollapsed}} = \sum_{i=1} \mathcal{L}_i - \text{KL}(q(\mathbf{u})||p(\mathbf{u})),$$

$$\text{where } \mathcal{L}_i = \log \mathcal{N}(y_i | \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{m}, \sigma_y^2) - \frac{1}{2\sigma_y^2} ([\mathbf{K}_f]_{ii} - \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{k}_i) - \frac{1}{2} \text{Tr} \left(\frac{1}{\sigma_y^2} \mathbf{S} \mathbf{K}_u^{-1} \mathbf{k}_i \mathbf{k}_i^\top \right)$$

$$\text{and } \mathbf{k}_i^\top = [\mathbf{K}_{fz}]_{i,:}.$$

→ Given a minibatch size B , computation and storage scale with $\mathcal{O}(BM^2 + M^3)$ and $\mathcal{O}(BM + M^2)$.

Stochastic Variation Inference for GPs



Observations

- The **ELBO does not increase monotonically** due to the use of minibatches.
- The **noisy** gradients of the ELBO can be computed **analytically**.
- More complex **stochastic optimization methods** may be used, e.g. Adam.

Sparse approximations

- In both cases, prediction is performed by integrating out the inducing variables \mathbf{u} using the **variational posterior** $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$:

$$\begin{aligned} q(f_*) &= \int p(f_* | \mathbf{u}) q(\mathbf{u}) d\mathbf{u} \\ &= \int \mathcal{N}(f_* | \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{u}, \mathbf{K}_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top) \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}) d\mathbf{u} \\ &= \mathcal{N}\left(f_* \mid \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{m}, \mathbf{K}_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} (\mathbf{K}_u - \mathbf{S}) \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top\right). \end{aligned}$$

- Recall that:
 - In the collapsed bound, the moments $\hat{\mathbf{m}}$ and $\hat{\mathbf{S}}$ are optimally obtained.
 - In the uncollapsed bound, the moments \mathbf{m} and \mathbf{S} are variational parameters updated iteratively using the ELBO.

Sparse approximations

- In both cases, prediction is performed by integrating out the inducing variables \mathbf{u} using the **variational posterior** $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$:

$$\begin{aligned} q(f_*) &= \int p(f_* | \mathbf{u}) q(\mathbf{u}) d\mathbf{u} \\ &= \int \mathcal{N}(f_* | \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{u}, \mathbf{K}_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top) \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}) d\mathbf{u} \\ &= \mathcal{N}\left(f_* \mid \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{m}, \mathbf{K}_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} (\mathbf{K}_u - \mathbf{S}) \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top\right). \end{aligned}$$

- Recall that:
 - In the **collapsed bound**, the moments $\hat{\mathbf{m}}$ and $\hat{\mathbf{S}}$ are **optimally obtained**.
 - In the **uncollapsed bound**, the moments \mathbf{m} and \mathbf{S} are variational parameters **updated iteratively** using the ELBO.

Sparse approximations

Summary

- Sparse models enable scaling GP models to **large datasets**.
 - We considered **inducing variables** and corresponding **pseudo-inputs**.
 - integrated out
 - optimized
 - A **variational approach** approximates the posterior of the inducing variables and obtains a lower bound to the model evidence, the **ELBO**.
 - Kernel hyperparameters and variational parameters are **optimized** using the ELBO.
 - Two ELBO's: **collapsed** (Titsias, 2009) and **uncollapsed** (Hensman *et al.*, 2013).
 - Predictions are performed using the **variational posterior** of the inducing points.

Table of Contents

① Part I

Fundamentals of Bayesian inference

Bayesian linear regression

The Gaussian process

② Part II

Learning the kernel and its hyperparameters

Beyond Gaussian likelihood

Sparse approximations

③ Part III

GPs for Bayesian optimization

Current trends on kernel design

From GPLVM to Deep GPs

Concluding remarks and additional topics

Global Bayesian optimization with GPs

- Consider the task of optimizing a function $f(\cdot)$ **whose expression is unknown**.
- We are able to evaluate $f(\cdot)$ in a **noisy** way by observing $y_i = f(\mathbf{x}_i) + \epsilon_i$ for any given input \mathbf{x}_i and a random noise ϵ_i .
- Some examples:
 - Locate the best position to extract some ore;
 - Minimize the resource consumption in an industrial process;
 - Configure a laboratory experiment;
 - Maximize the access speed in a database;
 - Select hyperparameters for a machine learning model.
- If the analytical form of $f(\cdot)$ is unknown, how will we maximize or minimize it?

Global Bayesian optimization with GPs

- Consider the task of optimizing a function $f(\cdot)$ **whose expression is unknown**.
- We are able to evaluate $f(\cdot)$ in a **noisy** way by observing $y_i = f(\mathbf{x}_i) + \epsilon_i$ for any given input \mathbf{x}_i and a random noise ϵ_i .
- Some examples:
 - Locate the best position to extract some ore;
 - Minimize the resource consumption in an industrial process;
 - Configure a laboratory experiment;
 - Maximize the access speed in a database;
 - Select hyperparameters for a machine learning model.
- If the analytical form of $f(\cdot)$ is unknown, how will we maximize or minimize it?

Global Bayesian optimization with GPs

- Consider the task of optimizing a function $f(\cdot)$ **whose expression is unknown**.
- We are able to evaluate $f(\cdot)$ in a **noisy** way by observing $y_i = f(\mathbf{x}_i) + \epsilon_i$ for any given input \mathbf{x}_i and a random noise ϵ_i .
- Some examples:
 - Locate the best position to extract some ore;
 - Minimize the resource consumption in an industrial process;
 - Configure a laboratory experiment;
 - Maximize the access speed in a database;
 - Select hyperparameters for a machine learning model.
- If the analytical form of $f(\cdot)$ is unknown, how will we maximize or minimize it?

Global Bayesian optimization with GPs

- Consider N experiments that result in N pairs $(\mathbf{x}_i, y_i)|_{i=1}^N$, where y_i is an observation added with Gaussian noise.
- We choose a GP prior with kernel function $k(\cdot, \cdot)$ for the function $f(\cdot)$ that we want to optimize:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}),$$

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}),$$

where $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$.

- As we have seen before, this model is **conjugate** and all the distributions of interest are **analytical**.

Global Bayesian optimization with GPs

- Predictions in this model are well defined Gaussian distributions:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}(f_* | \mu_*, \sigma_*^2), \\ \mu_* &= \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \\ \sigma_*^2 &= k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}. \end{aligned}$$

- The predictive distribution informs the **uncertainty** related to the evaluation of the unknown function in the input \mathbf{x}_* .
- In a global minimization task, for an input search domain \mathcal{X} , we seek:
 - Regions with lower predictive mean values μ_* (*exploitation*);
 - Regions with high predictive uncertainty σ_*^2 (*exploration*).
- We assume that good solutions are **close to each other** in the input space.

Global Bayesian optimization with GPs

- Predictions in this model are well defined Gaussian distributions:

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}(f_* | \mu_*, \sigma_*^2), \\ \mu_* &= \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, \\ \sigma_*^2 &= k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*}. \end{aligned}$$

- The predictive distribution informs the **uncertainty** related to the evaluation of the unknown function in the input \mathbf{x}_* .
- In a global minimization task, for an input search domain \mathcal{X} , we seek:
 - Regions with lower predictive mean values μ_* (*exploitation*);
 - Regions with high predictive uncertainty σ_*^2 (*exploration*).
- We assume that good solutions are **close to each other** in the input space.

Global Bayesian optimization with GPs

- How do we know if a given point is located in a region of interest?
- We want to use all the information provided by past evaluations when querying the next candidate.
- Idea: We use an **acquisition function** to quantify the relevance of a given candidate.
- The acquisition function $a : \mathcal{X} \rightarrow \mathbb{R}^+$ is chosen to be easily computed.
- The next point to be evaluated is given by

$$x_* = \arg \max a(x).$$

- Examples of usual acquisition functions:
 - Probability of Improvement;
 - Expected Improvement;
 - GP Lower Confidence Bound.

Global Bayesian optimization with GPs

- How do we know if a given point is located in a region of interest?
- We want to use all the information provided by past evaluations when querying the next candidate.
- **Idea:** We use an **acquisition function** to **quantify the relevance** of a given candidate.
- The acquisition function $a : \mathcal{X} \rightarrow \mathbb{R}^+$ is chosen to be easily computed.
- The next point to be evaluated is given by

$$x_* = \arg \max a(x).$$

- Examples of usual acquisition functions:
 - Probability of Improvement;
 - Expected Improvement;
 - GP Lower Confidence Bound.

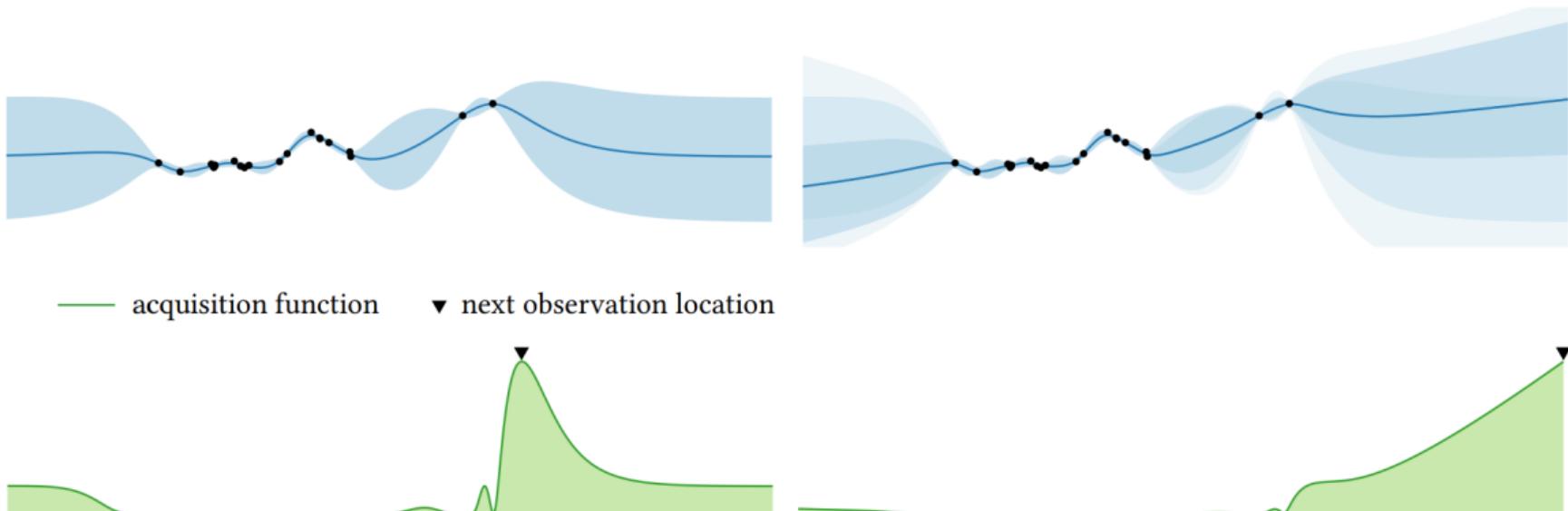
Global Bayesian optimization with GPs

- How do we know if a given point is located in a region of interest?
- We want to use all the information provided by past evaluations when querying the next candidate.
- **Idea:** We use an **acquisition function** to **quantify the relevance** of a given candidate.
- The acquisition function $a : \mathcal{X} \rightarrow \mathbb{R}^+$ is chosen to be easily computed.
- The next point to be evaluated is given by

$$\boldsymbol{x}_* = \arg \max a(\boldsymbol{x}).$$

- Examples of usual acquisition functions:
 - Probability of Improvement;
 - Expected Improvement;
 - GP Lower Confidence Bound.

Illustration of Bayesian optimization steps



BO for a maximization task. Taken from Garnett (2023).

Acquisition functions - Probability of Improvement

- Aims to maximize the probability of finding points better than the best current solution $\mu^- = \min_i \mu(\mathbf{x}_i)$, where $\mu(\mathbf{x}_i)$ denotes the mean GP prediction in \mathbf{x}_i .
- Focus in *exploitation*, with balance determined by the hyperparameter $\xi \geq 0$ (usual value $\xi = 0.01$):

$$\text{PI}(\mathbf{x}) = P(f(\mathbf{x}) \leq \mu^- - \xi) = \Phi\left(\frac{\mu^- - \xi - \mu(\mathbf{x})}{\sigma(\mathbf{x})}\right),$$

where $\sigma(\mathbf{x}_i)$ denotes the GP predictive standard deviation in \mathbf{x}_i and $\Phi(\cdot)$ is the cumulative distribution function of the normalized Gaussian:

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{1}{2}t^2\right) dt.$$

Acquisition functions - Expected Improvement

- Aims to maximize the magnitude of the expected progress when compared to the best current solution μ^- .
- Considering the hyperparameter $\xi \geq 0$ (usual value $\xi = 0.01$):

$$\text{EI}(\mathbf{x}) = \begin{cases} 0. & \text{if } \sigma(\mathbf{x}) = 0; \\ \tau(\mathbf{x})\Phi\left(\frac{\tau(\mathbf{x})}{\sigma(\mathbf{x})}\right) + \sigma(\mathbf{x})\phi\left(\frac{\tau(\mathbf{x})}{\sigma(\mathbf{x})}\right), & \text{if } \sigma(\mathbf{x}) > 0, \end{cases}$$

where $\tau(\mathbf{x}) = \mu^- - \xi - \mu(\mathbf{x})$.

Above, $\Phi(\cdot)$ and $\phi(\cdot)$ are respectively the cumulative and density functions of the normalized Gaussian.

Acquisition functions - GP Lower Confidence Bound

- Uses the GP confidence intervals to choose the next query point.
- Considering the hyperparameter $\kappa > 0$:

$$\text{LCB}(\boldsymbol{x}) = \mu(\boldsymbol{x}) - \kappa\sigma(\boldsymbol{x}),$$

where $\kappa = \sqrt{\nu\beta_t}$, $\beta_t = 2\log(t^{D/2+2}\pi^2/3\delta)$, t is the current optimization iteration, and $\nu, \delta > 0$.
Usual values: $\delta = 0.1$ and $\nu = 0.2$.

Global Bayesian optimization with GPs

Summary of the algorithm

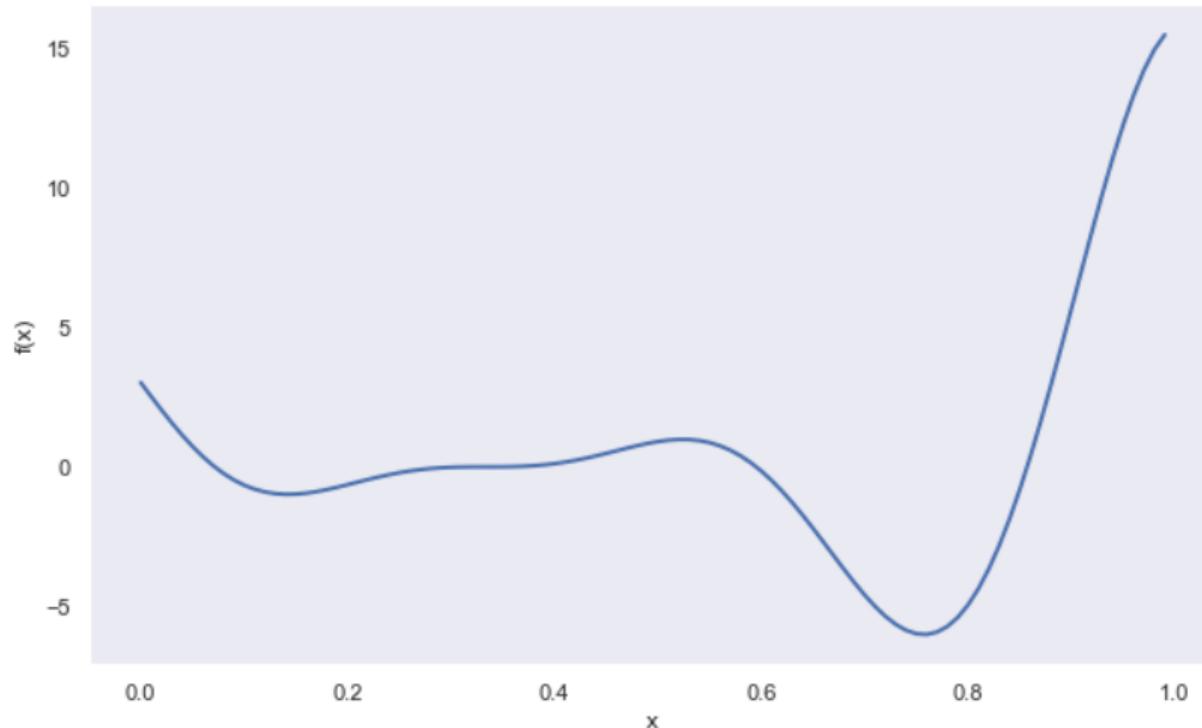
- ① Define the initial dataset $\mathcal{D}_0 = (\mathbf{x}_i, y_i) |_{i=1}^{N_0}$;
- ② Define the acquisition function $a(\mathbf{x})$;
- ③ Initialize the GP model using \mathcal{D}_0 ;
- ④ Repeat until convergence or until the maximum number of iterations is reached:
 - ① Find the next query point by maximizing the acquisition function:

$$\mathbf{x}_t = \arg \max_{\mathbf{x}} a(\mathbf{x} | \mathcal{D}_{t-1});$$

- ② Evaluate $y_t = f(\mathbf{x}_t) + \epsilon_t$;
- ③ Augment the dataset $\mathcal{D}_t = \{\mathcal{D}_{t-1}, (\mathbf{x}_t, y_t)\}$;
- ④ Update the GP model with the new dataset;
- ⑤ Return the point \mathbf{x}_t with the best corresponding y_t .

Global Bayesian optimization with GPs

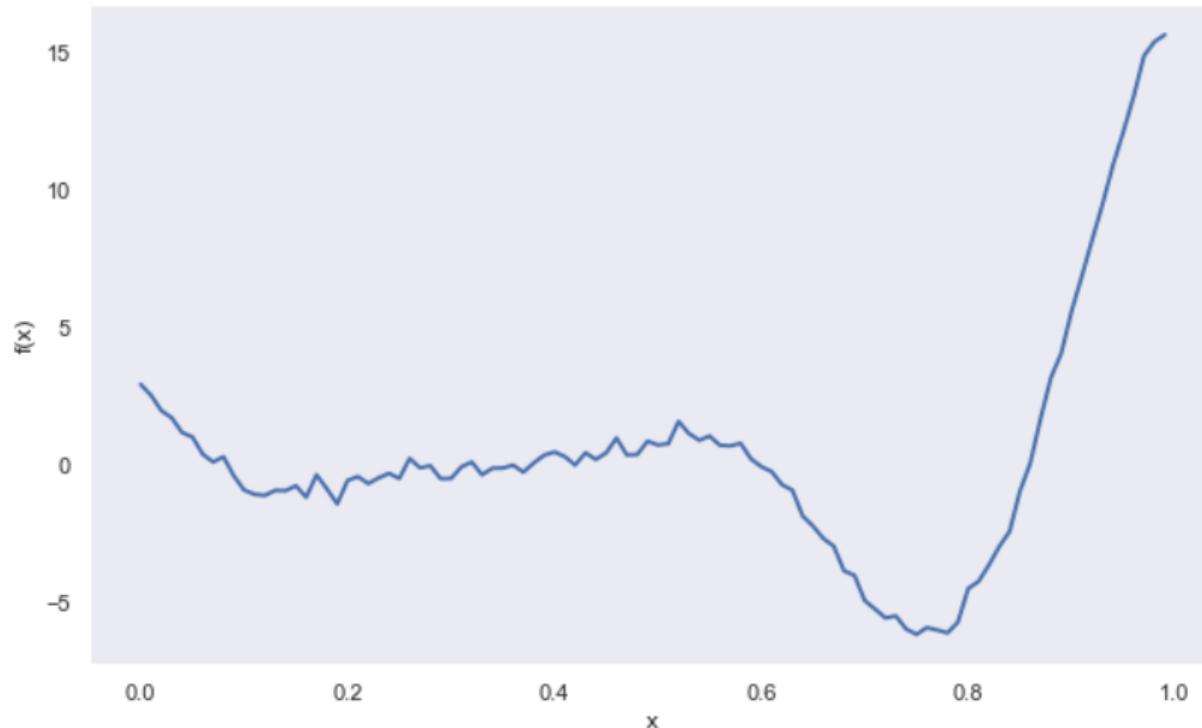
1D example



$$f(x) = (6x - 2)^2 \sin(12x - 4)$$

Global Bayesian optimization with GPs

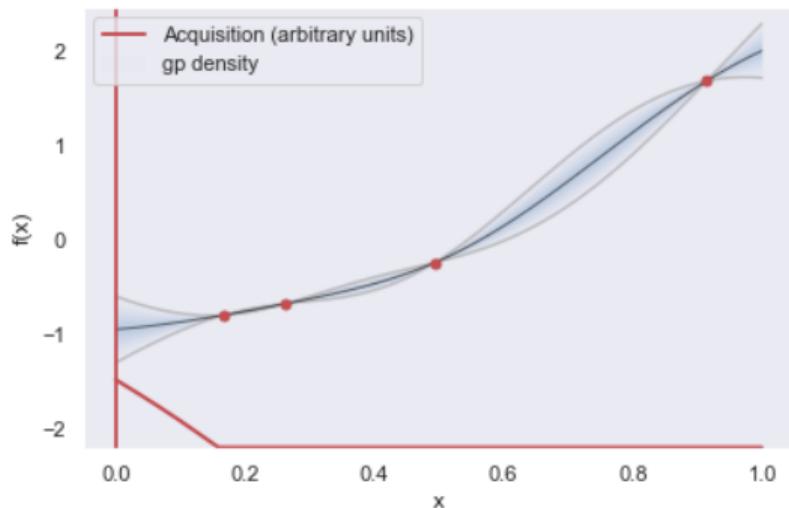
1D example



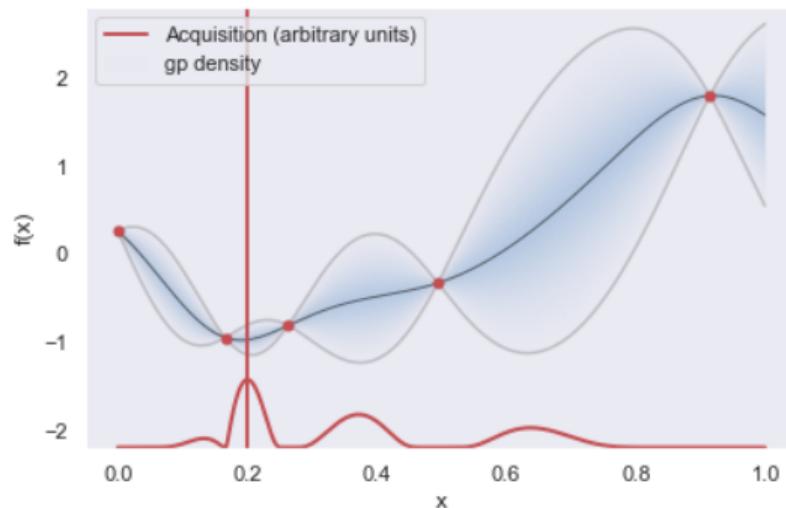
$$f(x) = (6x - 2)^2 \sin(12x - 4) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 0.25^2)$$

Global Bayesian optimization with GPs

1D example



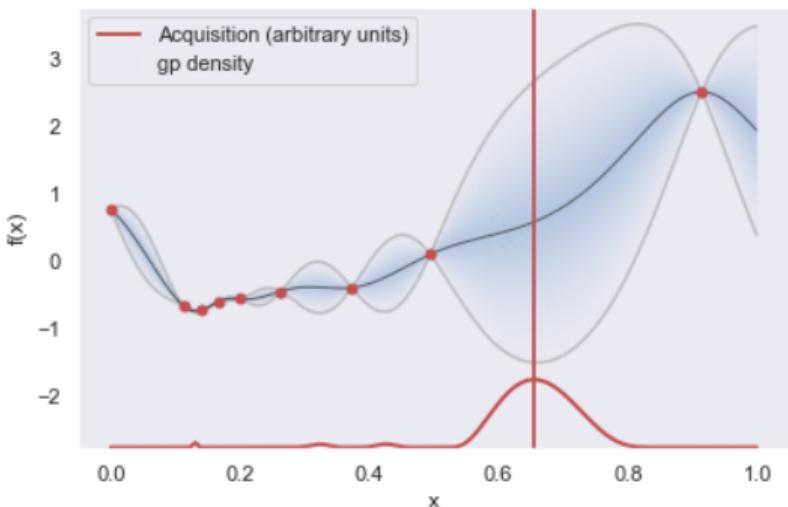
4 initial points



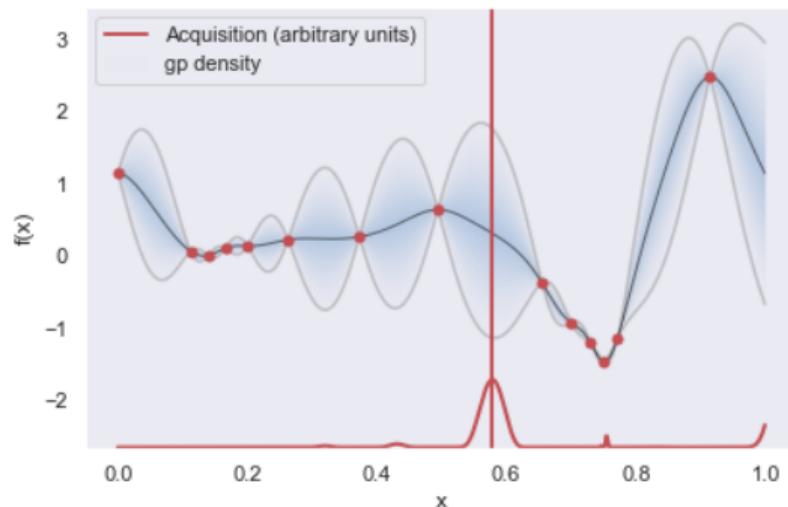
After 1 iteration

Global Bayesian optimization with GPs

1D example



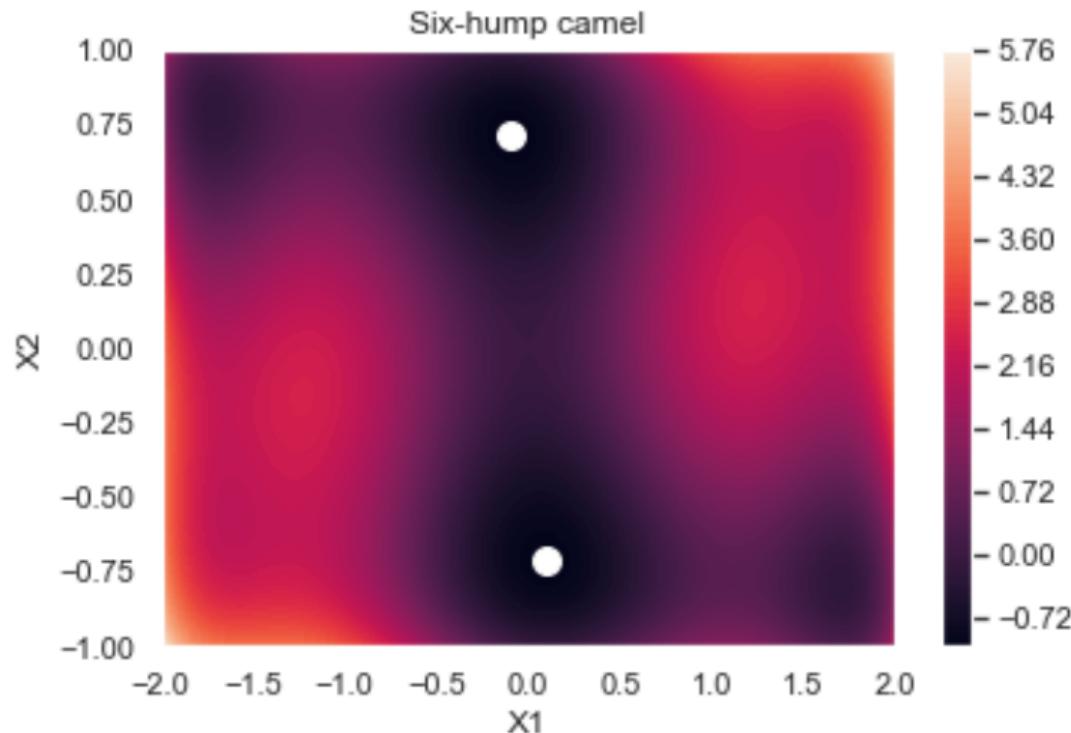
After 5 iterations



After 10 iterations

Global Bayesian optimization with GPs

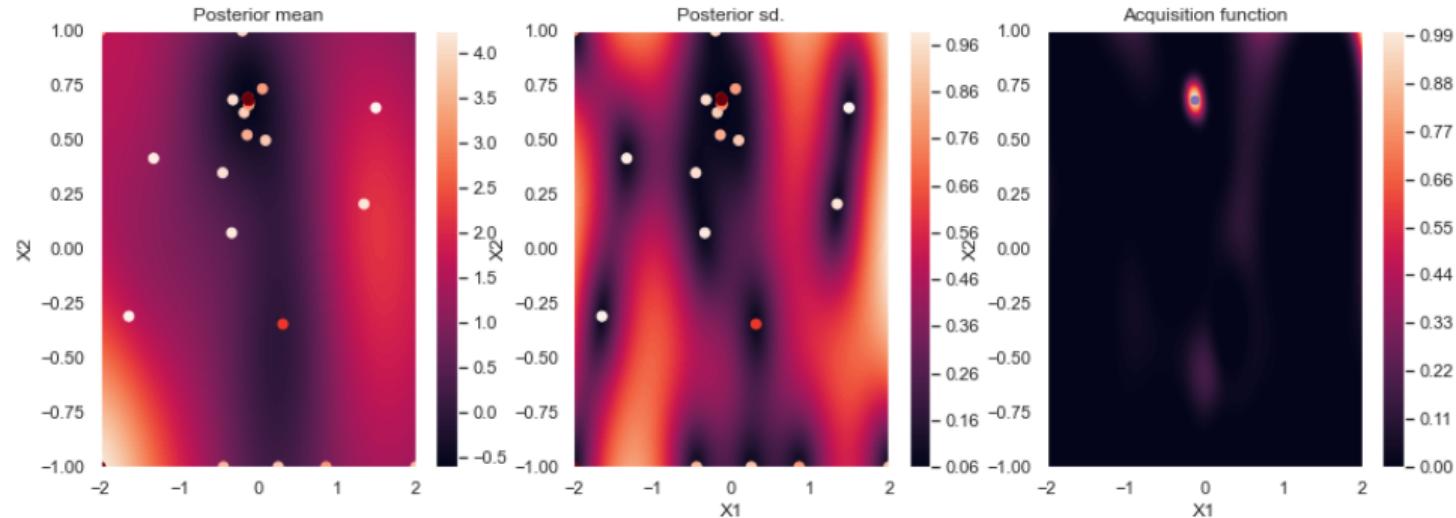
2D example



$$f(x_1, x_2) = \left(4 - 2.1x_1^2 + \frac{x_1^4}{3}\right)x_1^2 + x_1x_2 + (-4 + 4x_2^2)x_2^2$$

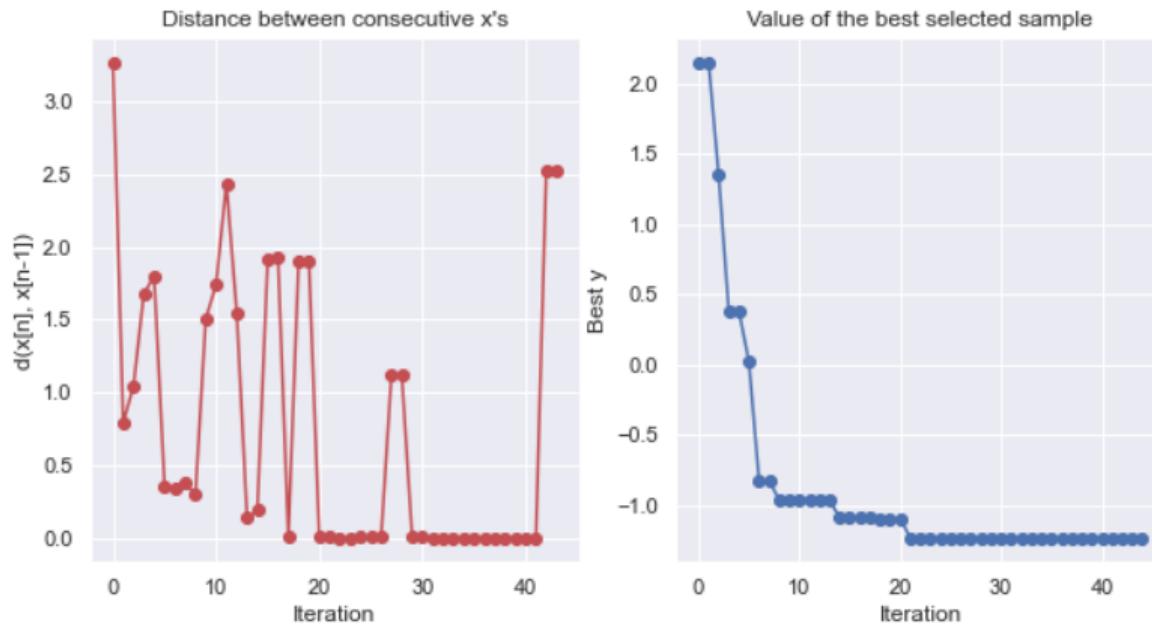
Global Bayesian optimization with GPs

2D example



Global Bayesian optimization with GPs

2D example

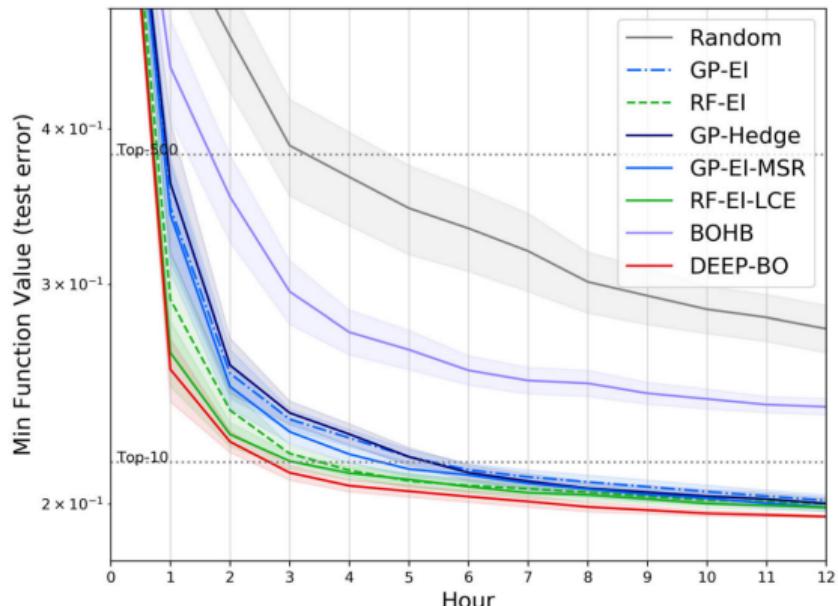


Global Bayesian optimization with GPs

Hyperparameter optimization (HPO) example

HPO of a CNN for the CIFAR-10 dataset.

Hyperparameter	Range
number of first convolution kernels	[8, 32]
number of second convolution kernels	[32, 64]
number of third convolution kernels	[64, 128]
number of forth convolution kernels	[64, 128]
number of neurons in a fully connected layer	[10, 1000]
square length of convolution kernel	[2, 3]
learning rate	[0.0001, 0.4] (log scale)
L2 regularization factor	[0.0, 1.0]
activation function type	ReLU, tanh, eLU
regularization method	None, Dropout, BatchNorm



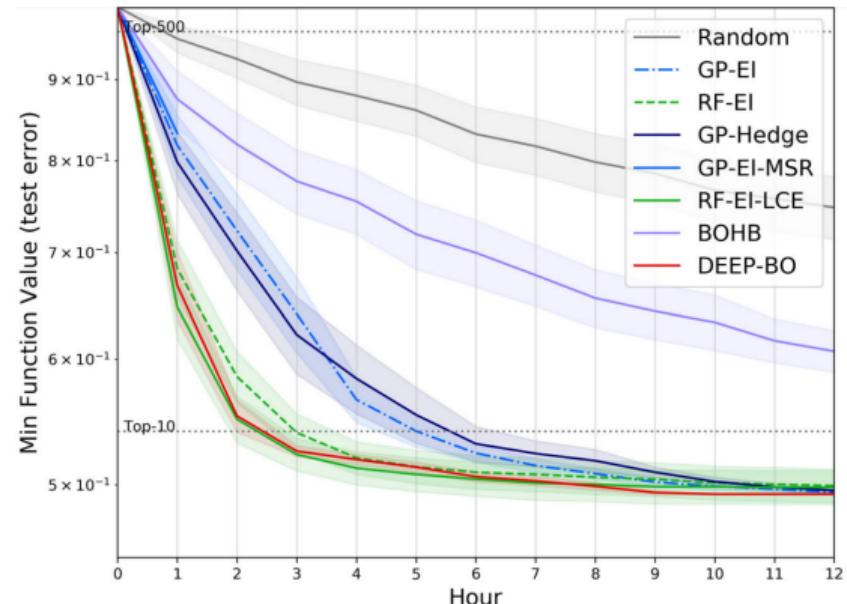
Taken from Chol *et al.* (2020).

Global Bayesian optimization with GPs

Hyperparameter optimization (HPO) example

HPO of a CNN for the CIFAR-100 dataset.

Hyperparameter	Range
number of first convolution kernels	[8, 32]
number of second convolution kernels	[32, 64]
number of third convolution kernels	[64, 128]
number of forth convolution kernels	[64, 128]
number of neurons in a fully connected layer	[10, 1000]
square length of convolution kernel	[2, 3]
learning rate	[0.0001, 0.4] (log scale)
L2 regularization factor	[0.0, 1.0]
activation function type	ReLU, tanh, sigmoid, eLU
regularization method	None, Dropout, BatchNorm



Taken from Chol *et al.* (2020).

Global Bayesian optimization with GPs

Practical aspects

- Use a kernel function less smooth than the RBF, such as the Matérn $\frac{5}{2}$ or $\frac{3}{2}$.
- Use efficient methods to optimize the acquisition function (search trees, hill-climbing, Monte Carlo, etc).
- Initialize the dataset \mathcal{D}_0 by using samples from the Sobol sequence or latin hypercube.
- If a priori knowledge about the task is available, choose informative initial points.
- Use a well established package: GPyOpt, GPflowOpt, Spearmint, BoTorch, Dragonfly, pybo...

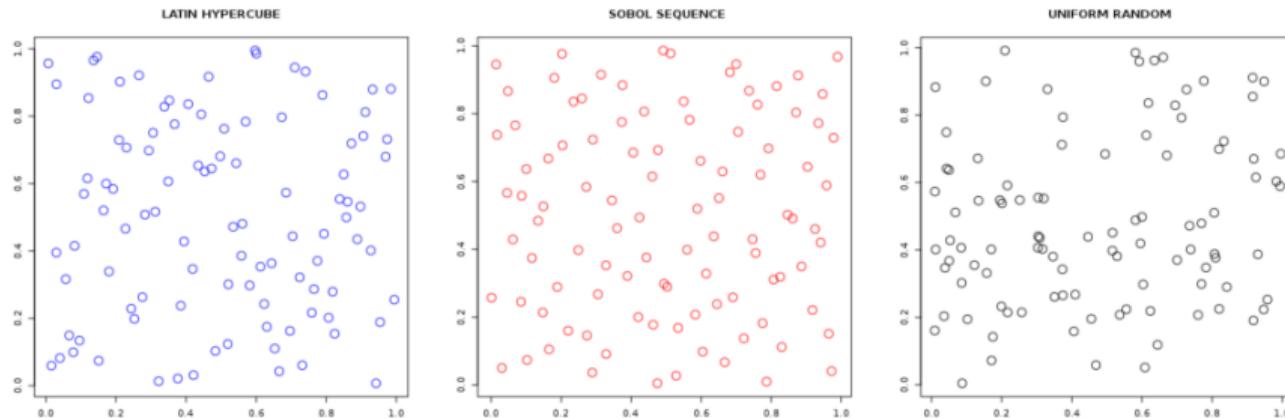


Table of Contents

① Part I

Fundamentals of Bayesian inference

Bayesian linear regression

The Gaussian process

② Part II

Learning the kernel and its hyperparameters

Beyond Gaussian likelihood

Sparse approximations

③ Part III

GPs for Bayesian optimization

Current trends on kernel design

From GPLVM to Deep GPs

Concluding remarks and additional topics

Current trends on kernel design

- The properties of a zero mean GP model come from its **kernel function**:

$$y_i = f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \quad i \in \{1, \dots, N\},$$

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$

$$[\mathbf{K}_f]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j).$$

- Functions that generate a **positive semidefinite** matrix \mathbf{K}_f are valid kernels.
- We have seen some of the most used kernels: Squared Exponential, Periodic, Matérn, Polynomial, etc.
- How to find suitable kernel functions for a given problem?

Current trends on kernel design

- The properties of a zero mean GP model come from its **kernel function**:

$$\begin{aligned}y_i &= f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \quad i \in \{1, \dots, N\}, \\ \mathbf{f} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f), \\ [\mathbf{K}_f]_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j).\end{aligned}$$

- Functions that generate a **positive semidefinite** matrix \mathbf{K}_f are valid kernels.
- We have seen some of the most used kernels: Squared Exponential, Periodic, Matérn, Polynomial, etc.
- How to find suitable kernel functions for a given problem?

Current trends on kernel design

- The properties of a zero mean GP model come from its **kernel function**:

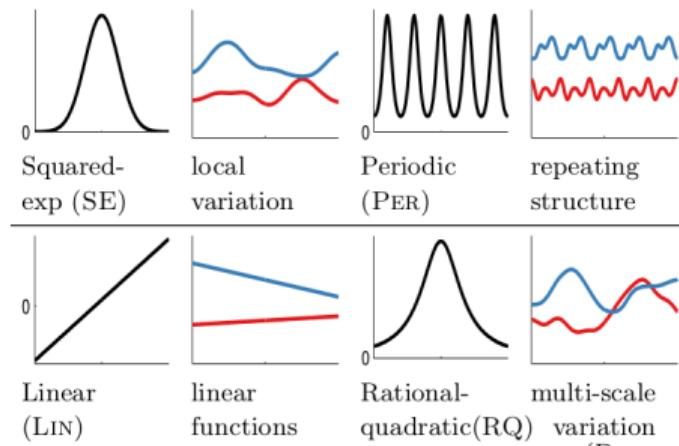
$$y_i = f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \quad i \in \{1, \dots, N\},$$

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$

$$[\mathbf{K}_f]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j).$$

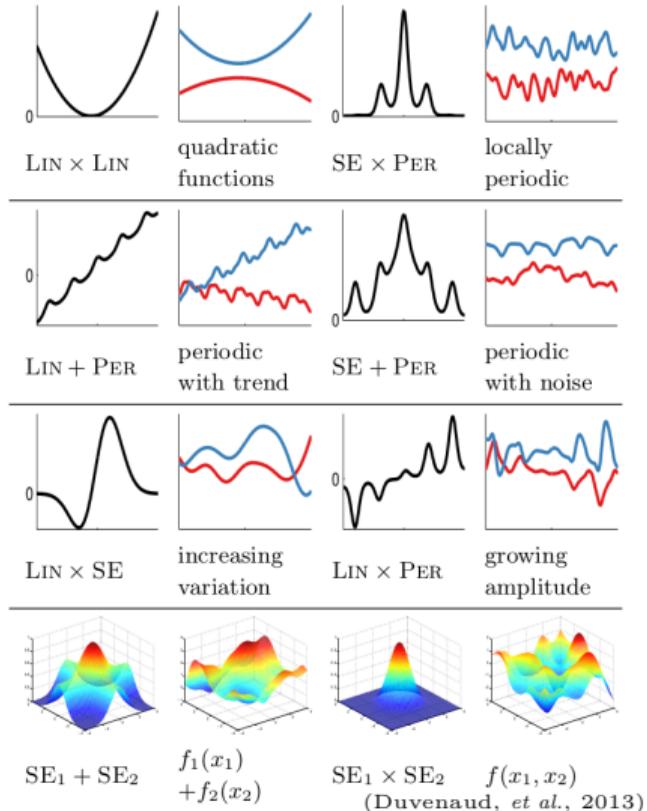
- Functions that generate a **positive semidefinite** matrix \mathbf{K}_f are valid kernels.
- We have seen some of the most used kernels: Squared Exponential, Periodic, Matérn, Polynomial, etc.
- How to find suitable kernel functions for a given problem?

Composition of base kernels



(Duvenaud, et al., 2013)

- Duvenaud, et al. “**Structure discovery in nonparametric regression through compositional kernel search**”, ICML, 2013.
- Malkomes, et al. “**Bayesian optimization for automated model selection**”. ICML AutoML Workshop, 2016.
- Kim and Teh. “**Scaling up the Automatic Statistician: Scalable structure discovery using Gaussian processes**”. AISTATS, 2018.
- Teng, et al. “**Scalable variational Bayesian kernel selection for sparse Gaussian process regression**”. AAAI, 2020.

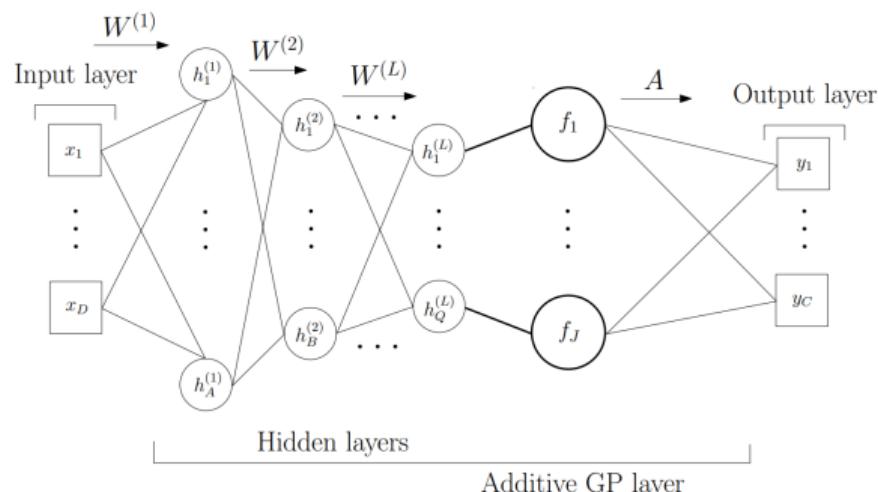


Deep Kernel Learning

- Mapping the input space by an arbitrary function $g(\cdot)$ results in a valid kernel:

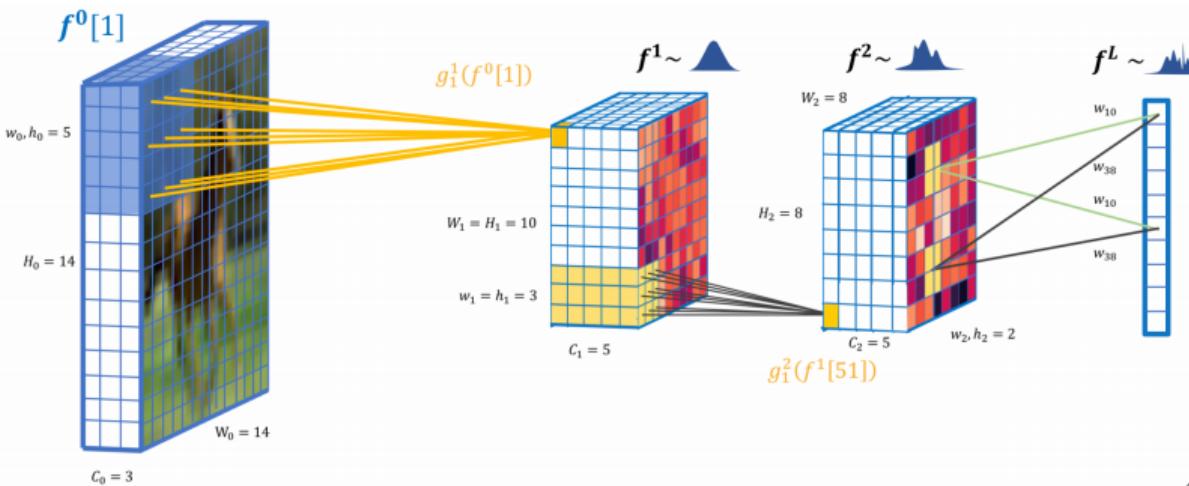
$$\tilde{k}(\mathbf{x}_i, \mathbf{x}_j) = k(g(\mathbf{x}_i), g(\mathbf{x}_j)).$$

- Calandra *et al.* (2016) model $g(\cdot)$ with a **neural network** parameterized by θ that is **jointly optimized** following the gradients of the marginal likelihood $\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial \tilde{K}} \frac{\partial \tilde{K}}{\partial g_\theta} \frac{\partial g_\theta}{\partial \theta}$.
- Wilson *et al.* (2016) scale the approach to large datasets.



(Wilson, *et al.*, 2016)

Convolutional Gaussian Processes



(Blomqvist, et al., 2018)

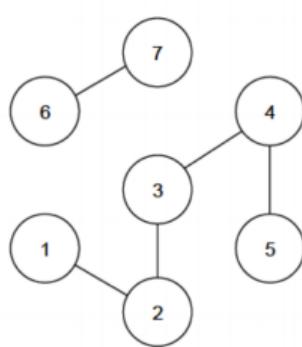
- Sum of patch responses $g(\mathbf{x}^{[p]})$ of a given image \mathbf{x} :

$$f(\mathbf{x}) = \sum_{p=1}^P g(\mathbf{x}^{[p]}), \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, k_g(\mathbf{x}, \mathbf{x}')),$$

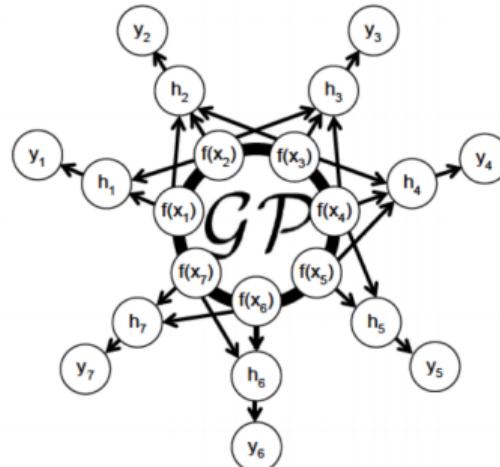
$$\mathbf{f} \sim \mathcal{N} \left(\mathbf{0}, \sum_{p=1}^P \sum_{p'=1}^P k_g(\mathbf{x}^{[p]}, \mathbf{x}^{[p']}) \right).$$

- Van der Wilk, et al. “Convolutional Gaussian processes”. NeurIPS, 2017.
- Blomqvist, et al. “Deep convolutional Gaussian processes”. ECMLPKDD, 2018.
- Dutordoir, et al. “Bayesian image classification with deep convolutional Gaussian processes”. AISTATS, 2020.

Gaussian Processes over graphs



Observed
Graph



- Let $\mathbf{A} \in \{0, 1\}^{N \times N}$ be an **adjacency matrix**, we have:

$$p(\mathbf{y}, \mathbf{h} | \mathbf{X}, \mathbf{A}) = p(\mathbf{h} | \mathbf{X}, \mathbf{A}) \prod_{i=1}^N p(y_i, h_i),$$

$$p(\mathbf{h} | \mathbf{X}, \mathbf{A}) = \mathcal{N}(\mathbf{0}, \mathbf{P} \mathbf{K}_f \mathbf{P}^\top), \quad p(\mathbf{f} | \mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$

$$\mathbf{P} = (\mathbf{I} + \mathbf{D})^{-1}(\mathbf{I} + \mathbf{A}),$$

where \mathbf{D} is the **degree matrix**.

(Ng, et al., 2018)

- Ng et al. “**Bayesian Semi-supervised Learning with Graph Gaussian Processes**”. NeurIPS, 2018.
- Walker and Glocker. “**Graph convolutional Gaussian processes**”. ICML, 2019.
- Borovitskiy, et al. “**Matérn Gaussian processes on graphs**”. AISTATS, 2021.

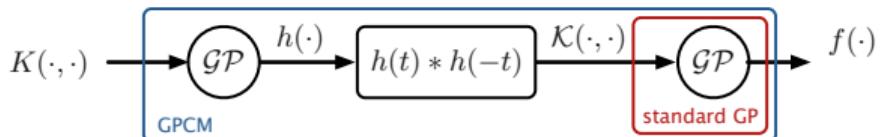
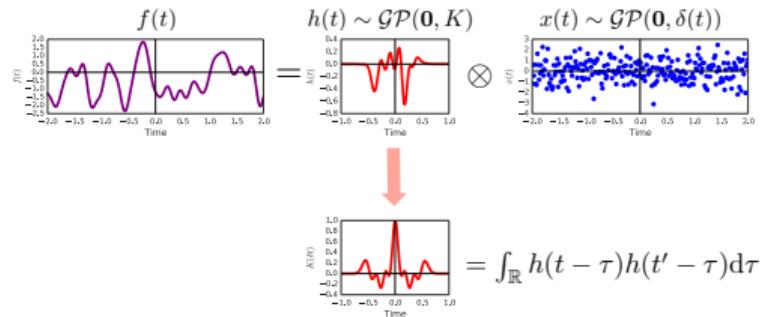
Nonparametric kernels

- Gaussian Process Convolution Model (GPCM)

- i) draw from the filter function $h \sim \mathcal{GP}(\mathbf{0}, K)$;
- ii) convolve it with a white noise process $x(t)$:

$$f(t) = \int_{\mathbb{R}} h(t - \tau) x(\tau) d\tau.$$

- Conditionally, $f|h \sim \mathcal{GP}(\mathbf{0}, \mathcal{K})$ is a GP with a nonparametric covariance function.
- Closely related to linear time-invariant (LTI) systems, with $h(t)$ acting as an impulse response.



Tobar, et al. “Learning stationary time series using Gaussian processes with nonparametric kernels”, Neurips, 2015.

Current trends on kernel design

Summary

- The choice of the **kernel function** defines the properties of a GP model.
- New kernels can be designed by **manually or automatically combining** base kernels.
- **Deep kernel learning** enables flexible kernels jointly optimized with the other GP hyperparameters.
- Specific data structures can be handled with **convolutional or graph kernels**.
- It is possible to use convolutions of processes to obtain **nonparametric kernels** and work with GPs in the spectral domain.

Table of Contents

① Part I

Fundamentals of Bayesian inference

Bayesian linear regression

The Gaussian process

② Part II

Learning the kernel and its hyperparameters

Beyond Gaussian likelihood

Sparse approximations

③ Part III

GPs for Bayesian optimization

Current trends on kernel design

From GPLVM to Deep GPs

Concluding remarks and additional topics

From GPLVM to Deep GPs

- So far we have considered supervised learning tasks of the form:

$$y_i = f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i),$$
$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f).$$

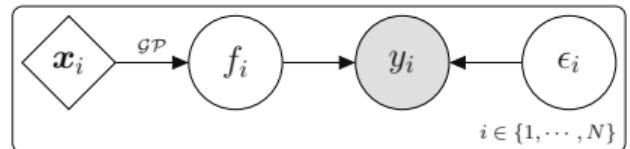
- $y_i|_{i=1}^N$ are **noisy observations**;
- $f_i|_{i=1}^N$ are **latent random variables**;
- $\mathbf{x}_i|_{i=1}^N$ are **deterministic observed inputs**.

- For multidimensional observations $\mathbf{Y} \in \mathbb{R}^{N \times D_y}$ and random inputs \mathbf{X} :



Unobserved inputs X

nonlinear unsupervised projection of Y



Observed inputs X

supervised learning with noisy inputs

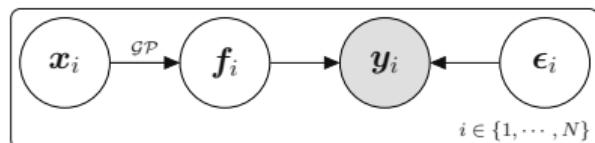
From GPLVM to Deep GPs

- So far we have considered supervised learning tasks of the form:

$$y_i = f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \\ p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f).$$

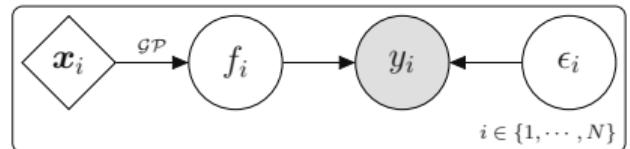
- $y_i|_{i=1}^N$ are **noisy observations**;
- $f_i|_{i=1}^N$ are **latent random variables**;
- $\mathbf{x}_i|_{i=1}^N$ are **deterministic observed inputs**.

- For multidimensional observations $\mathbf{Y} \in \mathbb{R}^{N \times D_y}$ and random inputs \mathbf{X} :



Unobserved inputs X

nonlinear unsupervised projection of \mathbf{Y}



Observed inputs X

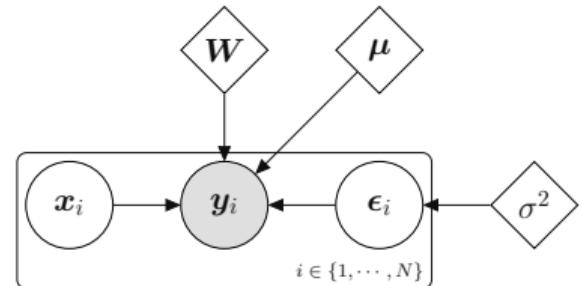
supervised learning with noisy inputs

PPCA - Probabilistic PCA

- The PPCA is one of the simplest models with **continuous latent variables**:

$$\underbrace{p(\mathbf{X})}_{\text{prior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}), \quad p(\boldsymbol{\epsilon}_i) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{likelihood}} = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}).$$



- Since the PPCA is linear, we have analytical marginal likelihood and posterior:

$$\underbrace{p(\mathbf{Y}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{marg. likelihood}} = \int \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}) d\mathbf{x}_i = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{X}|\mathbf{Y}, \hat{\mathbf{W}}, \hat{\boldsymbol{\mu}}, \hat{\sigma}^2)}_{\text{posterior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \hat{\mathbf{M}}^{-1} \hat{\mathbf{W}}^\top (\mathbf{y}_i - \hat{\boldsymbol{\mu}}), \hat{\sigma}^2 \hat{\mathbf{M}}^{-1}), \quad \text{where } \hat{\mathbf{M}} = (\hat{\mathbf{W}}^\top \hat{\mathbf{W}} + \hat{\sigma}^2 \mathbf{I})^{-1}.$$

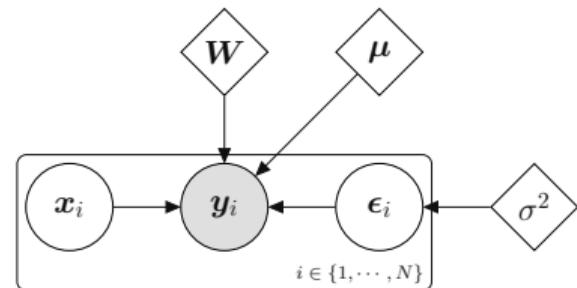
- ML estimates $\hat{\boldsymbol{\mu}}$, $\hat{\mathbf{W}}$ and $\hat{\sigma}^2$ can be obtained in closed form or via EM algorithm.

PPCA - Probabilistic PCA

- The PPCA is one of the simplest models with **continuous latent variables**:

$$\underbrace{p(\mathbf{X})}_{\text{prior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}), \quad p(\boldsymbol{\epsilon}_i) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{likelihood}} = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}).$$



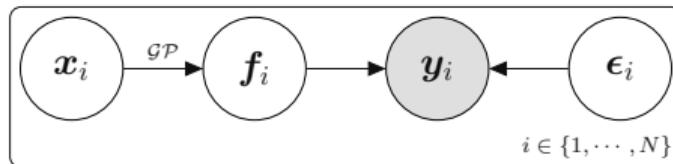
- Since the PPCA is **linear**, we have analytical marginal likelihood and posterior:

$$\underbrace{p(\mathbf{Y}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{marg. likelihood}} = \int \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}) d\mathbf{x}_i = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{X}|\mathbf{Y}, \hat{\mathbf{W}}, \hat{\boldsymbol{\mu}}, \hat{\sigma}^2)}_{\text{posterior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \hat{\mathbf{M}}^{-1} \hat{\mathbf{W}}^\top (\mathbf{x}_i - \hat{\boldsymbol{\mu}}), \sigma^2 \hat{\mathbf{M}}^{-1}), \quad \text{where } \hat{\mathbf{M}} = (\hat{\mathbf{W}}^\top \hat{\mathbf{W}} + \hat{\sigma}^2 \mathbf{I})^{-1}.$$

- ML estimates $\hat{\boldsymbol{\mu}}$, $\hat{\mathbf{W}}$ and $\hat{\sigma}^2$ can be obtained in closed form or via EM algorithm.

GPLVM - Gaussian Process Latent Variable Model

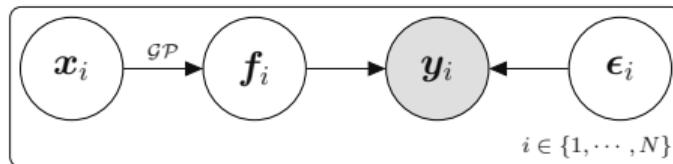


- Inspired by the PPCA formulation, Lawrence (2005) proposed the GPLVM:

$$\begin{aligned} p(\mathbf{X}) &= \prod_{i=1}^N \underbrace{\mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})}_{\text{input prior}}, & p(\mathbf{F} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}}, \\ p(\mathbf{Y} | \mathbf{F}, \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{f}_{:d}, \sigma^2 \mathbf{I})}_{\text{likelihood}} \mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f), \\ p(\mathbf{Y} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{0}, \mathbf{K}_f + \sigma^2 \mathbf{I})}_{\text{marginal likelihood}}. \end{aligned}$$

- Marginalization of \mathbf{F} is analytic, but of the random input \mathbf{X} is not.
→ \mathbf{X} appears in a nonlinear way inside the covariance matrix \mathbf{K}_f .

GPLVM - Gaussian Process Latent Variable Model

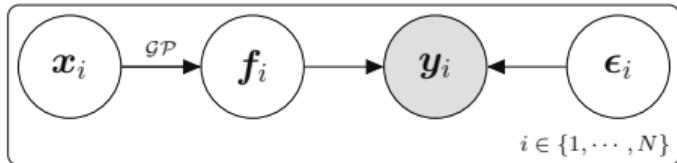


- Inspired by the PPCA formulation, Lawrence (2005) proposed the GPLVM:

$$\begin{aligned} p(\mathbf{X}) &= \prod_{i=1}^N \underbrace{\mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})}_{\text{input prior}}, & p(\mathbf{F} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}}, \\ p(\mathbf{Y} | \mathbf{F}, \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{f}_{:d}, \sigma^2 \mathbf{I})}_{\text{likelihood}} \mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f), \\ p(\mathbf{Y} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{0}, \mathbf{K}_f + \sigma^2 \mathbf{I})}_{\text{marginal likelihood}}. \end{aligned}$$

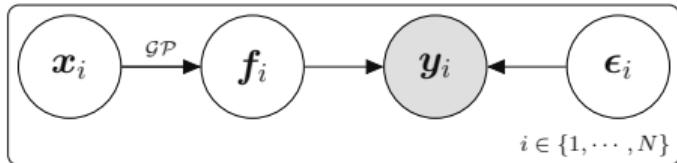
- Marginalization of \mathbf{F} is analytic, but of the random input \mathbf{X} is not.
→ \mathbf{X} appears in a nonlinear way inside the covariance matrix \mathbf{K}_f .

GPLVM - Gaussian Process Latent Variable Model



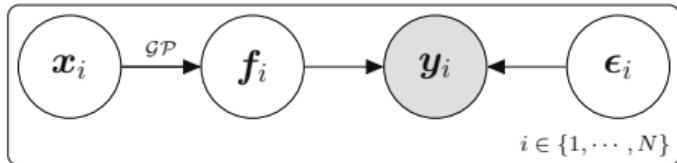
- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(x_i | m_i, S_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from x_* to y_* .
- Latent projection given a new observation y_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

GPLVM - Gaussian Process Latent Variable Model



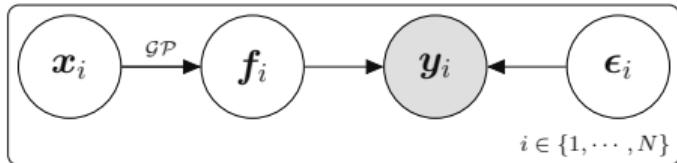
- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{m}_i, \mathbf{S}_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from \mathbf{x}_* to \mathbf{y}_* .
- Latent projection given a new observation \mathbf{y}_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

GPLVM - Gaussian Process Latent Variable Model



- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{m}_i, \mathbf{S}_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from \mathbf{x}_* to \mathbf{y}_* .
- Latent projection given a new observation \mathbf{y}_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

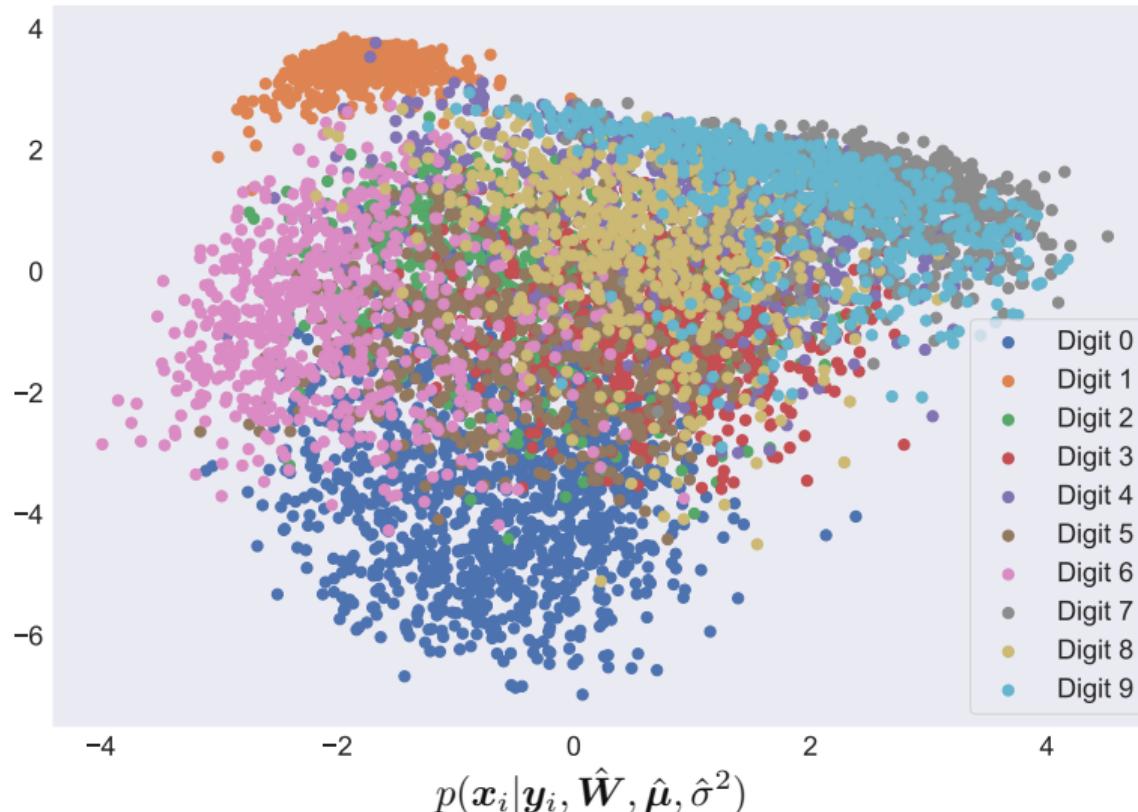
GPLVM - Gaussian Process Latent Variable Model



- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{m}_i, \mathbf{S}_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from \mathbf{x}_* to \mathbf{y}_* .
- Latent projection given a new observation \mathbf{y}_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

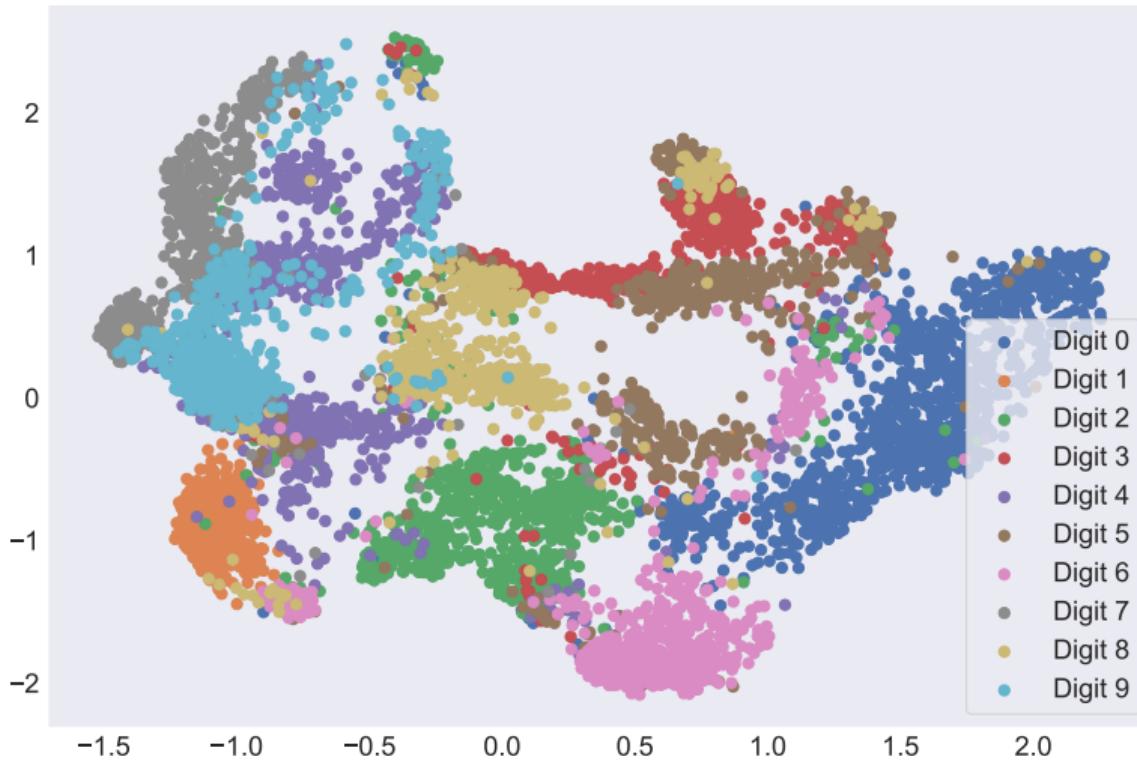
USPS digits - 7291 training samples

256 → 2 dimension projection with PCA



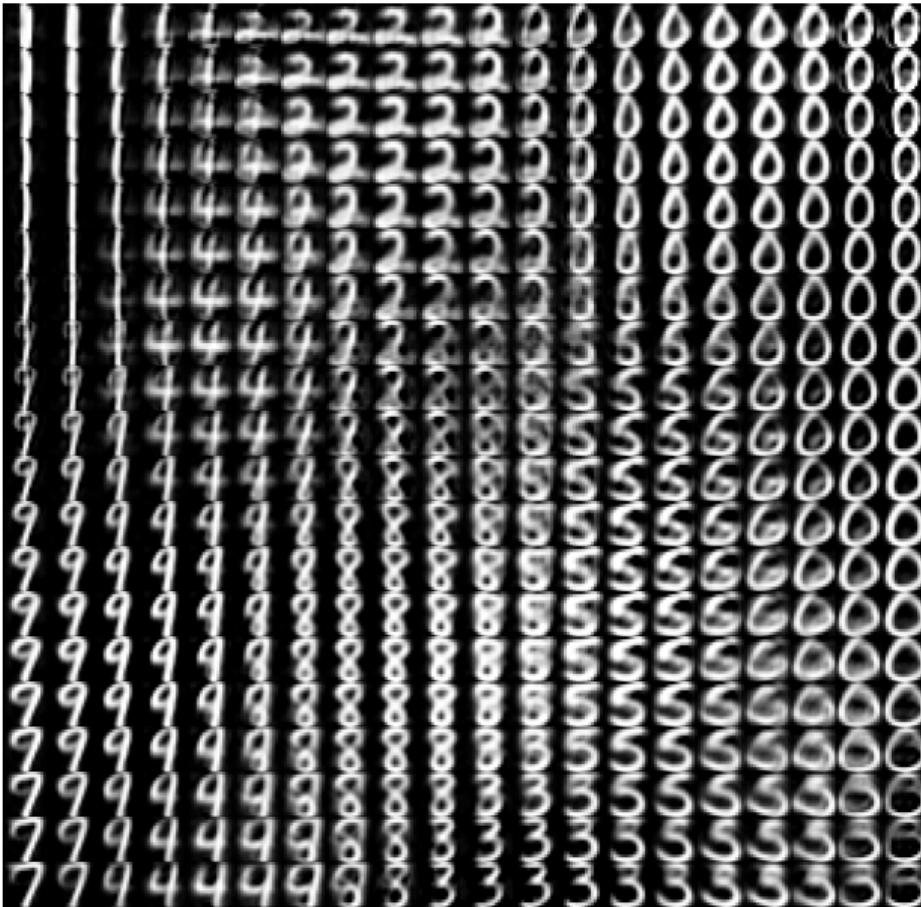
USPS digits - 7291 training samples

256 → 2 dimension projection with GPLVM (50 inducing points)

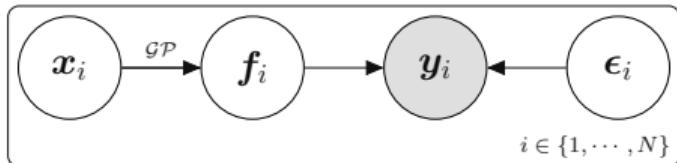


$$p(\mathbf{x}_i|\mathbf{y}_i) \approx q(\mathbf{x}_i|\mathbf{y}_i) = \mathcal{N}(\mathbf{x}_i|\mathbf{m}_i, \mathbf{S}_i).$$

USPS digits - GPLVM samples

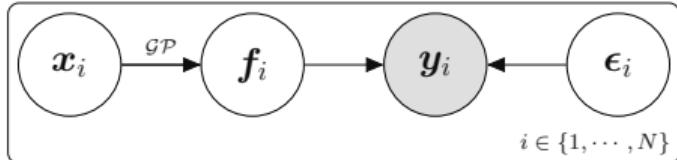


GPLVM - Gaussian Process Latent Variable Model



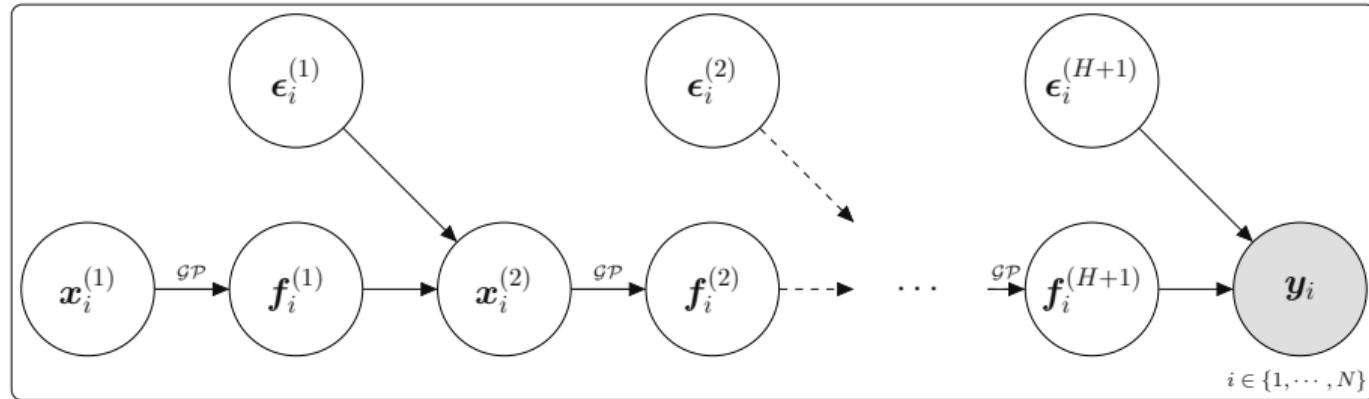
- Instead of a simple prior $p(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})$, we could have different choices:
 - **Discriminative priors** for classification (Urtasun and Darrell, 2007);
 - **Dynamical priors** with temporal dependence (Wang *et al.*, 2007; Damianou *et al.*, 2011; Frigola *et al.*, 2014; Mattos *et al.*, 2016);
 - Mixture of **visible and missing** data (Damianou and Lawrence, 2015).
 - **Conditional** variables (Dutordoir *et al.*, 2016);
 - Other(s) GP prior(s) (Lawrence and Moore, 2007; Damianou and Lawrence, 2013; Mattos *et al.*, 2016; Bui *et al.*, 2016; Salimbeni and Deisenroth, 2017).

GPLVM - Gaussian Process Latent Variable Model



- Instead of a simple prior $p(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})$, we could have different choices:
 - **Discriminative priors** for classification (Urtasun and Darrell, 2007);
 - **Dynamical priors** with temporal dependence (Wang *et al.*, 2007; Damianou *et al.*, 2011; Frigola *et al.*, 2014; Mattos *et al.*, 2016);
 - Mixture of **visible and missing** data (Damianou and Lawrence, 2015).
 - **Conditional** variables (Dutordoir *et al.*, 2016);
 - **Other(s) GP prior(s)** (Lawrence and Moore, 2007; Damianou and Lawrence, 2013; Mattos *et al.*, 2016; Bui *et al.*, 2016; Salimbeni and Deisenroth, 2017).

Deep Gaussian Processes

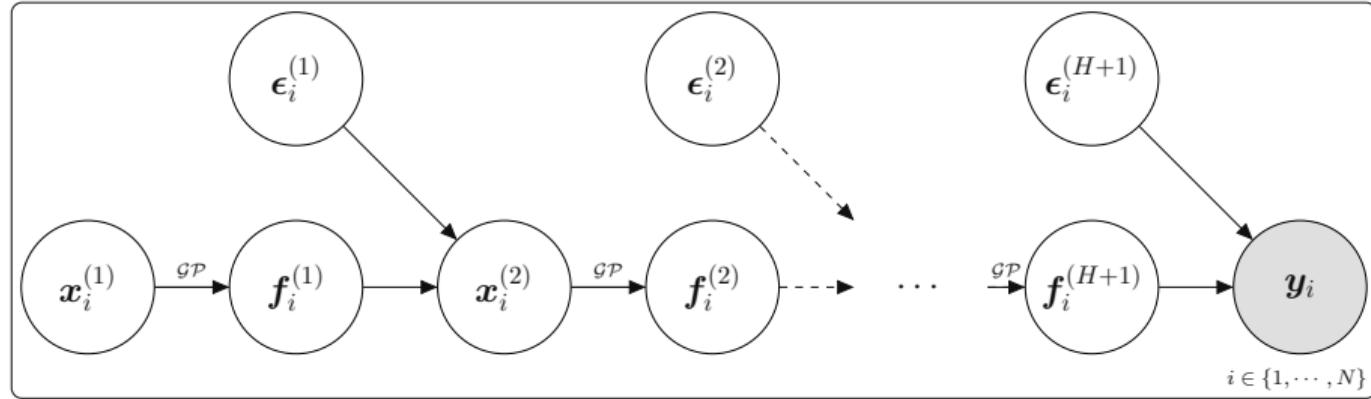


- A **hierarchy** of H GP priors gives rise to a **deep** GP model with H hidden layers:

$$\begin{aligned} \mathbf{y}_i &= f^{(H+1)}(\mathbf{x}_i^{(H+1)}) + \epsilon_i^{(H+1)}, & \mathbf{f}_{:d}^{(H+1)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(H+1)}) \\ \mathbf{x}_i^{(h+1)} &= f^{(h)}(\mathbf{x}_i^{(h)}) + \epsilon_i^{(h)}, & \mathbf{f}_{:d}^{(h)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(h)}), \quad 1 \leq h \leq H. \end{aligned}$$

- The intractabilities require approximate inference and inducing point approaches.

Deep Gaussian Processes

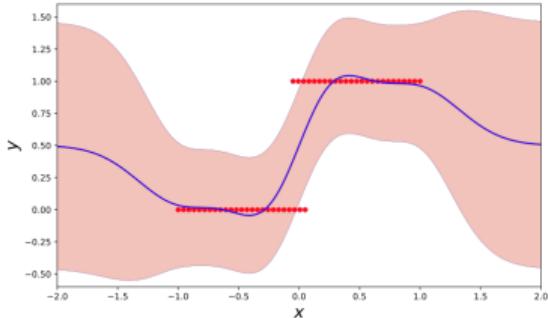


- A **hierarchy** of H GP priors gives rise to a **deep** GP model with H hidden layers:

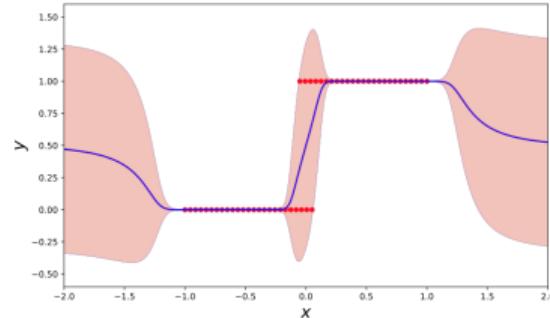
$$\begin{aligned} \mathbf{y}_i &= f^{(H+1)}(\mathbf{x}_i^{(H+1)}) + \epsilon_i^{(H+1)}, & \mathbf{f}_{:d}^{(H+1)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(H+1)}) \\ \mathbf{x}_i^{(h+1)} &= f^{(h)}(\mathbf{x}_i^{(h)}) + \epsilon_i^{(h)}, & \mathbf{f}_{:d}^{(h)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(h)}), \quad 1 \leq h \leq H. \end{aligned}$$

- The intractabilities require approximate inference and inducing point approaches.

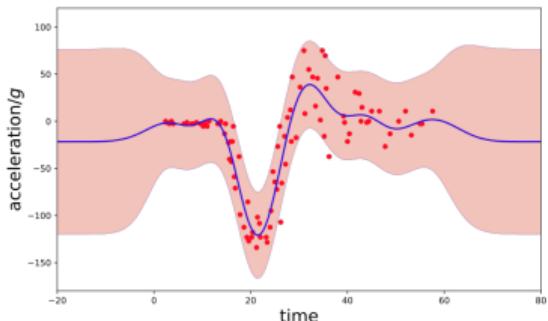
Deep Gaussian Processes



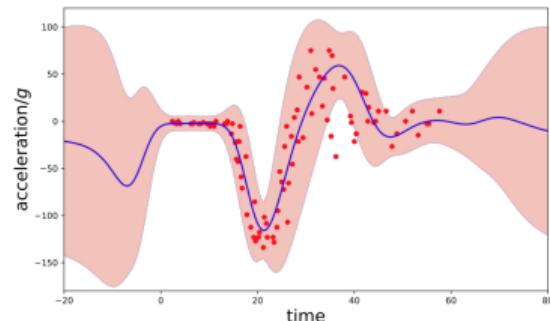
(a) Shallow GP.



(b) Deep GP.



(c) Shallow GP.



(d) Deep GP.

(Lawrence, 2020)

Deep Gaussian Processes

- But what are the advantages of deep GPs?
 - The resulting model **cannot be explained** by a single GP.
 - Multiple layers **alleviate the choice** of the kernel function.
 - Enables **hierarchical** feature learning.
 - Deep GPs are **related to deep neural networks** with infinite hidden units.
- Supported by modern frameworks, such as **GPyTorch** and **GPflux**.
- Research on deep GPs has a wide range of topics:
 - Variational Auto-encoded Deep Gaussian Processes (Dai *et al.*, 2016);
 - Recurrent Gaussian Processes (Mattos *et al.*, 2016);
 - Doubly Stochastic Variational Inference for Deep GPs (Salimbeni and Deisenroth, 2017);
 - Deep Convolutional GPs (Blomqvist *et al.*, 2018);
 - Deep GPs with Importance-Weighted Variational Inference (Salimbeni *et al.*, 2019);
 - Stochastic Deep Gaussian Processes over Graphs (Li *et al.*, 2020);
 - Global Inducing Point Variational Posteriors for BNNs and DGPs (Ober and Aitchison, 2021).

Deep Gaussian Processes

- But what are the advantages of deep GPs?
 - The resulting model **cannot be explained** by a single GP.
 - Multiple layers **alleviate the choice** of the kernel function.
 - Enables **hierarchical** feature learning.
 - Deep GPs are **related to deep neural networks** with infinite hidden units.
- Supported by modern frameworks, such as **GPyTorch** and **GPflux**.
- Research on deep GPs has a wide range of topics:
 - Variational Auto-encoded Deep Gaussian Processes (Dai *et al.*, 2016);
 - Recurrent Gaussian Processes (Mattos *et al.*, 2016);
 - Doubly Stochastic Variational Inference for Deep GPs (Salimbeni and Deisenroth, 2017);
 - Deep Convolutional GPs (Blomqvist *et al.*, 2018);
 - Deep GPs with Importance-Weighted Variational Inference (Salimbeni *et al.*, 2019);
 - Stochastic Deep Gaussian Processes over Graphs (Li *et al.*, 2020);
 - Global Inducing Point Variational Posteriors for BNNs and DGPs (Ober and Aitchison, 2021).

Deep Gaussian Processes

- But what are the advantages of deep GPs?
 - The resulting model **cannot be explained** by a single GP.
 - Multiple layers **alleviate the choice** of the kernel function.
 - Enables **hierarchical** feature learning.
 - Deep GPs are **related to deep neural networks** with infinite hidden units.
- Supported by modern frameworks, such as **GPyTorch** and **GPflux**.
- Research on deep GPs has a wide range of topics:
 - Variational Auto-encoded Deep Gaussian Processes (Dai *et al.*, 2016);
 - Recurrent Gaussian Processes (Mattos *et al.*, 2016);
 - Doubly Stochastic Variational Inference for Deep GPs (Salimbeni and Deisenroth, 2017);
 - Deep Convolutional GPs (Blomqvist *et al.*, 2018);
 - Deep GPs with Importance-Weighted Variational Inference (Salimbeni *et al.*, 2019);
 - Stochastic Deep Gaussian Processes over Graphs (Li *et al.*, 2020);
 - Global Inducing Point Variational Posteriors for BNNs and DGPs (Ober and Aitchison, 2021).

From GPLVM to Deep GPs

Summary

- GPLVM enables **unsupervised learning** with GP priors in a **generative setting**.
- GP learning with **missing input data** uses the same GPLVM inference tools.
- A **hierarchy** of GPLVM blocks results in a **deep GP model**.
- The multilayer **composition of processes** cannot be represented by a single GP.
- Bayesian inference with GPLVMs and deep GPs is hard, but some **open source modern frameworks** are available.

Table of Contents

① Part I

- Fundamentals of Bayesian inference
- Bayesian linear regression
- The Gaussian process

② Part II

- Learning the kernel and its hyperparameters
- Beyond Gaussian likelihood
- Sparse approximations

③ Part III

- GPs for Bayesian optimization
- Current trends on kernel design
- From GPLVM to Deep GPs
- Concluding remarks and additional topics

Concluding Remarks

- GPs are **Bayesian nonparametric** models, meaning that they:
 - i) naturally combine existing biases and the knowledge extracted from data;
 - ii) provide **probabilistic estimates**;
 - iii) the more data they see, the better they become.
- GPs can be applied to a plethora of settings: denoising, regression, interpolation, extrapolation, classification, optimization, unsupervised learning, hierarchical modelling, etc.
- **Hyperparameter optimization** can be performed via marginal likelihood, and jointly with sparse approximations, if they are required.
- The GP community is very active, come on over and join us!

Additional GP topics

Gaussian Process Summer Schools - GPSS

Superb presentations and workshops on several GP related subjects: <http://gpss.cc/>.

- More thorough discussion on sparse and scalable GPs
 - Bauer, van der Wilk and Rasmussen. “Understanding probabilistic sparse Gaussian process approximations”. NeurIPS, 2016.
 - Bui, Yan and Turner. “A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation”. JMLR, 2017.
 - Liu, *et al.* “When Gaussian process meets big data: A review of scalable GPs”. IEEE TNNLS, 2020.
 - Leibfried, *et al.* “A tutorial on sparse Gaussian processes and variational inference”. arXiv preprint, 2021.
- Bayesian optimization
 - Frazier, “Bayesian Optimization”. INFORMS Tutorials, 2018.
 - Shahriari, Swersky, Wang, Adams and de Freitas, “Taking the Human Out of the Loop: A Review of Bayesian Optimization”. Proc. of the IEEE, 2016.
 - Snoek, Larochelle and Adams, “Practical Bayesian Optimization of Machine Learning Algorithms”. NeurIPS, 2012.

Additional GP topics

Gaussian Process Summer Schools - GPSS

Superb presentations and workshops on several GP related subjects: <http://gpss.cc/>.

- More thorough discussion on sparse and scalable GPs
 - Bauer, van der Wilk and Rasmussen. “**Understanding probabilistic sparse Gaussian process approximations**”. NeurIPS, 2016.
 - Bui, Yan and Turner. “**A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation**”. JMLR, 2017.
 - Liu, *et al.* “**When Gaussian process meets big data: A review of scalable GPs**”. IEEE TNNLS, 2020.
 - Leibfried, *et al.* “**A tutorial on sparse Gaussian processes and variational inference**”. arXiv preprint, 2021.
- Bayesian optimization
 - Frazier, “**Bayesian Optimization**”. INFORMS Tutorials, 2018.
 - Shahriari, Swersky, Wang, Adams and de Freitas, “**Taking the Human Out of the Loop: A Review of Bayesian Optimization**”. Proc. of the IEEE, 2016.
 - Snoek, Larochelle and Adams, “**Practical Bayesian Optimization of Machine Learning Algorithms**”. NeurIPS, 2012.

Additional GP topics

Gaussian Process Summer Schools - GPSS

Superb presentations and workshops on several GP related subjects: <http://gpss.cc/>.

- More thorough discussion on sparse and scalable GPs
 - Bauer, van der Wilk and Rasmussen. “**Understanding probabilistic sparse Gaussian process approximations**”. NeurIPS, 2016.
 - Bui, Yan and Turner. “**A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation**”. JMLR, 2017.
 - Liu, *et al.* “**When Gaussian process meets big data: A review of scalable GPs**”. IEEE TNNLS, 2020.
 - Leibfried, *et al.* “**A tutorial on sparse Gaussian processes and variational inference**”. arXiv preprint, 2021.
- Bayesian optimization
 - Frazier, “**Bayesian Optimization**”. INFORMS Tutorials, 2018.
 - Shahriari, Swersky, Wang, Adams and de Freitas, “**Taking the Human Out of the Loop: A Review of Bayesian Optimization**”. Proc. of the IEEE, 2016.
 - Snoek, Larochelle and Adams, “**Practical Bayesian Optimization of Machine Learning Algorithms**”. NeurIPS, 2012.

Additional GP topics

- Deep GPs and beyond
 - Salimbeni and Deisenroth. “**Doubly stochastic variational inference for deep Gaussian processes**”. NeurIPS, 2017.
 - Salimbeni, *et al.* “**Deep Gaussian processes with importance-weighted variational inference**”. ICML, 2019.
 - Bui, *et al.* “**Deep Gaussian processes for regression using approximate expectation propagation**”. ICML, 2016.
 - Dunlop, *et al.* “**How deep are deep Gaussian processes?**”. JMLR, 2018.
 - Aitchison, Yang and Ober. “**Deep kernel processes**”. ICML, 2021.
- GPs, deep learning and Neural Tangent Kernel (NTK)
 - Matthews, *et al.* “**Gaussian process behaviour in wide deep neural networks**”. ICLR, 2018.
 - Lee, *et al.* “**Deep neural networks as Gaussian processes**”. ICLR, 2018.
 - Jacot, Gabriel and Hongler. “**Neural tangent kernel: Convergence and generalization in neural networks**”. NeurIPS, 2018.
 - Dutordoir, *et al.* “**Deep Neural Networks as Point Estimates for Deep Gaussian Processes**”. arXiv preprint, 2021.

Additional GP topics

- Deep GPs and beyond
 - Salimbeni and Deisenroth. “**Doubly stochastic variational inference for deep Gaussian processes**”. NeurIPS, 2017.
 - Salimbeni, *et al.* “**Deep Gaussian processes with importance-weighted variational inference**”. ICML, 2019.
 - Bui, *et al.* “**Deep Gaussian processes for regression using approximate expectation propagation**”. ICML, 2016.
 - Dunlop, *et al.* “**How deep are deep Gaussian processes?**”. JMLR, 2018.
 - Aitchison, Yang and Ober. “**Deep kernel processes**”. ICML, 2021.
- GPs, deep learning and Neural Tangent Kernel (NTK)
 - Matthews, *et al.* “**Gaussian process behaviour in wide deep neural networks**”. ICLR, 2018.
 - Lee, *et al.* “**Deep neural networks as Gaussian processes**”. ICLR, 2018.
 - Jacot, Gabriel and Hongler. “**Neural tangent kernel: Convergence and generalization in neural networks**”. NeurIPS, 2018.
 - Dutordoir, *et al.* “**Deep Neural Networks as Point Estimates for Deep Gaussian Processes**”. arXiv preprint, 2021.

Questions?

César Lincoln C. Mattos
@cesarlincoln
cesarlincoln@dc.ufc.br