

# Gaussian Processes

Definition, applications and deep extension

---

Luis Antonio Ortega Andrés

March 26, 2021

# Table of contents

## 1. Definition

Characterization

Examples

## 2. Regression Problem

Computational complexity

## 3. Inducing points

## 4. Deep Gaussian processes

# Definitions

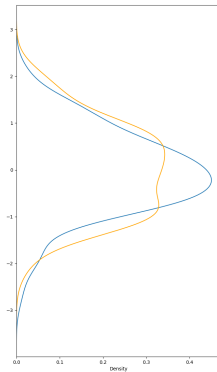
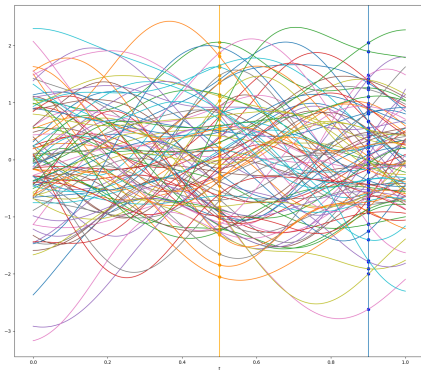
---

## Gaussian process

A **Gaussian process** is a **stochastic process**, i.e., a collection of random variables  $\{X_t\}_{t \in \mathcal{T}}$  indexed by a set  $\mathcal{T}$ , such that any finite subset is Gaussian.

$$\{X_{t_1}, \dots, X_{t_N}\} \sim \mathcal{N}(\cdot, \cdot)$$

For example,  $X_{t_i} \sim \mathcal{N}(\cdot, \cdot)$ .



Gaussian processes are completely determined by their first and second order moments<sup>1</sup>.

Given  $\mathbf{t} = (t_1, \dots, t_N)$ ,  $N \in \mathbb{N}$ :

$$X(\mathbf{t}) = (X_{t_1}, \dots, X_{t_N}) \sim \mathcal{N}(m(\mathbf{t}), K(\mathbf{t}, \mathbf{t}))$$

$$\text{where } \begin{cases} m(\mathbf{t}) &= \mathbb{E}[X(\mathbf{t})] \\ K(\mathbf{t}, \mathbf{t}) &= \text{Cov}(X(\mathbf{t}), X(\mathbf{t})) \end{cases}$$

Defining  $m$  and  $K$  we get a Gaussian process.

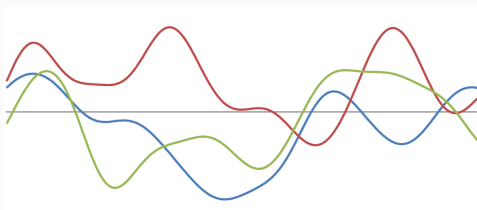
---

<sup>1</sup>Bishop, Christopher M. Pattern recognition and machine learning. Springer, 2006.

# Examples

It is usual to take a **zero mean function**,  $m(\mathbf{t}) = 0$  and **kernel functions**:

- *RBF*:  $K(\mathbf{t}, \mathbf{t}') = \exp\left(-\frac{\|\mathbf{t} - \mathbf{t}'\|^2}{2\sigma^2}\right)$ .
- *Matérn*: Family of kernels, parameterized by  $\nu$ . Generalize several kernels.



Define a **probability distribution over functions**: *Given a function, which is its probability when interpreted as a sample of a Gaussian process?*

# Regression Problem

---



# Problem statement

Given dataset  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  and an unknown function  $f$  such that

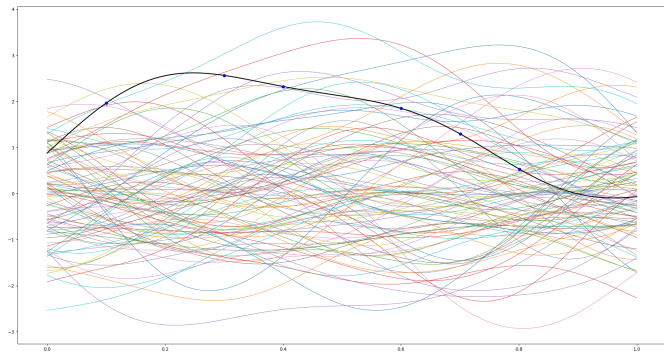
$$y_n = f(\mathbf{x}_n) + \epsilon \quad \forall n = 1, \dots, N, \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

## Assumption

Function  $f$  is a Gaussian process of unknown mean function  $m$  and kernel function  $K$

## Assumption

Function  $f$  is a Gaussian process of unknown mean function  $m$  and kernel function  $K$



Naming  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  and  $\mathbf{y} = (y_1, \dots, y_N)$ :

$$\mathbf{y} = f(\mathbf{X}) + \mathcal{N}(0, \sigma^2 \mathbf{I}) \implies \mathbf{y} \sim \mathcal{N}(m(\mathbf{X}), K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}).$$

Remark: A *distribution is assumed over function points* **but not over**  $\mathbf{x}$ .

Let  $\mathbf{X}^*$  be a test case where  $\mathbf{y}^* = f(\mathbf{X}^*) + \mathcal{N}(0, \sigma^2 \mathbf{I})$ :

$$\begin{pmatrix} f(\mathbf{X}) \\ f(\mathbf{X}^*) \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} m(\mathbf{X}) \\ m(\mathbf{X}^*) \end{pmatrix}, \begin{pmatrix} K(\mathbf{X}, \mathbf{X}) & K(\mathbf{X}, \mathbf{X}^*) \\ K(\mathbf{X}^*, \mathbf{X}) & K(\mathbf{X}^*, \mathbf{X}^*) \end{pmatrix} \right)$$

An **usual assumption** is that  $m = 0$ .

### Remark

Typically,  $\mathbf{x}$  is erased from the notation.

Naming  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  and  $\mathbf{y} = (y_1, \dots, y_N)$ :

$$\mathbf{y} = \mathbf{f} + \mathcal{N}(0, \sigma^2 \mathbf{I}) \implies \mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{f}, \mathbf{f}} + \sigma^2 \mathbf{I}).$$

Remark: A *distribution is assumed over  $\mathbf{f}$  points but not over  $\mathbf{X}$* .

Let  $\mathbf{X}^*$  be a test case where  $\mathbf{y}^* = \mathbf{f}^* + \mathcal{N}(0, \sigma^2 \mathbf{I})$ :

$$\begin{pmatrix} \mathbf{f} \\ \mathbf{f}^* \end{pmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \mathbf{K}_{\mathbf{f}, \mathbf{f}} & \mathbf{K}_{\mathbf{f}, \mathbf{f}^*} \\ \mathbf{K}_{\mathbf{f}^*, \mathbf{f}} & \mathbf{K}_{\mathbf{f}^*, \mathbf{f}^*} \end{pmatrix}\right)$$

$$\begin{pmatrix} \mathbf{f} \\ \mathbf{f}^* \end{pmatrix} \sim \mathcal{N} \left( 0, \begin{pmatrix} \mathbf{K}_{\mathbf{f},\mathbf{f}} & \mathbf{K}_{\mathbf{f},\mathbf{f}^*} \\ \mathbf{K}_{\mathbf{f}^*,\mathbf{f}} & \mathbf{K}_{\mathbf{f}^*,\mathbf{f}^*} \end{pmatrix} \right) \quad \mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$$

$\Downarrow$

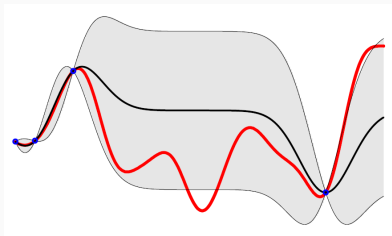
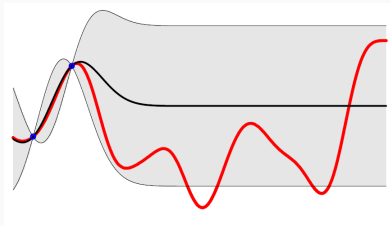
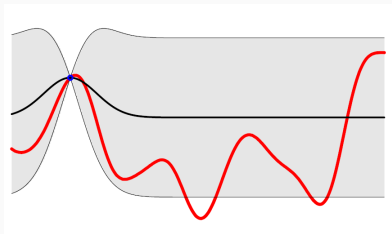
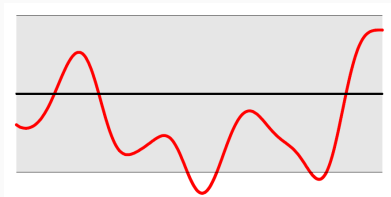
$$P(\mathbf{y} \mid \mathbf{f}, \mathbf{f}^*) = P(\mathbf{y} \mid \mathbf{f}) \implies \mathbf{y} \mid \mathbf{f}, \mathbf{f}^* \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$$

$\Downarrow$

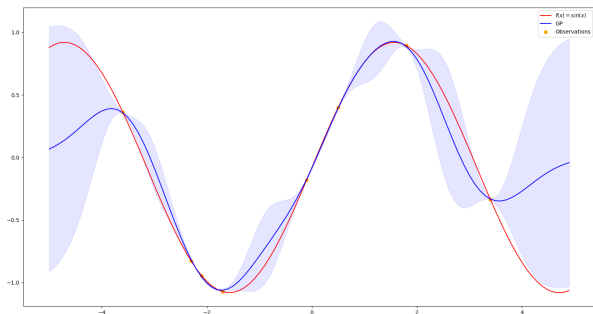
$$\mathbf{f}, \mathbf{f}^* \mid \mathbf{y} \sim \mathcal{N}(\cdot, \cdot) \implies \mathbf{f}^* \mid \mathbf{y} \sim \mathcal{N}(\cdot, \cdot)$$

$\Downarrow$

$$\mathbf{f}^* \mid \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \begin{cases} \boldsymbol{\mu} = \mathbf{K}_{\mathbf{f}^*,\mathbf{f}}(\mathbf{K}_{\mathbf{f},\mathbf{f}} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \\ \boldsymbol{\Sigma} = \mathbf{K}_{\mathbf{f}^*,\mathbf{f}^*} - \mathbf{K}_{\mathbf{f}^*,\mathbf{f}}(\mathbf{K}_{\mathbf{f},\mathbf{f}} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{\mathbf{f},\mathbf{f}^*} \end{cases}$$



Unknown function  $f(x) = \sin(x)$ ,  $\mathcal{D}$  is a sample of 8 points in  $(-5, 5)$ ,  
RBF kernel.



# Computational complexity

Several computations are done, assuming  $\mathbf{X}$  has  $N$  points and  $\mathbf{X}^*$  has  $M$ :

$$\begin{aligned} \mathbf{K}_{f,f} &\implies \mathcal{O}(N^2) \\ (\mathbf{K}_{f,f} + \sigma^2 \mathbf{I})^{-1} &\implies \mathcal{O}(N^3) \\ \mathbf{K}_{f,f^*} = \mathbf{K}_{f^*,f}^T &\implies \mathcal{O}(NM) \\ \mathbf{K}_{f^*,f^*} &\implies \mathcal{O}(M^2) \end{aligned}$$

Training:  $\mathcal{O}(N^3)$  and Test:  $\mathcal{O}(NM^2)$ . They are **computationally inefficient!!**

This can be slightly reduced using the *Cholesky decomposition* for the matrix inversion.



# Advantages

They give a prediction  $\mu = K_{f^*,f}(K_{f,f} + \sigma^2 I)^{-1}y$ . **Equals the kernel ridge regression estimator!!**.

Implicit **confidence interval**,  $(\mu - 3\Sigma, \mu + 3\Sigma)$ .

Full **probabilistic approach**.

## Inducing points

---

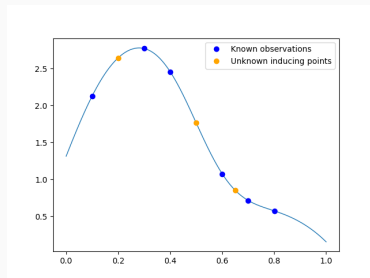
# Inducing points

**Main idea:** Use a *smaller* and *hidden* set of points.

Let  $\mathbf{X}_u \subset \mathbb{R}^D$  be a set of *known* points (commonly computed from  $\mathbf{X}$ ). We make the assumption that  $\mathbf{u} = f(\mathbf{X}_u)$  is representative of  $\mathbf{y}$ .

**Remark.** The *inducing points*  $\mathbf{u}$  are *unknown* and must be *marginalized*.

$$\begin{aligned} P(\mathbf{f}, \mathbf{f}^*) &= \int P(\mathbf{f}, \mathbf{f}^*, \mathbf{u}) d\mathbf{u} \\ &= \underbrace{\int P(\mathbf{f}, \mathbf{f}^* | \mathbf{u}) P(\mathbf{u}) d\mathbf{u}}_{\text{Intractable}} \end{aligned}$$



Where  $\mathbf{u}$  are taken from a Gaussian process:

$$\mathbf{u} \sim \mathcal{N}(0, \mathbf{K}(\mathbf{X}_u, \mathbf{X}_u)).$$

- Exact inference in approximated model:

$$P(\mathbf{f}, \mathbf{f}^* | \mathbf{u}) = P(\mathbf{f} | \mathbf{u}) P(\mathbf{f}^* | \mathbf{u})$$

$\Downarrow$

$$P(\mathbf{f}, \mathbf{f}^*) = \int P(\mathbf{f} | \mathbf{u}) P(\mathbf{f}^* | \mathbf{u}) P(\mathbf{u}) d\mathbf{u}$$

And further approximate  $P(\mathbf{f} | \mathbf{u})$  and  $P(\mathbf{f}^* | \mathbf{u})$

- Variational inference:

$$Q(\mathbf{u}) \approx P(\mathbf{u} | \mathbf{y}).$$

## Exact conditionals

$$P(\mathbf{f} \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \mathbf{K}_{f,f} - \mathbf{Q}_{f,f})$$

$$P(\mathbf{f}^* \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f^*,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \mathbf{K}_{f^*,f^*} - \mathbf{Q}_{f^*,f^*})$$

$$\mathbf{Q}_{a,b} = \mathbf{K}_{a,u} \mathbf{K}_{u,u}^{-1} \mathbf{K}_{u,b}$$

- The Subset of Regressors approximation

$$Q_{SOR}(\mathbf{f} \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, 0)$$

$$Q_{SOR}(\mathbf{f}^* \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f^*,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, 0)$$

- The Deterministic Training Conditional approximation

$$Q_{DTC}(\mathbf{f} \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, 0)$$

$$Q_{DTC}(\mathbf{f}^* \mid \mathbf{u}) = P(\mathbf{f}^* \mid \mathbf{u})$$

## Exact conditionals

$$P(\mathbf{f} \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \mathbf{K}_{f,f} - \mathbf{Q}_{f,f})$$

$$P(\mathbf{f}^* \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f^*,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \mathbf{K}_{f^*,f^*} - \mathbf{Q}_{f^*,f^*})$$

$$\mathbf{Q}_{a,b} = \mathbf{K}_{a,u} \mathbf{K}_{u,u}^{-1} \mathbf{K}_{u,b}$$

- The Fully Independent Training Conditional approximation

$$Q_{FITC}(\mathbf{f} \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \text{diag}(\mathbf{K}_{f,f} - \mathbf{Q}_{f,f}))$$

$$Q_{FITC}(\mathbf{f}^* \mid \mathbf{u}) = P(\mathbf{f}^* \mid \mathbf{u})$$

- The Partially Independent Training Conditional approximation

$$Q_{PITC}(\mathbf{f} \mid \mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \text{blockdiag}(\mathbf{K}_{f,f} - \mathbf{Q}_{f,f}))$$

$$Q_{PITC}(\mathbf{f}^* \mid \mathbf{u}) = P(\mathbf{f}^* \mid \mathbf{u})$$

# Variational bounds

Using Jensen's inequality:

$$\log P(\mathbf{y} \mid \mathbf{u}) = \log \mathbb{E}_{P(\mathbf{f} \mid \mathbf{u})}[P(\mathbf{u} \mid \mathbf{f})] \geq \mathbb{E}_{\log P(\mathbf{f} \mid \mathbf{u})}[P(\mathbf{u} \mid \mathbf{f})] \equiv \mathcal{L}_1,$$

raises Titsias' bound <sup>2</sup>

$$\log P(\mathbf{y}) = \log \int P(\mathbf{y} \mid \mathbf{u})P(\mathbf{u})d\mathbf{u} \geq \log \int \exp \mathcal{L}_1 P(\mathbf{u})d\mathbf{u} \equiv \mathcal{L}_2.$$

But it is not suitable for **stochastic optimization**. Appears a new bound<sup>3</sup>

$$\log P(\mathbf{y}) \geq \mathbb{E}_{Q(\mathbf{u})} [\mathcal{L}_1 + \log P(\mathbf{u}) - \log Q(\mathbf{u})] \equiv \mathcal{L}_3.$$

---

<sup>2</sup>Titsias, Michalis. "Variational learning of inducing variables in sparse Gaussian processes." In Artificial intelligence and statistics, 2009.

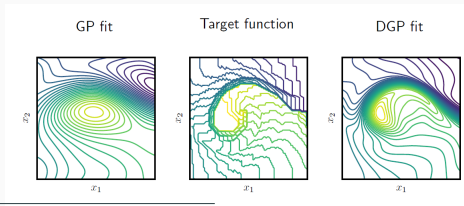
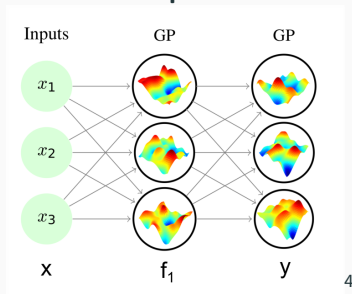
<sup>3</sup>Hensman, James, Nicolo Fusi, and Neil D. Lawrence. "Gaussian processes for big data." 2013.

# Deep Gaussian processes

---



## Connect Gaussian processes in a chain.



<sup>4</sup>Image reference: <https://www.groundai.com/project/inference-in-deep-gaussian-processes-using-stochastic-gradient-hamiltonian-monte-carlo/>

# Definition

Let  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  be a dataset of an unknown function  $f$ .

## Deep Gaussian process

A **deep Gaussian process** of length  $L$  considers  $L$  independent Gaussian processes  $f^1, \dots, f^L$  such that the input of a Gaussian process is the output of the previous one.

$$\mathbf{X}^1 = f^1(\mathbf{X}) \implies \mathbf{X}^2 = f^2(\mathbf{X}^1) \implies \dots \implies \mathbf{X}^L = f^L(\mathbf{X}^{L-1}) \approx \mathbf{Y}$$

## Problem

In the Gaussian process, the input **did not** follow any distribution.

$$\mathbf{X} \text{ no distribution} \implies \mathbf{X}^1 \sim \mathcal{N}(\cdot, \cdot) \implies \mathbf{X}^2 \text{ no longer Gaussian}$$

**Distribution in inner layers cannot be computed in closed form.**

## **Distribution in inner layers cannot be computed in closed form.**

- Variational inference is used to train the model.
- Different evidence lower bounds depending on the assumptions made (inducing points might be considered in each inner layer).
- Distribution is intractable but samples can be taken easily  $\implies$  Monte Carlo.
- Expectation propagation algorithm is used.

**Questions?**



C. M. Bishop.

**Pattern recognition and machine learning.**



T. D. Bui, J. M. Hernández-Lobato, D. Hernández-Lobato, Y. Li, and R. E. Turner.

**Deep Gaussian Processes for Regression using Approximate Expectation Propagation.**



J. Hensman, N. Fusi, and N. D. Lawrence.

**Gaussian Processes for Big Data.**



J. Quiñonero-Candela and C. E. Rasmussen.

**A Unifying View of Sparse Approximate Gaussian Process Regression.**



R. Ranganath, S. Gerrish, and D. M. Blei.

**Black Box Variational Inference.**



H. Salimbeni and M. Deisenroth.

**Doubly Stochastic Variational Inference for Deep Gaussian Processes.**



E. Snelson and Z. Ghahramani.

**Sparse Gaussian Processes using Pseudo-inputs.**



M. K. Titsias.

**Variational Learning of Inducing Variables in Sparse Gaussian Processes.**