

Gaussian Processes

Andrea Casalino

May 31, 2022

1 What is a Gaussian Process?

Gaussian Processes [3], [1], a.k.a. **GPs**, are predictive models able to approximate a multivariate scalar function:

$$g : \mathcal{I} \rightarrow \mathcal{O} \quad (1)$$

$$\mathcal{I} \subseteq \mathbb{R}^i \quad \mathcal{O} \subseteq \mathbb{R} \quad (2)$$

or a multivariate vectorial one:

$$G : \mathcal{I} \rightarrow \mathcal{O} \quad (3)$$

$$\mathcal{I} \subseteq \mathbb{R}^i \quad \mathcal{O} \subseteq \mathbb{R}^o \quad (4)$$

GPs are essentially defined by a **Training Set** and a **Kernel Function**. The **Training Set** is a collection of points pertaining to \mathcal{I} , for which the corresponding output inside \mathcal{O} is known¹.

The concept of **Kernel Function** is extensively detailed at Section 4.

2 Scalar case

The training set of a scalar **GP** is formally defined in this way:

$$S = \left\langle \begin{bmatrix} X^1 \\ y^1 \end{bmatrix} \cdots \begin{bmatrix} X^N \\ y^N \end{bmatrix} \right\rangle \quad (5)$$

$$\mathcal{X}^S = \{X^1 \dots X^N\} \subseteq \mathcal{I} \quad (6)$$

$$\mathcal{Y}^S = \{y^1 \dots y^N\} \subseteq \mathbb{R} \quad (7)$$

GPs consider values inside the training set to be somehow correlated, as they were generated from the same underlying function. In particular, the joint probability distribution describing such correlation is assumed to be a **Gaussian Distribution**:

$$\begin{bmatrix} y^1 \\ \vdots \\ y^N \end{bmatrix} \sim \mathcal{N}\left(0, K(\mathcal{X}^S)\right) \quad (8)$$

K is the kernel covariance, whose values depend on the definition of a kernel function k (refer to Section 4):

$$K = K(\mathcal{X}^S, \Theta) = \begin{bmatrix} k(X^1, X^1, \Theta) & \dots & k(X^1, X^N, \Theta) \\ \vdots & \ddots & \vdots \\ k(X^N, X^1, \Theta) & \dots & k(X^N, X^N, \Theta) \end{bmatrix} \quad (9)$$

where Θ is a vector of hyperparameters that can be tuned by training (see Sections 2.2 and 3.2) the model over a specific training set:

$$\Theta = [\theta_1 \quad \dots \quad \theta_m]^T \quad (10)$$

¹To be precise, the exact value might be unknown due to noise, but a also a close one is fine.

2.1 Predictions

The aim of **GPs** is to be able to make predictions about the output value $y(X)$ of an input X that is outside the training set. This is done assuming a joint correlation between such point and the ones in the training set:

$$\begin{bmatrix} y(X) \\ y^1 \\ \vdots \\ y^N \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k(X, X, \Theta) & K_x^T \\ K_x(X, \mathcal{X}^S, \Theta) & K(\mathcal{X}^S, \Theta) \end{bmatrix}\right) \quad (11)$$

where K_x is a vector obtained in the following way:

$$K_x(X, \mathcal{X}^S, \Theta) = [k(X, X^1, \Theta) \quad \dots \quad k(X, X^N, \Theta)]^T \quad (12)$$

As the joint distribution is **Gaussian**, the conditioned distribution can be obtained as follows:

$$y(X|\mathcal{X}^S) \sim \mathcal{N}\left(K_x^T K^{-1} \begin{bmatrix} y^1 \\ \vdots \\ y^N \end{bmatrix}, k(X, X) - K_x^T K^{-1} K_x\right) \quad (13)$$

2.2 Training

Training is done maximizing the likelihood L of the training set w.r.t. Θ . Since eq. (8) describes a **Gaussian** distribution, the likelihood can be computed as follows:

$$L(\mathcal{Y}^S) = \frac{1}{\sqrt{(2\pi)^N |K(\mathcal{X}^S)|}} \exp\left(-\frac{1}{2} [y^1 \quad \dots \quad y^N] K^{-1} \begin{bmatrix} y^1 \\ \vdots \\ y^N \end{bmatrix}\right) \quad (14)$$

At this point, the property described in appendix A can be exploited to rewrite the above equation as follows:

$$L(\mathcal{Y}^S) = \frac{1}{\sqrt{(2\pi)^N |K|}} \exp\left(-\frac{1}{2} \text{Tr}\left[K^{-1} \begin{bmatrix} y^1 \\ \vdots \\ y^N \end{bmatrix} [y^1 \quad \dots \quad y^N]\right]\right) \quad (15)$$

$$= \frac{1}{\sqrt{(2\pi)^N |K|}} \exp\left(-\frac{1}{2} \text{Tr}\left[K^{-1} M_{YY}\right]\right) \quad (16)$$

Passing to the logarithm we obtain:

$$\mathcal{L} = \log(L) = -\frac{N}{2} (2\pi) - \frac{1}{2} \log(|K|) - \frac{1}{2} \text{Tr}\left[K^{-1} M_{YY}\right] \quad (17)$$

keeping in mind that \mathcal{L} is a function of the hyperparameters Θ :

$$\mathcal{L}(\Theta) = -\frac{N}{2}(2\pi) - \frac{1}{2}\log(|K(\Theta)|) - \frac{1}{2}\text{Tr}\left[K(\Theta)^{-1}M_{YY}\right] \quad (18)$$

The gradient of \mathcal{L} w.r.t. the generic hyperparameter θ_t is computed as follows (refer to the properties detailed at [2]):

$$\frac{\partial \mathcal{L}}{\partial \theta_t} = -\frac{1}{2}\text{Tr}\left[K^{-1}\frac{\partial K}{\partial \theta_t}\right] - \frac{1}{2}\text{Tr}\left[\frac{\partial}{\partial \theta_t}(K^{-1}M_{YY})\right] \quad (19)$$

$$= -\frac{1}{2}\text{Tr}\left[K^{-1}\frac{\partial K}{\partial \theta_t}\right] - \frac{1}{2}\text{Tr}\left[\frac{\partial(K^{-1})}{\partial \theta_t}M_{YY}\right] \quad (20)$$

$$= -\frac{1}{2}\text{Tr}\left[K^{-1}\frac{\partial K}{\partial \theta_t}\right] + \frac{1}{2}\text{Tr}\left[K^{-1}\frac{\partial K}{\partial \theta_t}K^{-1}M_{YY}\right] \quad (21)$$

$$= \frac{1}{2}\text{Tr}\left[K^{-1}\frac{\partial K}{\partial \theta_t}\left(K^{-1}M_{YY} - I_N\right)\right] \quad (22)$$

Choosing your favourite gradient-based approach you can then tune the model.

3 Vectorial case

Vectorial **GP**s are defined as similarly done for scalar **GP**s. The training set should account for the multi-dimensionality of the process and is therefore defined in this way:

$$S = \left\langle \begin{bmatrix} X^1 \\ y_1^1 \\ \vdots \\ y_o^1 \end{bmatrix} \cdots \begin{bmatrix} X^N \\ y_1^N \\ \vdots \\ y_o^N \end{bmatrix} \right\rangle = \left\langle \begin{bmatrix} X^1 \\ Y^1 \end{bmatrix} \cdots \begin{bmatrix} X^N \\ Y^N \end{bmatrix} \right\rangle \quad (23)$$

3.1 Predictions

A vectorial **GP** is actually a composition of independent scalar **GP**s. The prediction is done as similarly discussed in Section 2.1, doing o predictions at the same time. Indeed, for each $i \in \{0, \dots, o\}$ holds that:

$$\begin{bmatrix} y_i(X) \\ y_i^1 \\ \vdots \\ y_i^N \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k(X, X, \Theta) & K_x^T \\ K_x(X, \mathcal{X}^S, \Theta) & K(\mathcal{X}^S, \Theta) \end{bmatrix}\right) \quad (24)$$

Then, the complete prediction is obtained in this way:

$$Y(X|X^S, \Theta) = \begin{bmatrix} y_1(X|X^S, \Theta) \\ \vdots \\ y_o(X|X^S, \Theta) \end{bmatrix} \quad (25)$$

$$\sim \begin{bmatrix} \mathcal{N}\left(K_x^T K^{-1} \begin{bmatrix} y_1^1 \\ \vdots \\ y_1^N \end{bmatrix}, k(X, X) - K_x^T K^{-1} K_x\right) \\ \vdots \\ \mathcal{N}\left(K_x^T K^{-1} \begin{bmatrix} y_o^1 \\ \vdots \\ y_o^N \end{bmatrix}, k(X, X) - K_x^T K^{-1} K_x\right) \end{bmatrix} \quad (26)$$

$$\sim \mathcal{N}\left(\left(K_x^T K^{-1} \begin{bmatrix} y_1^1 & \cdots & y_o^1 \\ \vdots & \ddots & \vdots \\ y_1^N & \cdots & y_o^N \end{bmatrix}\right)^T, K_x^T K^{-1} K_x I_{o,o}\right) \quad (27)$$

3.2 Training

The logarithmic likelihood is the summation of the logarithmic likelihood of each process that compose the vectorial \mathbf{GP} , which leads, omitting constant terms, to:

$$\mathcal{L} = \sum_{i=0}^o \left(-\frac{1}{2} \log(|K|) - \frac{1}{2} \text{Tr} \left[K^{-1} \begin{bmatrix} y_i^1 \\ \vdots \\ y_i^N \end{bmatrix} \begin{bmatrix} y_i^1 & \cdots & y_i^N \end{bmatrix} \right] \right) \quad (28)$$

$$= -\frac{o}{2} \log(|K|) - \frac{1}{2} \sum_{i=0}^o \left(\text{Tr} \left[K^{-1} \begin{bmatrix} y_i^1 \\ \vdots \\ y_i^N \end{bmatrix} \begin{bmatrix} y_i^1 & \cdots & y_i^N \end{bmatrix} \right] \right) \quad (29)$$

$$= -\frac{o}{2} \log(|K|) - \frac{1}{2} \text{Tr} \left[K^{-1} \sum_{i=0}^o \left(\begin{bmatrix} y_i^1 \\ \vdots \\ y_i^N \end{bmatrix} \begin{bmatrix} y_i^1 & \cdots & y_i^N \end{bmatrix} \right) \right] \quad (30)$$

$$= -\frac{o}{2} \log(|K|) - \frac{1}{2} \text{Tr} \left[K^{-1} M_{YY}^o \right] \quad (31)$$

with:

$$M_{YY}^o = \sum_{i=0}^o \left(\begin{bmatrix} y_i^1 \\ \vdots \\ y_i^N \end{bmatrix} \begin{bmatrix} y_i^1 & \cdots & y_i^N \end{bmatrix} \right) = \begin{bmatrix} \langle Y^1, Y^1 \rangle & \cdots & \langle Y^1, Y^N \rangle \\ \vdots & \ddots & \vdots \\ \langle Y^N, Y^1 \rangle & \cdots & \langle Y^N, Y^N \rangle \end{bmatrix} \quad (32)$$

keeping again in mind that \mathcal{L} is a function of the hyperparameters Θ :

$$\mathcal{L}(\Theta) = -\frac{o}{2} \log(|K(\Theta)|) - \frac{1}{2} \text{Tr} \left[K^{-1}(\Theta) M_{YY}^o \right] \quad (33)$$

The gradient can be computed with the same steps that led to eq. (22), leading to:

$$\frac{1}{2} \text{Tr} \left[K^{-1} \frac{\partial K}{\partial \theta_t} \left(K^{-1} M_{YY}^o - o I_N \right) \right] \quad (34)$$

4 The Kernel function

The kernel function describes correlation between inputs. Ideally, it should assume a low value for inputs that are "far", w.r.t. a certain metrics, from each other and high values for those inputs that are close.

The kernel function k should be designed in order to produce a symmetric positive definite matrix K , as this latter should be representative of a covariance matrix, see eq. (8).

Clearly, the gradient of K can be computed element by element:

$$\frac{\partial K}{\partial \theta_t} = \begin{bmatrix} \frac{\partial k(X^1, X^1)}{\partial \theta_t} & \cdots & \frac{\partial k(X^1, X^N)}{\partial \theta_t} \\ \vdots & \ddots & \vdots \\ \frac{\partial k(X^N, X^1)}{\partial \theta_t} & \cdots & \frac{\partial k(X^N, X^N)}{\partial \theta_t} \end{bmatrix} \quad (35)$$

In the following of this Section, the most popular kernel functions ² will be discussed.

4.1 Linear function

Hyperparameters:

$$\Theta = [\theta_0 \quad \theta_1 \quad \mu_1 \quad \dots \quad \mu_o] \quad (36)$$

Fuction evaluation:

$$k(x, y, \Theta) = \theta_0^2 + \theta_1^2 (x - \mu)^T (x - \mu) \quad (37)$$

$$= \theta_0^2 + \theta_1^2 x^T y + \theta_1^2 \mu^T \mu - \mu^T \theta_1^2 (x + y) \quad (38)$$

²Which are also the ones default supported by this package.

Function gradient:

$$\frac{\partial k(x, y)}{\partial \theta_0} = 2\theta_0 \quad (39)$$

$$\frac{\partial k(x, y)}{\partial \theta_1} = 2\theta_1(x - \mu)(y - \mu) \quad (40)$$

$$\begin{bmatrix} \frac{\partial k(x, y)}{\partial \mu_1} \\ \vdots \\ \frac{\partial k(x, y)}{\partial \mu_o} \end{bmatrix} = 2\theta_1^2 \mu - \theta_1^2(x + y) \quad (41)$$

4.2 Squared exponential

Hyperparameters:

$$\Theta = [\theta_0 \quad \theta_1] \quad (42)$$

Fuction evaluation:

$$k(x, y, \Theta) = \theta_0^2 \exp(-\theta_1^2 \|x - y\|_2^2) \quad (43)$$

$$= \theta_0^2 \exp(-\theta_1^2 (x - y)^T (x - y)) \quad (44)$$

Function gradient:

$$\frac{\partial k(x, y)}{\partial \theta_0} = 2\theta_0 \exp(-\theta_1^2 \|x - y\|_2^2) \quad (45)$$

$$\frac{\partial k(x, y)}{\partial \theta_1} = \theta_0^2 \exp(-\theta_1^2 \|x - y\|_2^2) (-2\theta_1 (x - y)^T (x - y)) \quad (46)$$

5 What to expect from the predictions

TODO spiegare che incertezza aumenta tanto piu sonon lontano da punti in training set

A Trace property

Take an (n, n) matrix A and a vector X , the scalar quantity $x^T A x$ is equal to:

$$x^T A x = \text{Tr} \begin{bmatrix} A x x^T \end{bmatrix} \quad (47)$$

$$= \text{Tr} \begin{bmatrix} x x^T A^T \end{bmatrix} \quad (48)$$

Clearly, in case of symmetric matrix, the following holds:

$$x^T A x = \text{Tr} \begin{bmatrix} x x^T A \end{bmatrix} \quad (49)$$

We will now prove equation (47).
 $x^T Ax$ can be also expressed as follows:

$$x^T Ax = x^T \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} x \quad (50)$$

$$= x^T \begin{bmatrix} a_1^T x \\ \vdots \\ a_n^T x \end{bmatrix} = x^T \begin{bmatrix} \langle a_1, x \rangle \\ \vdots \\ \langle a_n, x \rangle \end{bmatrix} \quad (51)$$

$$= \sum_{i=0}^n x_i \langle a_i, x \rangle \quad (52)$$

where a_i is the i^{th} row of A . At the same time, the following fact is also true:

$$Tr[Axx^T] = Tr \left[\begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} [xx_1 \quad \dots \quad xx_n] \right] \quad (53)$$

$$= Tr \left[\begin{bmatrix} a_1^T xx_1 & & \\ & \ddots & \\ & & a_n^T xx_n \end{bmatrix} \right] \quad (54)$$

$$= \sum_{i=0}^n x_i \langle a_i, x \rangle \quad (55)$$

where we recognize that eq. (52) and (55) are identical.

References

- [1] Richard A. Davis. Gaussian process: Theory, 2014.
- [2] K. B. Petersen and M. S. Pedersen. The matrix cookbook, 2008.
- [3] Carl Edward Rasmussen. Gaussian processes in machine learning. In *Summer school on machine learning*, pages 63–71. Springer, 2003.