

Gaussian Processes

Andrea Casalino

June 17, 2022

1 What is a Gaussian Process?

Gaussian Processes [3], [1], a.k.a. **GPs**, are data driven probabilistic models able to approximate generic multivariate and possibly vectorial real functions G defined like that:

$$G : \mathcal{I} \subseteq \mathbb{R}^i \rightarrow \mathcal{O} \subseteq \mathbb{R}^o \quad (1)$$

In essence, **GPs** are defined by their own **Training Set** and **Kernel Function**. The **Training Set** is a collection of points pertaining to \mathcal{I} , for which the corresponding output, or at least that value summed with noise, inside \mathcal{O} is known.

The **Kernel Function** is something that has to be chosen and typical of the kind of function to approximate. Definitions and meanings of possible **Kernel Function** are extensively detailed in Section 4.

The aim of **GPs** is to be able to predict the value of G for a point inside \mathcal{I} that is not inside of the training set. In particular, this is done in probabilistic terms, as the result of the prediction is not just a value, but a conditioned Gaussian distribution.

Section 2 discusses the formulation of scalar **GP**, i.e. cases for which $\mathcal{O} \subseteq \mathbb{R}$. Instead, Section 3 focuses on the more general cases. The reader will notice that the two formulations are not in contradiction and the decision to discuss them into separate Sections is only for the purpose of a better readability.

2 Scalar case

The training set of a scalar **GP** is a collection of points inside \mathcal{I} for which the corresponding value inside \mathcal{O} is known. More formally:

$$S = \left\{ \begin{bmatrix} X_1 \\ \vdots \\ y_1 \end{bmatrix}, \dots, \begin{bmatrix} X_N \\ \vdots \\ y_N \end{bmatrix} \right\} \quad (2)$$

$$\{X_1, \dots, X_N\} \subset \mathbb{R}^i \quad (3)$$

$$\{y_1, \dots, y_N\} \subset \mathbb{R} \quad (4)$$

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 the following 0 mean **Gaussian Distribution**:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \sim \mathcal{N}\left(0, K(X_{1,\dots,N})\right) \quad (5)$$

K is a covariance matrix induced by the choice of a certain kernel function

k (refer also to Section 4):

$$K = \begin{bmatrix} k(X_1, X_1, \Theta) & \dots & k(X_1, X_N, \Theta) \\ \vdots & \ddots & \vdots \\ k(X_N, X_1, \Theta) & \dots & k(X_N, X_N, \Theta) \end{bmatrix} \quad (6)$$

where Θ is a vector of hyperparameters that are typical of the chosen kernel function and whose values can be tuned also by training (see Sections 2.2 and 3.2):

$$\Theta = [\theta_1 \quad \dots \quad \theta_m]^T \quad (7)$$

It is generally a good practice to add to the covariance K an additional term modelling noise in the training set. This is done by simply summing an isotropic standard deviation σ_{noise} :

$$K' = K + \sigma_{noise} I_{N \times N} \quad (8)$$

2.1 Predictions

The aim of a **GP** is to be able to make predictions about the output value $y = G(X)$ pertaining o an input X that is outside of the training set. This is done assuming again a joint Gaussian correlation between such an additional point and all the ones in the training set:

$$\begin{bmatrix} y = G(X) \\ y_1 \\ \vdots \\ y_N \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k(X, X, \Theta) & K_x(X, X_{1,\dots,N}, \Theta) \\ K_x(X, X_{1,\dots,N}, \Theta)^T & K \end{bmatrix} \right) \quad (9)$$

where K_x is a vector assembled in this way:

$$K_x(X, X_{1,\dots,N}, \Theta) = [k(X, X_1, \Theta) \quad \dots \quad k(X, X_N, \Theta)]^T \quad (10)$$

Since eq. 10 describes a Gaussian distribution, the conditioned distribution involving only X can be obtained as follows:

$$y(X|X_{1,\dots,N}) \sim \mathcal{N} \left(K_x^T K^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \sigma_K(X) \right) \quad (11)$$

where $\sigma_K(X)$ is the covariance of the conditioned distribution and can be computed as follows:

$$\sigma_K(X) = k(X, X) - K_x^T K^{-1} K_x \quad (12)$$

The distribution described by eq. 11, actually represents the prediction made by the **GP**. Eq. 11 can also be rewritten as follows:

$$y(X|X_{1,\dots,N}) \sim \mathcal{N} \left([y_1 \quad \dots \quad y_N] K^{-1} K_x, \sigma_K(X) \right) \quad (13)$$

2.2 Training

Training is done maximizing the likelihood L of the training set w.r.t. the hyperparameters Θ of the kernel function. Since eq. (5) describes a **Gaussian** distribution, the likelihood can be computed as follows:

$$L = \frac{1}{\sqrt{(2\pi)^N |K|}} \exp\left(-\frac{1}{2} [y_1 \ \dots \ y_N] K^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}\right) \quad (14)$$

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

$$L = \frac{1}{\sqrt{(2\pi)^N |K|}} \exp\left(-\frac{1}{2} \text{Tr}\left[K^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} [y_1 \ \dots \ y_N]\right]\right) \quad (15)$$

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

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]):

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

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

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

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

Choosing your favourite gradient-based approach you can tune the model, byt computing the gradient as described by the above equation.

Check Section 3.2.1 to understand how priors about the hyperparameters Θ can be handled.

3 Vectorial case

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

$$S = \left\{ \begin{bmatrix} X_1 \\ \downarrow \\ Y_1 \end{bmatrix}, \dots, \begin{bmatrix} X_N \\ \downarrow \\ Y_N \end{bmatrix} \right\} \quad (23)$$

$$\{X_1, \dots, X_N\} \subset \mathbb{R}^i \quad (24)$$

$$\{Y_1, \dots, Y_N\} \subset \mathbb{R}^o \quad (25)$$

The generic Y_k is a vector made of o components:

$$Y_k = [y_k^1 \quad \dots \quad y_k^o]^T \quad (26)$$

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 component $y^{k \in \{1, \dots, o\}}$ of the prediction holds equation 11. Therefore, the complete prediction can be obtained in this way:

$$Y(X|X_{1,\dots,N}) = \begin{bmatrix} \mathcal{N}\left([y_1^1 \quad \dots \quad y_N^1] K^{-1} K_x, \sigma_K(X)\right) \\ \vdots \\ \mathcal{N}\left([y_1^o \quad \dots \quad y_N^o] K^{-1} K_x, \sigma_K(X)\right) \end{bmatrix} \quad (27)$$

the above expression can be further elaborated, leading to the distribution of an isotropic Gaussian:

$$Y(X|X_{1,\dots,N}) \sim \mathcal{N}\left(\begin{bmatrix} y_1^1 & \dots & y_N^1 \\ \vdots & \ddots & \vdots \\ y_1^o & \dots & y_N^o \end{bmatrix} K^{-1} K_x, \sigma_K(X) I_{o \times o}\right) \quad (28)$$

$$\sim \mathcal{N}\left([Y_1 | \quad \dots \quad | Y_N] K^{-1} K_x, \sigma_K(X) I_{o \times o}\right) \quad (29)$$

3.2 Training

The logarithmic likelihood is the summation of the logarithmic likelihood of each process that compose the vectorial **GP**, which leads, omitting constant

terms, to:

$$\mathcal{L} = \sum_{k=1}^o \left(-\frac{1}{2} \log(|K|) - \frac{1}{2} \text{Tr} \left[K^{-1} \begin{bmatrix} y_1^i \\ \vdots \\ y_N^i \end{bmatrix} [y_1^i \dots y_N^i] \right] \right) \quad (30)$$

$$= -\frac{o}{2} \log(|K|) - \frac{1}{2} \sum_{k=1}^o \left(\text{Tr} \left[K^{-1} \begin{bmatrix} y_1^i \\ \vdots \\ y_N^i \end{bmatrix} [y_1^i \dots y_N^i] \right] \right) \quad (31)$$

$$= -\frac{o}{2} \log(|K|) - \frac{1}{2} \text{Tr} \left[K^{-1} \sum_{k=1}^o \left(\begin{bmatrix} y_1^i \\ \vdots \\ y_N^i \end{bmatrix} [y_1^i \dots y_N^i] \right) \right] \quad (32)$$

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

with:

$$M_{YY}^o = \sum_{k=1}^o \left(\begin{bmatrix} y_1^i \\ \vdots \\ y_N^i \end{bmatrix} [y_1^i \dots y_N^i] \right) \quad (34)$$

$$= \sum_{k=1}^o \begin{bmatrix} y_1^i y_1^i & \dots & y_1^i y_N^i \\ \vdots & \ddots & \vdots \\ y_N^i y_1^i & \dots & y_N^i y_N^i \end{bmatrix} \quad (35)$$

$$= \begin{bmatrix} \sum_{k=1}^o y_1^i y_1^i & \dots & \sum_{k=1}^o y_1^i y_N^i \\ \vdots & \ddots & \vdots \\ \sum_{k=1}^o y_N^i y_1^i & \dots & \sum_{k=1}^o y_N^i y_N^i \end{bmatrix} \quad (36)$$

$$= \begin{bmatrix} \langle Y_1, Y_1 \rangle & \dots & \langle Y_1, Y_N \rangle \\ \vdots & \ddots & \vdots \\ \langle Y_N, Y_1 \rangle & \dots & \langle Y_N, Y_N \rangle \end{bmatrix} \quad (37)$$

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

$$\mathcal{L}(\Theta) = -\frac{o}{2} \log(|K(\Theta)|) - \frac{1}{2} \text{Tr} \left[K^{-1}(\Theta) M_{YY}^o \right] \quad (38)$$

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

$$\frac{\partial \mathcal{L}}{\partial \theta_t} = \frac{1}{2} \text{Tr} \left[K^{-1} \frac{\partial K}{\partial \theta_t} \left(K^{-1} M_{YY}^o - o I_{N \times N} \right) \right] \quad (39)$$

3.2.1 Hyperparameters prior knowledge

It is possible to also take into account an a-priori knowledge of the process, in terms of a prior distribution describing the hyperparameters values. Such priors

modify the likelihood function to optimize in this way:

$$L'(\Theta) = L(\Theta)L(\Theta)_{prior} \quad (40)$$

In principle, any distribution can be used as a prior for the hyperparameters. Without loss of generality, assume the prior is modelled with a multivariate gaussian distribution defined by a certain mean μ_{prior} and covariance Σ_{prior} . In such a case, equation (40) became:

$$L'(\Theta) \propto L(\Theta) \exp\left(-\frac{1}{2}(\Theta - \mu_{prior})^T \Sigma_{prior}^{-1} (\Theta - \mu_{prior})\right) \quad (41)$$

passing to the logarithm and neglecting constant values we get:

$$\mathcal{L}'(\Theta) = \mathcal{L}(\Theta) - \frac{1}{2}(\Theta - \mu_{prior})^T \Sigma_{prior}^{-1} (\Theta - \mu_{prior}) \quad (42)$$

Clearly, the gradient of the complete likelihood $\mathcal{L}'(\Theta)$ can be computed as follows:

$$\frac{\partial \mathcal{L}'}{\partial \Theta} = \frac{\partial \mathcal{L}}{\partial \Theta} - \Sigma_{prior}^{-1} (\Theta - \mu_{prior}) \quad (43)$$

4 Kernel functions

The kernel function describes the correlation between the 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. (5).

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

$$\frac{\partial K}{\partial \theta_t} = \begin{bmatrix} \frac{\partial k(X_1, X_1)}{\partial \theta_t} & \dots & \frac{\partial k(X_1, X_N)}{\partial \theta_t} \\ \vdots & \ddots & \vdots \\ \frac{\partial k(X_N, X_1)}{\partial \theta_t} & \dots & \frac{\partial k(X_N, X_N)}{\partial \theta_t} \end{bmatrix} \quad (44)$$

In the following of this Section, some of the most popular kernel functions ¹ will be discussed.

4.1 Linear function

Hyperparameters:

$$\Theta = [\theta_0 \quad \theta_1 \quad \mu_1 \quad \dots \quad \mu_o] \quad (45)$$

¹Which are also the ones default supported by this package.

Fuction evaluation:

$$k(a, b, \Theta) = \theta_0^2 + \theta_1^2 (a - \mu)^T (b - \mu) \quad (46)$$

$$= \theta_0^2 + \theta_1^2 (\langle \mu, \mu \rangle + \langle a, b \rangle - \langle \mu, a + b \rangle) \quad (47)$$

Function gradient:

$$\frac{\partial k(a, b)}{\partial \theta_0} = 2\theta_0 \quad (48)$$

$$\frac{\partial k(a, b)}{\partial \theta_1} = 2\theta_1 (a - \mu)(b - \mu) \quad (49)$$

$$\begin{bmatrix} \frac{\partial k(a, b)}{\partial \mu_1} \\ \vdots \\ \frac{\partial k(a, b)}{\partial \mu_o} \end{bmatrix} = \theta_1^2 (2\mu - a - b) \quad (50)$$

4.2 Squared exponential

Hyperparameters:

$$\Theta = [\sigma \quad d] \quad (51)$$

Fuction evaluation:

$$k(a, b, \Theta) = \sigma^2 \exp\left(-\frac{\|a - b\|_2^2}{d^2}\right) \quad (52)$$

Function gradient:

$$\frac{\partial k(a, b)}{\partial \sigma} = 2\sigma \exp\left(-\frac{\|a - b\|_2^2}{d^2}\right) \quad (53)$$

$$\frac{\partial k(a, b)}{\partial d} = \sigma^2 \exp\left(-\frac{\|a - b\|_2^2}{d^2}\right) 2 \frac{\|a - b\|_2^2}{d^3} \quad (54)$$

4.3 Periodic kernel

This function is able to well express the periodicity in the function to approximate. The hyperparameters are:

$$\Theta = [\sigma \quad d \quad p] \quad (55)$$

p represents the period of the kernel function.

Fuction evaluation:

$$k(a, b, \Theta) = \sigma^2 \exp\left(-\frac{\sin^2\left(\frac{2\pi}{p} \|a - b\|_2\right)}{d^2}\right) = \sigma^2 \exp\left(-\frac{\sin^2(\alpha)}{d^2}\right) \quad (56)$$

Function gradient:

$$\frac{\partial k(a, b)}{\partial \sigma} = 2\sigma \exp\left(-\frac{\sin^2(\alpha)}{d^2}\right) \quad (57)$$

$$\frac{\partial k(a, b)}{\partial d} = \sigma^2 \exp\left(-\frac{\sin^2(\alpha)}{d^2}\right) 2\frac{\sin^2(\alpha)}{d^3} \quad (58)$$

$$\begin{aligned} \frac{\partial k(a, b)}{\partial p} &= \sigma^2 \exp\left(-\frac{\sin^2(\alpha)}{d^2}\right) \left(-\frac{2\sin(\alpha)\cos(\alpha)}{d^2}\right) \frac{\partial \alpha}{\partial p} \\ &= \sigma^2 \exp\left(-\frac{\sin^2(\alpha)}{d^2}\right) \left(-\frac{2\sin(\alpha)\cos(\alpha)}{d^2}\right) \left(-\frac{2\pi \|a-b\|_2}{p^2}\right) \\ &= \sigma^2 \exp\left(-\frac{\sin^2(\alpha)}{d^2}\right) \frac{2\sin(\alpha)\cos(\alpha)\alpha}{pd^2} \end{aligned} \quad (59)$$

4.4 Combining kernel functions

Clearly, you can also combine the kernel functions described in this Section to create more complex one.

4.4.1 Summation

You can sum more kernel functions $k_{1,\dots,m}$ together, each having its own group of hyperparameters $\Theta_{1,\dots,m}$:

$$\Theta = [\Theta_1^T \ \dots \ \Theta_m^T]^T \quad (60)$$

$$k(a, b, \Theta) = \sum_{i=1}^m k_i(a, b, \Theta_i) \quad (61)$$

The gradient is clearly computed as follows:

$$\frac{\partial k(a, b)}{\partial \Theta} = \begin{bmatrix} \frac{\partial k_1(a, b)}{\partial \Theta_1}^T \\ \vdots \\ \frac{\partial k_m(a, b)}{\partial \Theta_m}^T \end{bmatrix} \quad (62)$$

4.4.2 Product

You can multiply bunch of kernel functions $k_{1,\dots,m}$ together, each having its own group of hyperparameters $\Theta_{1,\dots,m}$:

$$\Theta = [\Theta_1^T \ \dots \ \Theta_m^T]^T \quad (63)$$

$$k(a, b, \Theta) = \prod_{i=1}^m k_i(a, b, \Theta_i) \quad (64)$$

The gradient is clearly computed as follows:

$$\frac{\partial k(a, b)}{\partial \Theta} = \left(\prod_{i=1}^m k_i(a, b) \right) \begin{bmatrix} \frac{\partial k_1(a, b)}{\partial \Theta_1} \frac{1}{k_1(a, b)} \\ \vdots \\ \frac{\partial k_m(a, b)}{\partial \Theta_m} \frac{1}{k_m(a, b)} \end{bmatrix} \quad (65)$$

A Trace property

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

$$x^T A x = \text{Tr} \left[A x x^T \right] \quad (66)$$

$$= \text{Tr} \left[x x^T A \right] \quad (67)$$

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

$$x^T A x = \text{Tr} \left[x x^T A \right] \quad (68)$$

We will now prove equation (66).
 $x^T A x$ can be also expressed as follows:

$$x^T A x = x^T \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} x \quad (69)$$

$$= x^T \begin{bmatrix} a_1^T x \\ \vdots \\ a_n^T x \end{bmatrix} = x^T \begin{bmatrix} \langle a_1, x \rangle \\ \vdots \\ \langle a_n, x \rangle \end{bmatrix} \quad (70)$$

$$= \sum_{i=1}^n x_i \langle a_i, x \rangle \quad (71)$$

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

$$\text{Tr} \left[A x x^T \right] = \text{Tr} \left[\begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} \begin{bmatrix} x x_1 & \dots & x x_n \end{bmatrix} \right] \quad (72)$$

$$= \text{Tr} \left[\begin{bmatrix} a_1^T x x_1 & & \\ & \ddots & \\ & & a_n^T x x_n \end{bmatrix} \right] \quad (73)$$

$$= \sum_{i=1}^n x_i \langle a_i, x \rangle \quad (74)$$

where we recognize that eq. (71) and (74) 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.