

Learning from data: Gaussian processes

Christian Forssén

Department of Physics, Chalmers University of Technology, Sweden

Oct 4, 2020

1 Inference using Gaussian processes

Assume that there is a set of input vectors with independent, predictor, variables

$$\mathbf{X}_N \equiv \{\mathbf{x}^{(i)}\}_{i=1}^N$$

and a set of target values

$$\mathbf{t}_N \equiv \{t^{(i)}\}_{i=1}^N.$$

- Note that we will use the symbol t to denote the target, or response, variables in the context of Gaussian Processes. Here we will consider single, scalar outputs $t^{(i)}$. The extension to vector outputs $\mathbf{t}^{(i)}$ is straightforward.
- Furthermore, we will use the subscript N to denote a set of N vectors (or scalars): \mathbf{X}_N (\mathbf{t}_N),
- ... while a single instance i is denoted by a superscript: $\mathbf{x}^{(i)}$ ($t^{(i)}$).

We will consider two different *inference problems*:

1. The prediction of a *new target* $t^{(N+1)}$ given the data $\mathbf{X}_N, \mathbf{t}_N$ and a new input $\mathbf{x}^{(N+1)}$.
2. The inference of a *model function* $y(\mathbf{x})$ from the data $\mathbf{X}_N, \mathbf{t}_N$.

The former can be expressed with the pdf

$$p\left(t^{(N+1)}|\mathbf{t}_N, \mathbf{X}_N, \mathbf{x}^{(N+1)}\right)$$

while the latter can be written using Bayes' formula (in these notes we will not be including information I explicitly in the conditional probabilities)

$$p(y(\mathbf{x})|\mathbf{t}_N, \mathbf{X}_N) = \frac{p(\mathbf{t}_N|y(\mathbf{x}), \mathbf{X}_N)p(y(\mathbf{x}))}{p(\mathbf{t}_N|\mathbf{X}_N)}$$

The inference of a function will obviously also allow to make predictions for new targets. However, we will need to consider in particular the second term in the numerator, which is the **prior** distribution on functions assumed in the model.

- This prior is implicit in parametric models with priors on the parameters.
- The idea of Gaussian process modeling is to put a prior directly on the **space of functions** without parameterizing $y(\mathbf{x})$.
- A Gaussian process can be thought of as a generalization of a Gaussian distribution over a finite vector space to a **function space of infinite dimension**.
- Just as a Gaussian distribution is specified by its mean and covariance matrix, a Gaussian process is specified by a **mean and covariance function**.

Gaussian process

A Gaussian process is a stochastic process (a collection of random variables indexed by time or space), such that every finite collection of those random variables has a multivariate normal distribution

References:

1. [Gaussian Processes for Machine Learning](#), Carl Edward Rasmussen and Chris Williams, the MIT Press, 2006, [online version](#).
2. [GPy](#): a Gaussian Process (GP) framework written in python, from the Sheffield machine learning group.

1.1 Parametric approach

Let us express $y(\mathbf{x})$ in terms of a model function $y(\mathbf{x}; \boldsymbol{\theta})$ that depends on a vector of model parameters $\boldsymbol{\theta}$.

For example, using a set of basis functions $\{\phi^{(h)}(\mathbf{x})\}_{h=1}^H$ with linear weights $\boldsymbol{\theta}_H$ we have

$$y(\mathbf{x}, \boldsymbol{\theta}) = \sum_{h=1}^H \theta^{(h)} \phi^{(h)}(\mathbf{x})$$

Notice

The basis functions can be non-linear in \mathbf{x} such as Gaussians (aka *radial basis functions*)

$$\phi^{(h)}(\mathbf{x}) = \exp \left[-\frac{(\mathbf{x} - \mathbf{c}^{(h)})^2}{2(\sigma^{(h)})^2} \right].$$

Still, this constitutes a linear model since $y(\mathbf{x}, \boldsymbol{\theta})$ depends linearly on the parameters $\boldsymbol{\theta}$.

The inference of model parameters should be a well-known problem by now. We state it in terms of Bayes theorem

$$p(\boldsymbol{\theta} | \mathbf{t}_N, \mathbf{X}_N) = \frac{p(\mathbf{t}_N | \boldsymbol{\theta}, \mathbf{X}_N) p(\boldsymbol{\theta})}{p(\mathbf{t}_N | \mathbf{X}_N)}$$

Having solved this inference problem (note that the likelihood is Gaussian, cf linear regression) a prediction can be made through marginalization

$$p(t^{(N+1)} | \mathbf{t}_N, \mathbf{X}_N, \mathbf{x}^{(N+1)}) = \int d^H \boldsymbol{\theta} p(t^{(N+1)} | \boldsymbol{\theta}, \mathbf{x}^{(N+1)}) p(\boldsymbol{\theta} | \mathbf{t}_N, \mathbf{X}_N).$$

Here it is important to note that the final answer does not make any explicit reference to our parametric representation of the unknown function $y(\mathbf{x})$.

Assuming that we have a fixed set of basis functions and Gaussian prior distributions (with zero mean) on the weights $\boldsymbol{\theta}$ we will show that:

- The joint pdf of the observed data given the model $p(\mathbf{t}_N | \mathbf{X}_N)$, is a multivariate Gaussian with mean zero and with a covariance matrix that is determined by the basis functions.
- This implies that the conditional distribution $p(t^{(N+1)} | \mathbf{t}_N, \mathbf{X}_{N+1})$, is also a multivariate Gaussian whose mean depends linearly on \mathbf{t}_N .

Proof.

Sum of normally distributed random variables

If X and Y are independent random variables that are normally distributed (and therefore also jointly so), then their sum is also normally distributed. i.e., $Z = X + Y$ is normally distributed with its mean being the sum of the two means, and its variance being the sum of the two variances.

Consider the linear model and define the $N \times H$ design matrix \mathbf{R} with elements

$$R_{nh} \equiv \phi^{(h)}(\mathbf{x}^{(n)}).$$

Then $\mathbf{y}_N = \mathbf{R}\boldsymbol{\theta}$ is the vector of model predictions, i.e.

$$y^{(n)} = \sum_{h=1}^H R_{nh} \boldsymbol{\theta}^{(h)}.$$

Assume that we have a Gaussian prior for the linear model weights $\boldsymbol{\theta}$ with zero mean and a diagonal covariance matrix

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; 0, \sigma_\theta^2 \mathbf{I}).$$

Now, since y is a linear function of $\boldsymbol{\theta}$, it is also Gaussian distributed with mean zero. Its covariance matrix becomes

$$\mathbf{Q} = \langle \mathbf{y} \mathbf{y}^T \rangle = \langle \mathbf{R} \boldsymbol{\theta} \boldsymbol{\theta}^T \mathbf{R}^T \rangle = \sigma_\theta^2 \mathbf{R} \mathbf{R}^T,$$

which implies that

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; 0, \sigma_\theta^2 \mathbf{R} \mathbf{R}^T).$$

This will be true for any set of points \mathbf{X}_N ; which is the defining property of a **Gaussian process**.

- What about the target values \mathbf{t} ?

Well, if $t^{(n)}$ is assumed to differ by additive Gaussian noise, i.e.,

$$t^{(n)} = y^{(n)} + \varepsilon^{(n)},$$

where $\varepsilon^{(n)} \sim \mathcal{N}(0, \sigma_\nu^2)$; then \mathbf{t} also has a Gaussian distribution

$$p(\mathbf{t}) = \mathcal{N}(\mathbf{t}; 0, \mathbf{C}),$$

where the covariance matrix of this target distribution is given by

$$\mathbf{C} = \mathbf{Q} + \sigma_\nu^2 \mathbf{I} = \sigma_\theta^2 \mathbf{R} \mathbf{R}^T + \sigma_\nu^2 \mathbf{I}.$$

The covariance matrix as the central object. The covariance matrices are given by

$$Q_{nn'} = \sigma_\theta^2 \sum_h \phi^{(h)}(\mathbf{x}^{(n)}) \phi^{(h)}(\mathbf{x}^{(n')}),$$

and

$$C_{nn'} = Q_{nn'} + \delta_{nn'} \sigma_\nu^2.$$

This means that the correlation between target values $t^{(n)}$ and $t^{(n')}$ is determined by the points $\mathbf{x}^{(n)}$, $\mathbf{x}^{(n')}$ and the behaviour of the basis functions.

1.2 Non-parametric approach: Mean and covariance functions

In fact, we don't really need the basis functions and their parameters anymore. The influence of these appear only in the covariance matrix that describes the distribution of the targets, which is our key object. We can replace the parametric model altogether with a **covariance function** $C(\mathbf{x}, \mathbf{x}')$ which generates the elements of the covariance matrix

$$Q_{nn'} = C(\mathbf{x}^{(n)}, \mathbf{x}^{(n')}),$$

for any set of points $\{\mathbf{x}^{(n)}\}_{n=1}^N$.

Note, however, that \mathbf{Q} must be positive-definite. This constrains the set of valid covariance functions.

Once we have defined a covariance function, the covariance matrix for the target values will be given by

$$C_{nn'} = C(\mathbf{x}^{(n)}, \mathbf{x}^{(n')}) + \sigma_v^2 \delta_{nn'}.$$

A wide range of different covariance contributions can be **constructed**. These standard covariance functions are typically parametrized with hyperparameters $\boldsymbol{\alpha}$ so that

$$C_{nn'} = C(\mathbf{x}^{(n)}, \mathbf{x}^{(n')}, \boldsymbol{\alpha}) + \delta_{nn'} \Delta(\mathbf{x}^{(n)}; \boldsymbol{\alpha}),$$

where Δ is usually included as a flexible noise model.

Stationary kernels. The most common types of covariance functions are stationary, or translationally invariant, which implies that

$$C(\mathbf{x}^{(n)}, \mathbf{x}^{(n')}, \boldsymbol{\alpha}) = D(\mathbf{x} - \mathbf{x}'; \boldsymbol{\alpha}),$$

where the function D is often referred to as a *kernel*.

A very standard kernel is the RBF (also known as Exponentiated Quadratic or Gaussian kernel) which is differentiable infinitely many times (hence, very smooth),

$$C_{\text{RBF}}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\alpha}) = \alpha_0 + \alpha_1 \exp \left[-\frac{1}{2} \sum_{i=1}^I \frac{(x_i - x'_i)^2}{r_i^2} \right]$$

where I denotes the dimensionality of the input space. The hyperparameters are: $\boldsymbol{\alpha} = \{\alpha_0, \alpha_1, \vec{r}\}$. Sometimes, a single correlation length r is used.

2 GP models for regression

Let us return to the problem of predicting $t^{(N+1)}$ given \mathbf{t}_N . The independent variables \mathbf{X}_{N+1} are also given, but will be omitted from the conditional pdfs below.

The joint density is

$$p\left(t^{(N+1)}, \mathbf{t}_N\right) = p\left(t^{(N+1)}|\mathbf{t}_N\right) p\left(\mathbf{t}_N\right) \Rightarrow p\left(t^{(N+1)}|\mathbf{t}_N\right) = \frac{p\left(t^{(N+1)}, \mathbf{t}_N\right)}{p\left(\mathbf{t}_N\right)}.$$

Since both $p\left(t^{(N+1)}, \mathbf{t}_N\right)$ and $p\left(\mathbf{t}_N\right)$ are Gaussian distributions, then the conditional distribution, obtained by the ratio, must also be a Gaussian. Let us use the notation \mathbf{C}_{N+1} for the $(N+1) \times (N+1)$ covariance matrix for $\mathbf{t}_{N+1} = (\mathbf{t}_N, t^{(N+1)})$. This implies that

$$p\left(t^{(N+1)}|\mathbf{t}_N\right) \propto \exp\left[-\frac{1}{2}\left(\mathbf{t}_N, t^{(N+1)}\right) \mathbf{C}_{N+1}^{-1} \begin{pmatrix} \mathbf{t}_N \\ t^{(N+1)} \end{pmatrix}\right]$$

Summary

The prediction of the (Gaussian) pdf for $t^{(N+1)}$ requires an inversion of the covariance matrix \mathbf{C}_{N+1} .

Elegant approach using linear algebra tricks. Let us split the \mathbf{C}_{N+1} covariance matrix into four different blocks

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & \kappa \end{pmatrix},$$

where \mathbf{C}_N is the $N \times N$ covariance matrix (which depends on the positions \mathbf{X}_N), \mathbf{k} is an $N \times 1$ vector (that describes the covariance of \mathbf{X}_N with $\mathbf{x}^{(N+1)}$), while κ is the single diagonal element obtained from $\mathbf{x}^{(N+1)}$.

We can use the partitioned inverse equations (Barnett, 1979) to rewrite \mathbf{C}_{N+1}^{-1} in terms of \mathbf{C}_N^{-1} and \mathbf{C}_N as follows

$$\mathbf{C}_{N+1}^{-1} = \begin{pmatrix} \mathbf{M}_N & \mathbf{m} \\ \mathbf{m}^T & \mu \end{pmatrix},$$

where

$$\begin{aligned} \mu &= (\kappa - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k})^{-1} \\ \mathbf{m} &= -\mu \mathbf{C}_N^{-1} \mathbf{k} \\ \mathbf{M}_N &= \mathbf{C}_N^{-1} + \frac{1}{\mu} \mathbf{m} \mathbf{m}^T. \end{aligned}$$

Question

What are the dimensions of the different blocks? Check that the answer.

This implies that we can make a prediction for the Gaussian pdf of $t^{(N+1)}$ (meaning that we predict its value with an associated uncertainty) for an N^3 computational cost (the inversion of an $N \times N$ matrix).

Summary

The prediction for $t^{(N+1)}$ is a Gaussian

$$p\left(t^{(N+1)}|\mathbf{t}_N\right) = \frac{1}{Z} \exp \left[-\frac{\left(t^{(N+1)} - \hat{t}^{(N+1)}\right)^2}{2\sigma_{\hat{t}_{N+1}}^2} \right]$$

with

$$\begin{aligned} \text{mean : } \hat{t}^{(N+1)} &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N \\ \text{variance : } \sigma_{\hat{t}_{N+1}}^2 &= \kappa - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \end{aligned}$$

In fact, since the prediction only depends on the N available data we might as well predict several new target values at once. Consider $\mathbf{t}_M = \{t^{(N+i)}\}_{i=1}^M$ so that

$$\mathbf{C}_{N+M} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & \boldsymbol{\kappa} \end{pmatrix},$$

where \mathbf{k} is now an $N \times M$ matrix and $\boldsymbol{\kappa}$ an $M \times M$ matrix.

The prediction becomes a multivariate Gaussian

$$p(\mathbf{t}_{N+M}|\mathbf{t}_N) = \frac{1}{Z} \exp \left[-\frac{1}{2} (\mathbf{t}_M - \hat{\mathbf{t}}_M)^T \boldsymbol{\Sigma}_M^{-1} (\mathbf{t}_M - \hat{\mathbf{t}}_M) \right],$$

where the $M \times 1$ mean vector and $M \times M$ covariance matrix are

$$\begin{aligned} \hat{\mathbf{t}}_M &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N \\ \boldsymbol{\Sigma}_M &= \boldsymbol{\kappa} - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \end{aligned}$$

Optimizing the GP model hyperparameters. Predictions can be made once we have

1. Chosen an appropriate covariance function.
2. Determined the hyperparameters.
3. Evaluated the relevant blocks in the covariance function and inverted \mathbf{C}_N .

How do we determine the hyperparameters $\boldsymbol{\alpha}$? Well, recall that

$$p(\mathbf{t}_N) = \frac{1}{Z} \exp \left[-\frac{1}{2} \mathbf{t}_N^T \mathbf{C}_N^{-1} \mathbf{t}_N \right].$$

This pdf is basically a data likelihood.

- The frequentist approach would be to find the set of hyperparameters $\boldsymbol{\alpha}^*$ that maximizes the data likelihood, i.e. that minimizes $\boldsymbol{t}_N^T \boldsymbol{C}_N^{-1} \boldsymbol{t}_N$.
- A Bayesian approach would be to assign a prior to the hyperparameters and seek a posterior pdf $p(\boldsymbol{\alpha}|\boldsymbol{t}_N)$ instead.

The former approach is absolutely dominating the literature on GP regression. The covariance function hyperparameters are first optimized and then used for regression.