



DEPARTMENT OF MATHEMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Seminar: Gaussian Processes for Machine Learning

Gaussian Processes for Regression

Author: Davit Papikyan

Lecturer: Prof. Dr. Michael Marc Wolf



Contents

1	Introduction	1
2	Bayesian Linear Regression	1
2.1	Likelihood	2
2.2	Prior	2
2.3	Posterior	3
2.4	Posterior predictive distribution	3
2.5	Addressing nonlinear dependency	4
2.6	Kernel trick	5
3	Gaussian Process Regression	6
3.1	Gaussian Process	7
3.2	Posterior predictive distribution	7
3.3	Equivalence to Bayesian Linear Regression	9
3.4	Hyperparameter optimization	10
3.5	Pros and Cons	11
	Bibliography	12

1 Introduction

Supervised Learning is an approach to creating a model in Machine Learning that is trained on input data that has been labeled for particular output. The main classes of problems are regression and classification. Classification is related to data categorization in which the labels are discrete. In simpler words, classification is the grouping of related data points into classes. In case of regression, the labels are continuous. The aim of regression is to explain the relationships between a dependent variable and one or more independent variables.

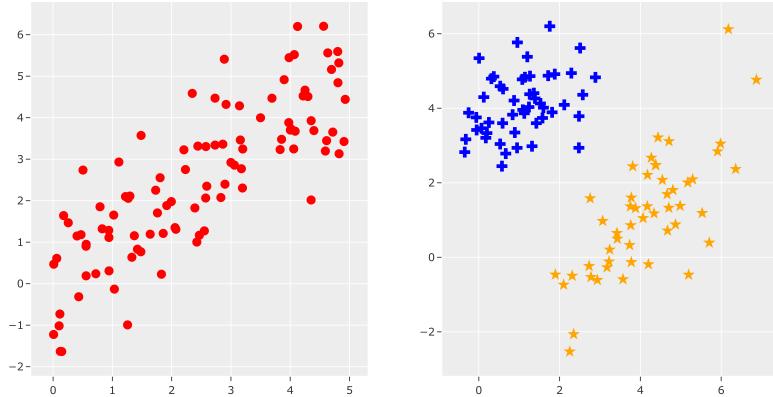


Figure 1: Regression (left) vs. classification (right) in terms of data given in Euclidean space

In this report we will analyze Gaussian Processes for solving the regression. Particularly, in section 2 we will analyze Bayesian interpretation of Linear Regression and will introduce the concept of kernel. Next, in section 3, we will introduce Gaussian Processes for regression and state how it is connected with Bayesian Linear Regression. In the end we will explore on a high level how hyperparameter tuning is done to achieve better results and will analyze the pros and cons of the model. [3]

2 Bayesian Linear Regression

Consider a training set with n observations, $\mathcal{D} = \{(x_i, y_i) \mid i = 1, 2, \dots, n\}$ where x is the input vector of dimension d and y is a scalar target value that we aim to learn

to predict. In a vectorized notation we have a data matrix (also called design matrix) $\mathbf{X} \in \mathbb{R}^{d \times n}$ which stores input data points as columns and a target vector $\mathbf{y} \in \mathbb{R}^n$. The goal of Bayesian Linear Regression is to model the relationship between \mathbf{X} and \mathbf{y} , that is, the conditional distribution of targets given inputs.

We assume that data is generated from the following process

$$y_i = f(\mathbf{x}_i) + \epsilon \text{ where } \epsilon \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2).$$

where \mathbf{w} is the weight vector we want to learn. Note that the variance of noise σ^2 is a hyperparameter (we will talk about it later). The second and main assumption is that f is linear, i.e. $f(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w}$. In order to make predictions for unseen data we need to derive the posterior predictive distribution over a single unseen data point for which we will need Likelihood, Prior and Posterior distributions.

2.1 Likelihood

Likelihood is the probability density of data given the parameters. Because of the independence assumption, the joint density can be factorized as given below

$$\begin{aligned} p(\mathbf{y} | \mathbf{X}, \mathbf{w}) &= \prod_{i=1}^n p(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2} \right] \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left[-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|^2 \right] = \mathcal{N}(\mathbf{X}^T \mathbf{w}, \sigma^2 \mathbf{I}). \end{aligned}$$

2.2 Prior

Prior is the probability density of parameters expressing our beliefs about them before observing the data. We will assume that parameters are driven from 0-centered Gaussian with some covariance matrix Σ_w , i.e.

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_w).$$

The choice of mean to be zero comes from the concept of Overfitting, a well-known problem in Machine Learning. To avoid that, one needs to lower the magnitude of weights of a model which is achieved by 0-centered assumption of prior.

2.3 Posterior

Posterior is the probability density of parameters given the data.

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \quad p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}) \times p(\mathbf{w})}{p(\mathbf{y} | \mathbf{X})}$$

The denominator of posterior, known as marginal likelihood, is independent of \mathbf{w} and plays a role of normalization constant. It is given by the following formula

$$p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}.$$

Since the density of a normally distributed random variable $Z \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ can up to a normalization constant be written as $\exp[-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})]$, then

$$\begin{aligned} p(\mathbf{w} | \mathbf{X}, \mathbf{y}) &\propto \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}^T \mathbf{w})^T (\mathbf{y} - \mathbf{X}^T \mathbf{w}) \right] \exp \left[\frac{1}{2} \mathbf{w}^T \boldsymbol{\Sigma}_w \mathbf{w} \right] \\ &\propto \exp \left[-\frac{1}{2} (\mathbf{w} - \tilde{\mathbf{w}})^T \underbrace{\left(\frac{1}{\sigma^2} \mathbf{X} \mathbf{X}^T + \boldsymbol{\Sigma}_w^{-1} \right)}_A (\mathbf{w} - \tilde{\mathbf{w}}) \right], \end{aligned}$$

where $\tilde{\mathbf{w}} = \frac{1}{\sigma^2} \mathbf{X} \mathbf{y}$ and $A = \frac{1}{\sigma^2} \mathbf{X} \mathbf{X}^T + \boldsymbol{\Sigma}_w^{-1}$. Therefore

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\tilde{\mathbf{w}}, A^{-1}).$$

2.4 Posterior predictive distribution

Now focus on the main problem – predicting a target y_* for an unseen data point \mathbf{x}_* . For so we need to compute the posterior predictive distribution $p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y})$:

$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_*, \mathbf{w} | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ &= \int p(y_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ &= \mathcal{N}\left(\frac{1}{\sigma^2} \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y}, \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{x}_*\right). \end{aligned}$$

Using the mean of posterior predictive as a predicted target we can also estimate the *uncertainty* using the covariance.

2.5 Addressing nonlinear dependency

In case when the dependency between y and x is not linear, we can transform data points to some other space where the dependency between the response variable and new features will be linear. Let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^k$ with $\Phi := \Phi(X) \in \mathbb{R}^{k \times n}$ be the aggregation of columns $\phi(x)$ for all samples in training set. Our model will now become $f(x) = \phi(x)^T w$ where $w \in \mathbb{R}^k$. The posterior predictive is given below

$$p(y_* | x_*, X, y) = \mathcal{N}\left(\frac{1}{\sigma^2} \phi(x_*)^T A^{-1} \Phi y, \phi(x_*)^T A^{-1} \phi(x_*)\right)$$

$$\text{where } A = \frac{1}{\sigma^2} \Phi \Phi^T + \Sigma_w^{-1}.$$

Note that we need to compute the inverse of a $k \times k$ matrix which is not tractable for large values of k . To address this problem, we will try to transform both mean and variance of the posterior predictive distribution to escape from computation of A^{-1} . Denote $K := \Phi^T \Sigma_w \Phi$ and $\phi_* := \phi(x_*)$. Let's start with rewriting the mean.

$$A \Sigma_w \Phi = \frac{1}{\sigma^2} \Phi \Phi^T \Sigma_w \Phi + \Phi = \frac{1}{\sigma^2} \Phi (\Phi^T \Sigma_w \Phi + \sigma^2 I) = \frac{1}{\sigma^2} \Phi (K + \sigma^2 I).$$

Multiplying both sides by A^{-1} from left and $(K + \sigma^2 I)^{-1}$ from right we will get

$$\Sigma_w \Phi (K + \sigma^2 I)^{-1} = \frac{1}{\sigma^2} A^{-1} \Phi,$$

by using which we can rewrite the mean of posterior predictive distribution as

$$\frac{1}{\sigma^2} \phi_*^T A^{-1} \Phi y = \phi_*^T \Sigma_w \Phi (K + \sigma^2 I)^{-1} y.$$

Next we rewrite the variance using the matrix inversion lemma

$$(B + UCV)^{-1} = B^{-1} - B^{-1}U(C^{-1} + VB^{-1}U)^{-1}VB^{-1}.$$

By denoting $B := \Sigma_w^{-1}$, $U := \Phi$, $C := \frac{1}{\sigma^2} I$, $V := \Phi^T$ we will get

$$A^{-1} = \left(\frac{1}{\sigma^2} \Phi \Phi^T + \Sigma_w^{-1} \right)^{-1} = \Sigma_w - \Sigma_w \Phi (\sigma^2 I + \underbrace{\Phi^T \Sigma_w \Phi}_K)^{-1} \Phi^T \Sigma_w.$$

After plugging the above equation for A^{-1} into the expression of variance we get

$$\boldsymbol{\phi}_*^T A^{-1} \boldsymbol{\phi}_* = \boldsymbol{\phi}_*^T \boldsymbol{\Sigma}_w \boldsymbol{\phi}_*^T - \boldsymbol{\phi}_*^T \boldsymbol{\Sigma}_w \boldsymbol{\Phi} (\boldsymbol{K} + \sigma^2 \boldsymbol{I})^{-1} \boldsymbol{\Phi}^T \boldsymbol{\Sigma}_w \boldsymbol{\phi}_*.$$

Thus, the posterior predictive distribution can be rewritten as

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}\left(\boldsymbol{\phi}_*^T \boldsymbol{\Sigma}_w \boldsymbol{\Phi} (\boldsymbol{K} + \sigma^2 \boldsymbol{I})^{-1} \mathbf{y}, \right. \\ \left. \boldsymbol{\phi}_*^T \boldsymbol{\Sigma}_w \boldsymbol{\phi}_*^T - \boldsymbol{\phi}_*^T \boldsymbol{\Sigma}_w \boldsymbol{\Phi} (\boldsymbol{K} + \sigma^2 \boldsymbol{I})^{-1} \boldsymbol{\Phi}^T \boldsymbol{\Sigma}_w \boldsymbol{\phi}_*\right)$$

where we need to invert matrix of size $n \times n$ which is more convenient when $n < k$.

In the following examples we have generated synthetic data points (in red) from $y = \sin(x)$ true signal by perturbing function values with Gaussian noise. In the first example we don't transform data, i.e. we use identity basis function $\phi(x) = x$.

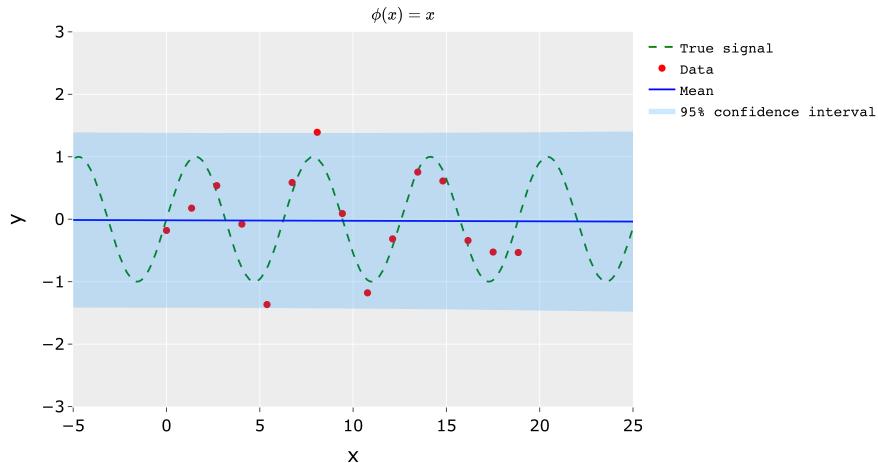


Figure 2: Bayesian Linear Regression applied on data generated from $\sin(x)$ function perturbed by a negligible Gaussian noise

In contrast to the first example, below we have applied $\sin(x)$ transformation onto original data space. One can easily see that transformed model fits the data better and approximates true signal with high accuracy (see Figure 3).

2.6 Kernel trick

From the expressions of mean and variance obtained for posterior predictive distribution, one can notice that feature vectors $\boldsymbol{\phi}(x)$ appear always with covariance matrix

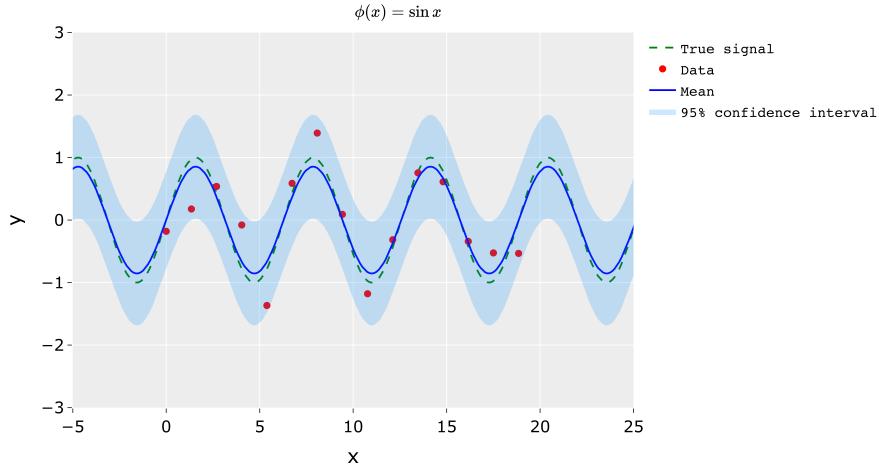


Figure 3: Bayesian Linear Regression applied on sine transformed data generated from $\sin(x)$ function perturbed by a negligible standard Gaussian noise

of prior, i.e. in the forms of $\Phi^T \Sigma_w \Phi$, $\phi_*^T \Sigma_w \Phi$ or $\phi_*^T \Sigma_w \phi_*$. An entry of each of these matrices has the form of $k(x, x') := \phi(x)^T \Sigma_w \phi(x')$ which is known as *kernel function*. If we define $\psi(x) := \Sigma_w^{1/2} \phi(x)$ we can rewrite kernel as a simple dot product

$$k(x, x') = \psi(x)^T \psi(x'),$$

where $\Sigma_w^{1/2} = U D^{1/2} U^T$ for $U D U^T$ being the SVD of Σ_w . That is, we can think of a kernel value to be a measure of interaction (similarity) of data points. Kernel function provides us with an efficient (computationally cheaper) way to transform data without having to compute the coordinates in a new space, i.e. feature vectors $\phi(x)$. This approach is known in the literature as **kernel trick**.

Being equipped with the preliminary knowledge we are now ready to discuss Gaussian Process Regression.

3 Gaussian Process Regression

Besides assuming that data is normally distributed with some 0-centered noise, we additionally assume that function f is a random variable drawn from a Gaussian Process conditioned on data.

3.1 Gaussian Process

Definition: Gaussian Process

Let $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be a positive definite kernel and $m : \mathbb{R}^d \rightarrow \mathbb{R}$ be a finite real-valued function. Then a random function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be a Gaussian Process (GP) with mean m and covariance kernel k , denoted by $GP(m, k)$, if the following holds: for any finite set $X = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \subset \mathbb{R}^d$ of any size $n \in \mathbb{N}$, the random vector

$$f_X = (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n))^T \in \mathbb{R}^n$$

follows a multivariate normal distribution $\mathcal{N}(m_X, k_{XX})$ with covariance $k_{XX} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^n$ and mean vector $m_X = [m(\mathbf{x}_1), m(\mathbf{x}_2), \dots, m(\mathbf{x}_n)]^T$. It is written as

$$f \sim GP(m, k).$$

The key point of the definition is that the function values $f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)$ are jointly Gaussian for any finite n .

In order to proceed with regression, as stated in the definition of Gaussian Process, we need to define mean and kernel (covariance) functions

$$\begin{aligned} m(\mathbf{x}) &= \mathbb{E}_{f \sim GP(m,k)}[f(\mathbf{x})], \\ k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{f \sim GP(m,k)}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]. \end{aligned}$$

Let's take the mean to be a constant 0 function and define the kernel to be RBF (Radial Bases Function) with length scale parameter of 1

$$k(\mathbf{x}, \mathbf{x}') = \text{cov}(f(\mathbf{x}), f(\mathbf{x}')) = \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|_2^2\right).$$

Assuming to have a finite point set X_* we can compute the corresponding covariance matrix $K(X_*, X_*)$ by applying the kernel function to each pair of points. Next, let's draw some samples from the prior distribution $f \sim \mathcal{N}(\mathbf{0}, K(X_*, X_*))$ (see Figure 4).

3.2 Posterior predictive distribution

Recall the standard assumptions $y = f(\mathbf{x}) + \epsilon$ where noise is i.i.d. $\mathcal{N}(0, \sigma^2)$. This implies that

$$\begin{aligned} \text{cov}(y_i, y_j) &= \text{cov}(f(\mathbf{x}_i) + \epsilon_i, f(\mathbf{x}_j) + \epsilon_j) = \\ &\quad \text{cov}(f(\mathbf{x}_i), f(\mathbf{x}_j)) + \text{cov}(f(\mathbf{x}_i), \epsilon_j) + \\ &\quad \text{cov}(\epsilon_i, f(\mathbf{x}_j)) + \text{cov}(\epsilon_i, \epsilon_j) = \\ &\quad k(\mathbf{x}_i, \mathbf{x}_j) + \sigma^2 \delta_{ij} \end{aligned}$$

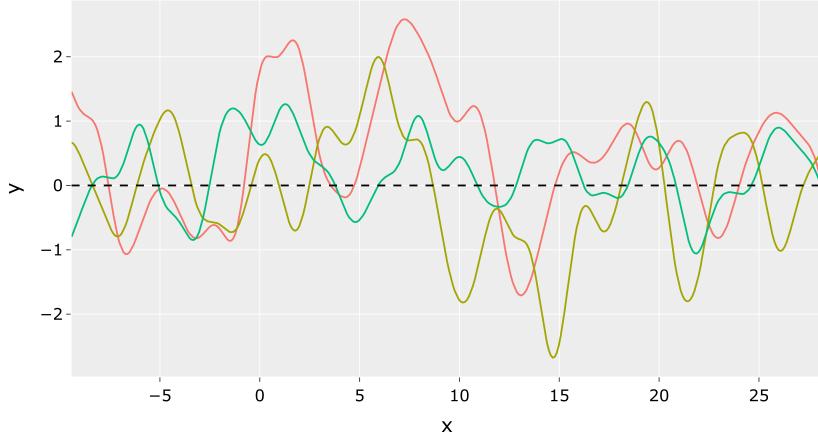


Figure 4: "Samples" drawn from prior distribution

or in a vectorized notation

$$\text{cov}(\mathbf{y}) = k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}.$$

According to the definition of Gaussian processes, the joint distribution of observed and test (unseen) values is again Gaussian

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & K(\mathbf{X}, \mathbf{X}_*) \\ K(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix}\right).$$

Using the property that Gaussian distribution is closed under conditioning, one can write that

$$\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*))$$

where

$$\begin{aligned} \bar{\mathbf{f}}_* &:= \mathbb{E}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] = K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \\ \text{cov}(\mathbf{f}_*) &= K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X}_*)(K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} K(\mathbf{X}, \mathbf{X}_*). \end{aligned}$$

Let's rewrite the result that we got for a single test point in order to compare it with Bayesian Linear Regression case. Denote $\mathbf{K} := K(\mathbf{X}, \mathbf{X})$, $\mathbf{K}_* := K(\mathbf{X}, \mathbf{X}_*)$, $\mathbf{k}_* := k(\mathbf{x}_*) = [k(\mathbf{x}_i, \mathbf{x}_*)]_{i=1}^n$ the covariance vector of a test point \mathbf{x}_* and n training points. Then,

$$\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \mathbb{V}[\mathbf{f}_*])$$

where

$$\bar{f}_* = \mathbf{k}_*^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\mathbb{V}[f_*] = \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_*.$$

The prediction of a Gaussian Process Regression is the value computed by the mean function of its posterior predictive distribution, i.e. by \bar{f}_* . $\mathbb{V}[f_*]$, as in the case of Bayesian Linear Regression, provides us with an estimate of uncertainty of predictions.

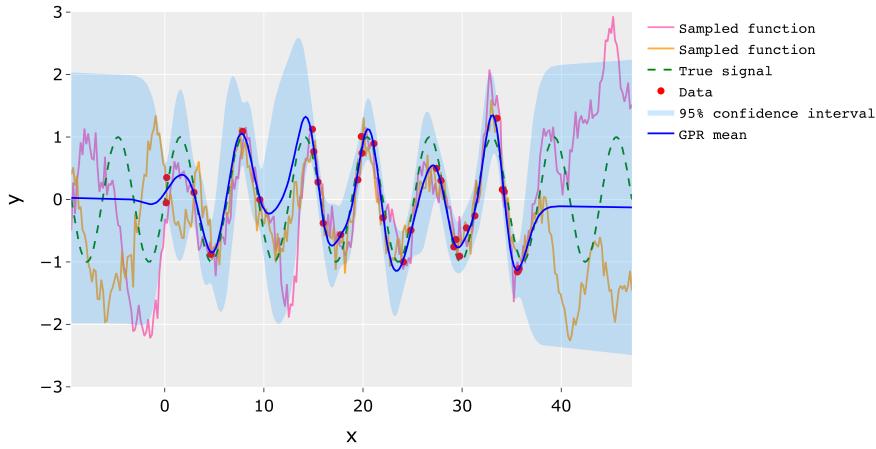


Figure 5: Gaussian Process Regression example

Interestingly, the model can also be used for solving a well known problem in numerical mathematics - interpolation. To do so, we only need to assume that the generated data is noise-free, i.e. $\epsilon = 0$ and do the corresponding changes in posterior predictive distribution (see Figure 6).

3.3 Equivalence to Bayesian Linear Regression

Recall that the prediction function of Bayesian Linear Regression is provided by its mean which has the following form

$$\boldsymbol{\phi}_*^T \boldsymbol{\Sigma}_w \boldsymbol{\Phi} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}.$$

If we define the kernel to be $k(\mathbf{x}_i, \mathbf{x}_j) := \boldsymbol{\phi}(\mathbf{x}_i)^T \boldsymbol{\Sigma}_w \boldsymbol{\phi}(\mathbf{x}_j)$, then the prediction function of Gaussian Process regression $\boldsymbol{\phi}(\mathbf{x}_*)^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$ will become equivalent to that of Bayesian Linear Regression.

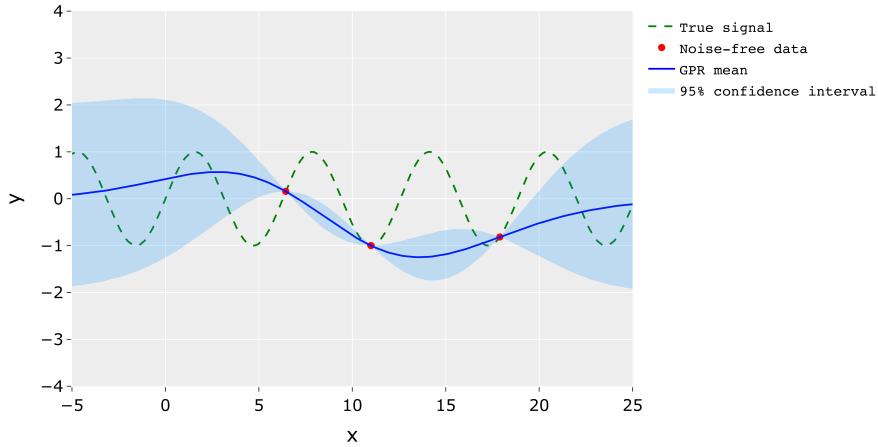


Figure 6: Gaussian Process Interpolation example

3.4 Hyperparameter optimization

In general, the performance of Gaussian Process Regression is highly dependent on hyperparameters - parameters that affect the model's behavior but are not learnt during the training process. An example of a hyperparameter is length scale l in RBF kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2l^2}\right).$$

Below we can see the effect of different values of length scale on prior functions sampled with constant 0 mean

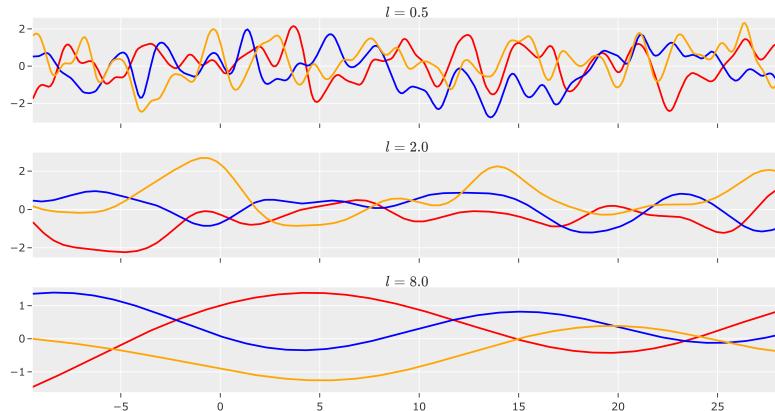


Figure 7: Fluctuations of prior functions increase with the increase of length scale

Hyperparameters can be optimized by maximizing so-called marginal likelihood (equivalently marginal log-likelihood)

$$p(\mathbf{y} | \mathbf{X}) = \int \overbrace{p(\mathbf{y} | f, \mathbf{X})}^{\text{likelihood}} \overbrace{p(f | \mathbf{X})}^{\text{prior}} df.$$

Because prior and likelihood above are Gaussian, $f | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$ and $\mathbf{y} | f \sim \mathcal{N}(f, \sigma^2 \mathbf{I})$, then the final marginal log-likelihood can be given by

$$\log p(\mathbf{y} | \mathbf{X}) = -\frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K} + \sigma^2 \mathbf{I}| - \frac{n}{2} \log 2\pi.$$

Denoting by Θ the set of all hyperparameters, then finding the optimal hyperparameters is equivalent to the following optimization problem

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \log p(\mathbf{y} | \mathbf{X}, \Theta),$$

which can be solved with any gradient-based optimization algorithm.

3.5 Pros and Cons

Gaussian Process Regression is flexible model since it can take any valid kernel function to define a prior distribution, i.e. the prediction function's shape. It is interpretable mostly because it allows to estimate model's uncertainty about any prediction (or region of domain). Another advantage is that these models can directly be optimized given the hyperparameters, as the computation of posterior predictive is pure Linear Algebra stuff.

Because the mean of the posterior predictive can be represented as a linear combination of targets

$$\bar{f}(\mathbf{x}_*) = \underbrace{\mathbf{k}(\mathbf{x}_*)^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1}}_{\alpha} \mathbf{y} = \sum_{i=1}^n \alpha_i \mathbf{y}_i,$$

then having outliers in the data can drive the mean function (and hence the predictions). Based on this we can state that Gaussian Process Regression is prone to outliers unless the hyperparameters and kernel are cleverly chosen. Moreover, since the mean and covariance of posterior predictive involve matrix/vector products and inversions, computing them when having high-dimensional data can be computationally expensive.

Bibliography

- [1] M. Kanagawa, P. Hennig, D. Sejdinovic, and B. K. Sriperumbudur. *Gaussian Processes and Kernel Methods: A Review on Connections and Equivalences*. 2018. arXiv: 1807.02582 [stat.ML].
- [2] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. “Scikit-learn: Machine Learning in Python.” In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [3] C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. Adaptive computation and machine learning. MIT Press, 2006, pp. I–XVIII, 1–248. ISBN: 026218253X.