

Lecture 22: Gaussian Processes

The notes below are based on references [?, Sec. 6.4] and [?, Ch. 2,4,5] and [?, Ch. 2].

Gaussian Process: Function-Space View

Gaussian Process Definition. A GP $f(x)$ is a *random function* of an argument $x \in \mathcal{X}$ where \mathcal{X} is an input domain. For example, if the GP is a temporal process then $\mathcal{X} = \mathbb{R}$ is the time domain, or if the GP is a spatial process then $\mathcal{X} = \mathbb{R}^2$ or \mathbb{R}^3 . If the GP models a dynamical system then the input may be the system state $\mathcal{X} = \mathbb{R}^n$. Often the output $f(x)$ is a scalar but this is not strictly required. A GP is completely specified by its mean function and covariance function. Let $m(x)$ denote the mean over the input space and let $k(x, x')$ denote the covariance function (also called the *kernel*). Formally,

$$m(x) = \mathbb{E}[f(x)] \quad (1)$$

$$k(x, x') = \mathbb{E}[\{f(x) - m(x)\}\{f(x') - m(x')\}] = \text{Cov}(x, x') . \quad (2)$$

From the above definition it follows that:

1. The kernel is positive since

$$k(x, x) = \mathbb{E}[\{f(x) - m(x)\}^2] = \text{Var}(x) \geq 0 \quad (3)$$

2. The kernel is symmetric since

$$k(x, x') = \mathbb{E}[\{f(x) - m(x)\}\{f(x') - m(x')\}] \quad (4)$$

$$= \mathbb{E}[\{f(x') - m(x')\}\{f(x) - m(x)\}] \quad (5)$$

$$= k(x', x) \quad (6)$$

These resemble our familiar definitions of mean and covariance, however they pertain to random *functions* rather than random variables or random vectors. The GP definition above is compactly summarized with the notation:

$$f(x) \sim \mathcal{GP}(m(x), k(x, x')) . \quad (7)$$

Covariance Functions (Kernels). At the heart of GPs is the covariance function or kernel $k(x, x')$ that encodes a sense of “smoothness” of the GP over the input space. If the two nearby points, x and x' , are highly correlated, then $k(x, x')$ will be large. Conversely, if the two points are far away then $k(x, x')$ will be small. The sense of “close” or “far” is encoded by the shape of the covariance function $k(x, x')$, which is often parameterized by a time or length scale and other constants called *hyperparameters*. The distance between two input points is sometimes called the *lag* $\mathbf{h} = x - x'$ (a term common in the geostatistics). When the covariance function depends only on the lag (i.e., $k(x, x') = k(\mathbf{h})$) it is called *stationary*. Moreover, if the covariance function depends only on the magnitude of the lag vector $k(x, x') = k(\|\mathbf{h}\|)$ it is called *isotropic* (i.e., direction-independent). Stationary, isotropic covariance functions are perhaps the most common and include:

- Squared exponential:

$$k_{SE}(\|\mathbf{h}\|) = \sigma^2 \exp\left(-\frac{\|\mathbf{h}\|^2}{2L}\right) \quad (8)$$

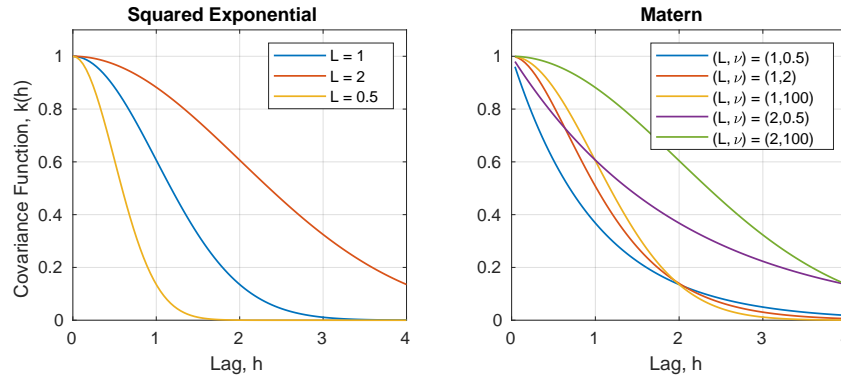
where the hyperparameters are $\boldsymbol{\theta} = [L, \sigma]^T$ and L is a length-scale parameter, and σ^2 is a variance parameter.

- Matérn class

$$k_{\text{Matern}}(\|\mathbf{h}\|) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu} \left(\frac{\|\mathbf{h}\|}{L} \right) \right)^2 K_\nu \left(\sqrt{2\nu} \left(\frac{\|\mathbf{h}\|}{L} \right) \right) \quad (9)$$

where the hyperparameters are $\boldsymbol{\theta} = [L, \sigma, \nu]^T$ and, as before, L is a length-scale parameter, σ^2 is a variance parameter, and ν is a positive constant shape parameter. In this expression $K_\nu(\cdot)$ is the modified Bessel function of the second kind and $\Gamma(\cdot)$ is the gamma function.

These covariance functions are plotted below for different values of the hyperparameters (with $\sigma = 1$ in both cases).



The main requirement for the covariance function is for it to be positive definite. Olea [?, Defn. 2.3] defines a positive definite function as follows. Let m be a positive integer, let $\{\alpha_1, \alpha_2, \dots, \alpha_m\}$ be a set of real or complex numbers, and let $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ be a set of points in an n -dimensional Euclidean space. The function $\phi(\mathbf{x}_i, \mathbf{x}_j)$ is a positive definite function if

$$\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \phi(\mathbf{x}_i, \mathbf{x}_j) \geq 0. \quad (10)$$

Since the product of two positive definite functions is also positive definite the kernels can be multiplied to form new (valid) kernels¹. This is often used when modeling spatiotemporal processes. For example, let the input space consist of vectors $\mathbf{x} = [s, t]^T$ where $s \in D \subset \mathbb{R}^2$ is a spatial coordinate and $t \in \mathbb{R}^2$ is the time. Spatial correlations can be modeled by a kernel of the form $k_s(s, s')$ whereas temporal correlations can be modeled by a kernel of the form $k_t(t, t')$. Then a spatiotemporal kernel can be formed as

$$k(\mathbf{x}, \mathbf{x}') = k_s(s, s') k_t(t, t'). \quad (11)$$

¹For an extensive overview see the Kernel Cookbook [Link]

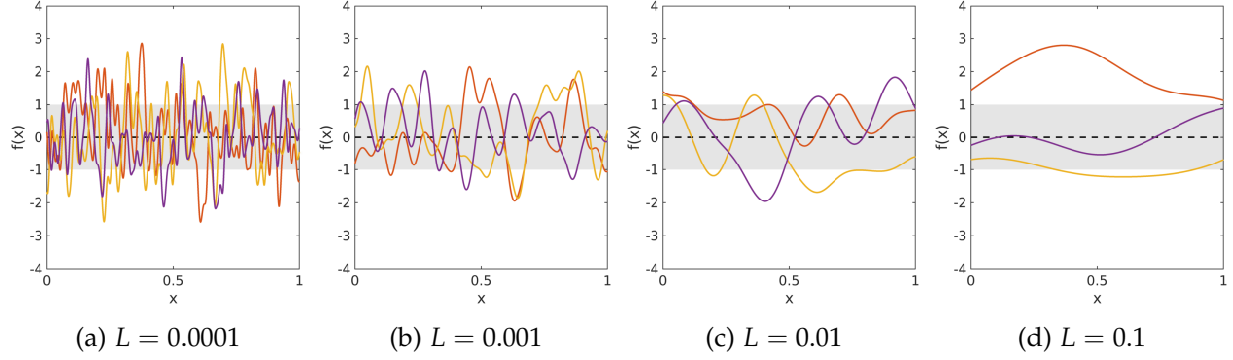


Figure 1: Gaussian process realizations on the interval of $x \in [0, 1]$ using a squared exponential kernel with $\sigma^2 = 1$ with varying lengthscales L .

Realizations. Since the GP is a random function there are infinitely many *realizations* (i.e., example functions) that can be sampled from it. On average, the ensemble of these realizations satisfies the mean $m(x)$ and covariance $k(x, x')$ of the GP. We can numerically simulate a GP over this input space as follows. First choose a set of Q points over which to define the realization $\bar{\mathbf{X}} = \{\bar{x}_1, \dots, \bar{x}_Q\}$ and evaluate the mean at each point:

$$\boldsymbol{\mu}(\bar{\mathbf{X}}) = \begin{bmatrix} m(\bar{x}_1) \\ \vdots \\ m(\bar{x}_Q) \end{bmatrix}. \quad (12)$$

Then, compute the covariance matrix relating all of the grid points to each other as

$$\mathbf{K}(\bar{\mathbf{X}}, \bar{\mathbf{X}}) = \begin{bmatrix} k(\bar{x}_1, \bar{x}_1) & k(\bar{x}_1, \bar{x}_2) & \cdots & k(\bar{x}_1, \bar{x}_Q) \\ k(\bar{x}_2, \bar{x}_1) & \vdots & \ddots & k(\bar{x}_2, \bar{x}_Q) \\ \vdots & \vdots & \ddots & \vdots \\ k(\bar{x}_Q, \bar{x}_1) & k(\bar{x}_Q, \bar{x}_2) & \cdots & k(\bar{x}_Q, \bar{x}_Q) \end{bmatrix} \quad (13)$$

A GP realization is then found by sampling from the distribution

$$\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu}(\bar{\mathbf{X}}), \mathbf{K}(\bar{\mathbf{X}}, \bar{\mathbf{X}})) \quad (14)$$

For example, in MATLAB the `mvrnd` function can be used.

GP Regression. Rather than simulation GP realizations, a more common usage of GPs is to perform regression—given some set of partial observations \mathbf{f} at locations \mathbf{X} we wish to use the assumed mean and covariance to predict the values of the process \mathbf{f}_* at the grid points \mathbf{X}_* (or a single point, if desired). Suppose that the number of observations is $|\mathbf{X}| = N$ and the number of grid points is $|\mathbf{X}_*| = M$. To perform the regression we begin by forming a new stacked vector of the observations $[\mathbf{f}^T, \mathbf{f}_*^T]^T$. For simplicity, assume the mean of the GP is zero. Then this new random vector is distributed according to the joint distribution:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix}\right) \quad (15)$$

where

- $K(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N}$ is the covariance matrix relating observation points to each other as in (??).
- $K(\mathbf{X}_*, \mathbf{X}) \in \mathbb{R}^{M \times N}$ is the covariance matrix relating the test/prediction grid points to the observation points, similar to (??). Note that $K(\mathbf{X}_*, \mathbf{X}) = K(\mathbf{X}, \mathbf{X}_*)^T$.
- $K(\mathbf{X}_*, \mathbf{X}_*) \in \mathbb{R}^{M \times M}$ is the covariance matrix relating the test/prediction grid points to each other.

Equation (??) represents the joint distribution between the random vectors \mathbf{f}_* (the prediction/grid values) and \mathbf{f} (the observed values). To perform the regression we first review the following fact.

Partitioned Gaussian Vectors. Consider a Gaussian random vector partitioned as $\mathbf{z} = [\mathbf{x}^T, \mathbf{y}^T]^T$ with mean

$$E[\mathbf{z}] = \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \quad (16)$$

and partitioned covariance matrix

$$\mathbf{P}_f = \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_{yy} \end{bmatrix} \quad (17)$$

The symmetry of \mathbf{P}_z implies that $\mathbf{P}_{xy} = \mathbf{P}_{yx}^T$. It can be shown that the conditional probability of \mathbf{y} given \mathbf{x} is also a Gaussian vector with mean and covariance:

$$\boldsymbol{\mu}_{y|x} = \boldsymbol{\mu}_y + \mathbf{P}_{yx} \mathbf{P}_{xx}^{-1} (\mathbf{x} - \boldsymbol{\mu}_x) \quad (18)$$

$$\mathbf{P}_{y|x} = \mathbf{P}_{yy} - \mathbf{P}_{yx} \mathbf{P}_{xx}^{-1} \mathbf{P}_{xy} . \quad (19)$$

For a derivation of this fact refer to the Appendix notes on Gaussian random vectors.

Applying the above conditioning principle to (??) we find that the mean value $\mathbf{f}_* \in \mathbb{R}^{M \times 1}$ at the grid points \mathbf{X}_* is conditioned on the observed values $\mathbf{f} \in \mathbb{R}^{N \times 1}$ is a Gaussian random vector

$$\mathbf{f}_* \sim \mathcal{N}(\boldsymbol{\mu}_{f_*|f}, \mathbf{P}_{f_*|f}) \quad (20)$$

where

$$\boldsymbol{\mu}_{f_*|f} = K(\mathbf{X}_*, \mathbf{X}) K(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f} \quad (21)$$

$$\mathbf{P}_{f_*|f} = K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X}) K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{X}_*) \quad (22)$$

Notice that the predicted values $\boldsymbol{\mu}_{f_*|f}$ in (??) are a linear combination of the existing data \mathbf{f} and the matrix $K(\mathbf{X}_*, \mathbf{X}) K(\mathbf{X}, \mathbf{X})^{-1}$ can be viewed as a matrix of weights that determine the contribution of each data point to each prediction. A useful feature of GP regression is that it is accompanied immediately by a covariance matrix $\mathbf{P}_{f_*|f}$ given by (??). The diagonal elements of this matrix describe the uncertainty at each estimation point. As an example in Fig. ??, the resulting mean and associated uncertainty are plotted with increasing number of measurements. Both (??)–(??) involve the term $K(\mathbf{X}, \mathbf{X})^{-1}$ which requires inverting a $N \times N$ matrix where N is the number of data points. For large data sets this can be prohibitive and numerous techniques have been proposed to approximate the regression.

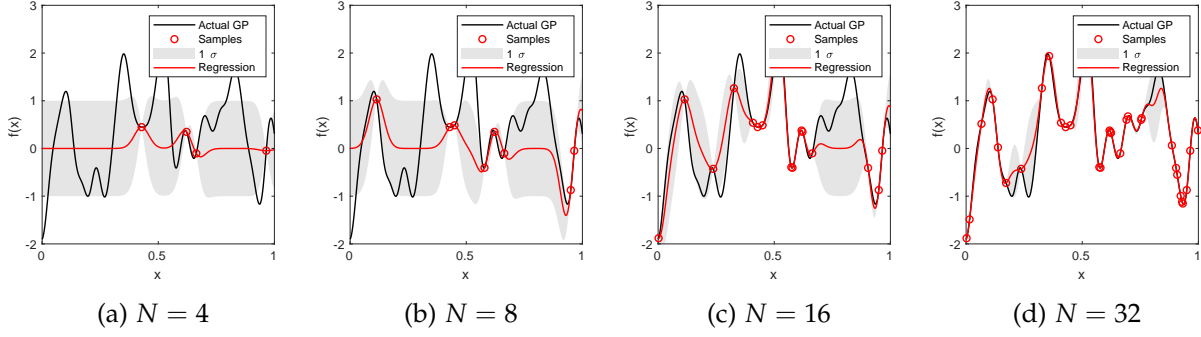


Figure 2: Gaussian process regression with an increasing number of measurements.

GP Regression with Noisy Measurements. Suppose now that the measurements of the GP are corrupted by Gaussian noise

$$y = f(x) + \epsilon \quad (23)$$

where $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$. In this setting, \mathbf{f} is considered a “latent” or “hidden” variable. The vector of noisy observations \mathbf{y} is dependent on the latent variable \mathbf{f} by

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_n^2 \mathbf{1}) \quad (24)$$

We assumed that \mathbf{f} is zero-mean and has covariance given by $\mathbf{K}(\mathbf{X}, \mathbf{X})$, thus

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X})) \quad (25)$$

The marginal distribution conditioned on the input values is

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} \quad (26)$$

where the terms $p(\mathbf{y}|\mathbf{f})$ and $p(\mathbf{f})$ are given above. Using properties of Gaussian random vectors it follows that the joint distribution is

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{1} & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix}\right) \quad (27)$$

Again, applying the conditioning principle to (27) we find that the mean value \mathbf{f} at the grid points \mathbf{X} is conditioned on the observed values \mathbf{f} is a Gaussian random vector

$$\mathbf{f}_* \sim \mathcal{N}(\mu_{\mathbf{f}_*|\mathbf{f}}, \mathbf{P}_{\mathbf{f}_*|\mathbf{f}}) \quad (28)$$

where

$$\mu_{\mathbf{f}_*|\mathbf{f}} = \mathbf{K}(\mathbf{X}_*, \mathbf{X})(\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{1})^{-1} \mathbf{y} \quad (29)$$

$$\mathbf{P}_{\mathbf{f}_*|\mathbf{f}} = \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) - \mathbf{K}(\mathbf{X}_*, \mathbf{X})(\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{1})^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \quad (30)$$

The figures below illustrate the GP regression with noisy data.

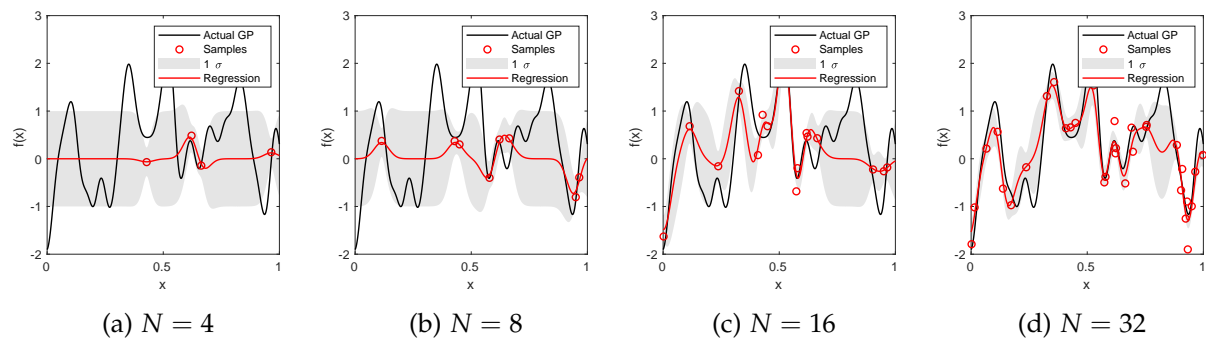


Figure 3: Gaussian process regression with an increasing number of measurements corrupted by additive Gaussian noise.

GP Realizations (Conditioned on Data). As before, we can generate sample realizations of the GP that are now conditioned on the data. This amounts to generate a random vector from the normal distribution given by the posterior (??) or (??) for the noise-free or noisy cases, respectively. Example GP realizations conditioned on data are shown below.

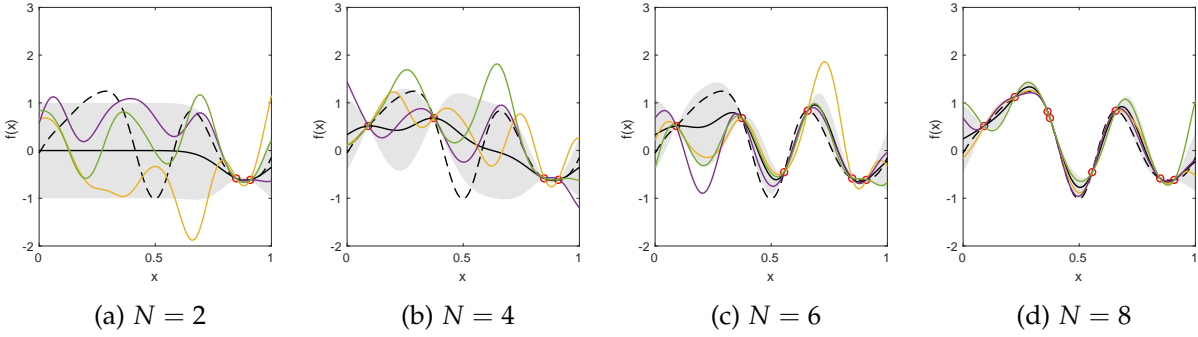


Figure 4: Gaussian process realizations conditioned on data for an increasing number of measurements. The solid black line is the regression, the dashed black line is the actual GP, and the colored lines are various GP realizations

Hyperparameter selection. Thusfar we have not commented on how the hyperparameters of the covariance function are chosen. However, as the figures below illustrate, the choice of hyperparameters can greatly affect both the mean and variance of the posterior GP regression. In the following we will overview an approach for optimizing the hyperparameters.

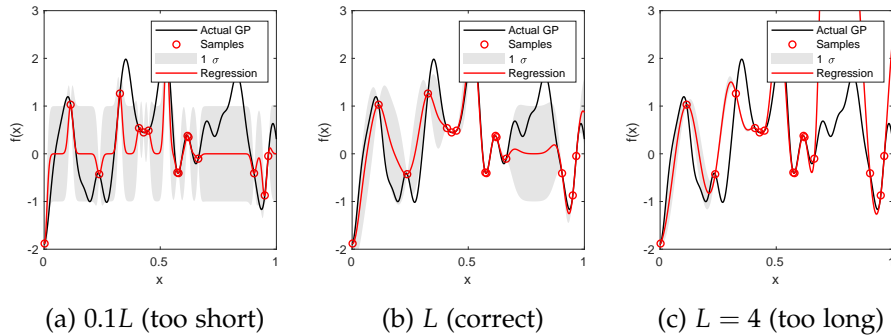


Figure 5: Hyperparameter selection

Hyperparameter Optimization. If the hyperparameters are unknown an optimization problem may be formulated to determine them, often with the aim of minimizing the the marginal likelihood $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$, which is the probability of the observed (noisy) data \mathbf{y} given the model. If the covariance function is fixed, then the statistical model relating the observations is only a function of the observation locations \mathbf{X} and the model hyperparameters $\boldsymbol{\theta}$. The marginal likelihood can be computed as follows

$$p(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X}; \boldsymbol{\theta})d\mathbf{f} \quad (31)$$

which is normally distributed with covariance

$$\mathbf{K}_y = \sigma_n^2 \mathbf{1} + \mathbf{K} \quad (32)$$

(see [?, Sec. 6.4.2-6.4.3] and [?, Sec. 5.4.1]). It is often more useful to work with the log likelihood

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{n}{2} \log 2\pi \quad (33)$$

For optimization, we need the partial derivatives (with respect to θ_j for all $j = 1, \dots, |\boldsymbol{\theta}|$) :

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{\partial}{\partial \theta_j} \left(\frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} \right) - \frac{\partial}{\partial \theta_j} \left(\frac{1}{2} \log |\mathbf{K}_y| \right) - \underbrace{\frac{\partial}{\partial \theta_j} \left(\frac{n}{2} \log 2\pi \right)}_{=0} \quad (34)$$

$$= -\frac{1}{2} \mathbf{y}^T \left(\frac{\partial}{\partial \theta_j} \mathbf{K}_y^{-1} \right) \mathbf{y} - \frac{1}{2} \frac{\partial}{\partial \theta_j} \log |\mathbf{K}_y| \quad (35)$$

$$= -\frac{1}{2} \mathbf{y}^T \left(-\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_j} \mathbf{K}_y^{-1} \right) \mathbf{y} - \frac{1}{2} \text{tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_j} \right) \quad (36)$$

where the following identities described below were used.

Aside: The follow identity holds when differentiating a matrix $\mathbf{U}(x)$ with respect to a scalar x :

$$\frac{\partial \mathbf{U}^{-1}}{\partial x} = -\mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \mathbf{U}^{-1} \quad (37)$$

or a scalar-by-scalar with matrices involved

$$\frac{\partial \log |\mathbf{U}|}{\partial x} = \text{tr} \left(\mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \right) \quad (38)$$

Link: For a more extensive list of identities see [here](#).

The choice of kernel will determine the partial derivative $\frac{\partial \mathbf{K}_y}{\partial \theta_j}$. We can now formulate an optimization problem to find the optimal hyperparameters $\boldsymbol{\theta}^*$ that maximize (??) (equivalently minimize negative of (??)) using any standard gradient-descent solver. The expression (??) gives the gradient of the cost function and is used to iteratively move the candidate hyperparameter $\boldsymbol{\theta}$ to a lower-cost location. However, since the marginal likelihood is nonconvex there may be multiple local minima.

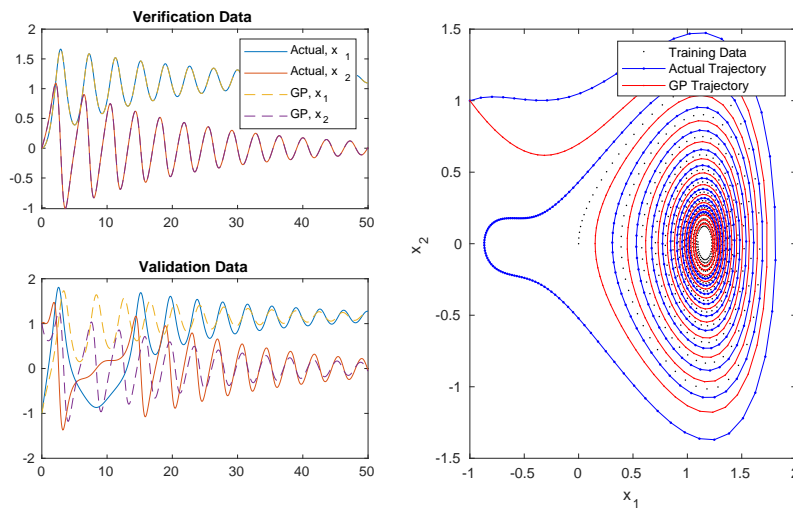
Cross Validation. As with any machine learning algorithm, the trained model should be *cross-validated*. The basic idea is to split the data into a training set used to learn the model and a validation set to test the learned model (i.e., estimate its performance on previously unobserved data). However, drawbacks of this approach are that only a fraction of the full data can be used for training and, if the validation data set is small, then the estimated performance can have a large variance. To address this issue k -fold cross validation is often used. The data is split into k disjoint subsets and then trained separately, k times, each time leaving out a different subset of the data for validation. The hyperparameters learned from each of the model are then combined, for example, by averaging.

Model Selection. Our previous discussion concerned hyperparameter optimization under a known model. However, we may wish to evaluate several different models (e.g., GP models with different kernels) from a library of hypothesized candidate models. This problem can be viewed as a nested optimization problem wherein at the outer level we optimize the choice of model and at the inner level we optimize the chosen model's hyperparameters.

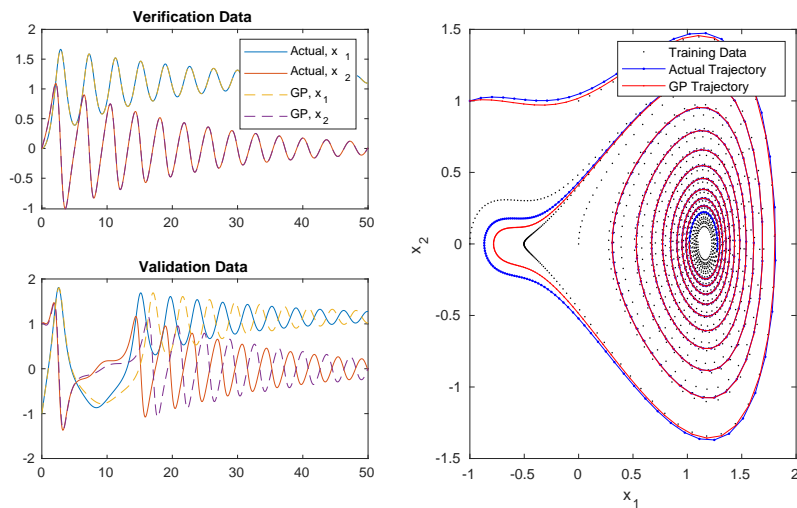
Example (Double Well, Strogatz Ch. 12.5). A well known example of a bi-stable system is the model:

$$\ddot{x} + \delta \dot{x} - x + x^3 = F \cos \omega t \quad (39)$$

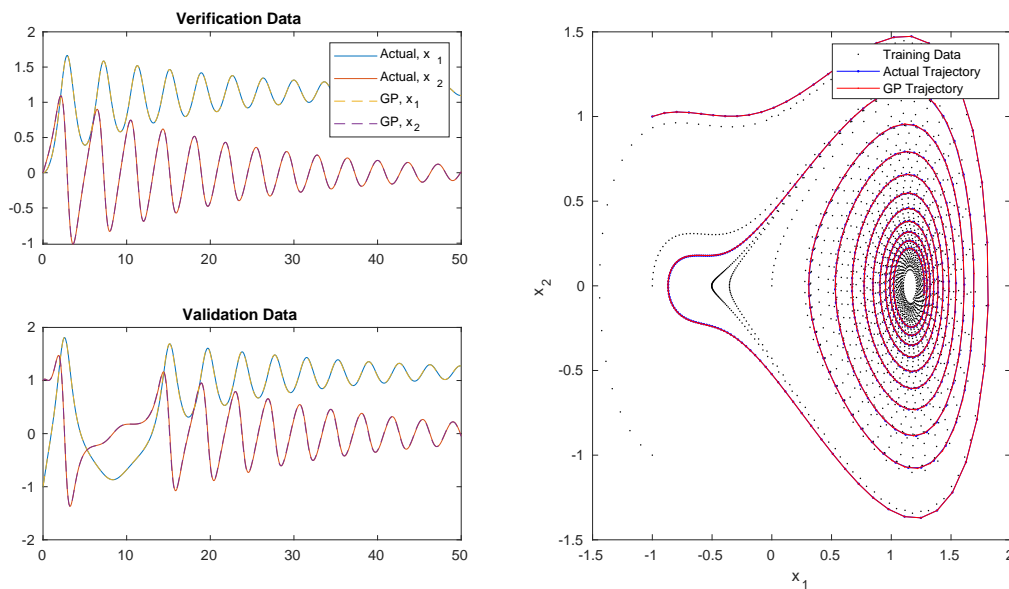
In the following we make the assumption that $\omega = 0$ (i.e., a constant force) and assume the parameter $\delta = 0.1$ and $F = 0.4$. The system is converted into a two-state model with $x_1 = x$ and $x_2 = \dot{x}$ and simulated for 50 seconds using ode45 in MATLAB. To generate the training data we begin by choosing the initial condition $(x_1, x_2) = (0, 0)$. For validation we use $(x_1, x_2) = (-1, 1)$. The GP is trained with a squared exponential kernel with hyperparameters of unit length and unit variance. The resulting GP-predicted trajectory is shown below. Since the training data does not cover the space of the validation trajectory we see a mismatch in the initial response.



To alleviate this issue we can add additional training data for the initial conditions $(x_1, x_2) = (0, 1)$ and $(x_1, x_2) = (-1, 0)$. The results are shown below. There is an improvement however the predicted response appears out of phase with the actual response.



Adding two more initials to the training data set, $(x_1, x_2) = (1, 1)$ and $(x_1, x_2) = (-1, -1)$, appears to alleviate this issue for this validation data set. For the GP predictions to be accurate over a wide range of initial conditions additional training data is needed.



Gaussian Processes: Weight-Space View

In the functions-space view described above GP-based predictions were obtained by conditioning the joint probability distribution of the training data and prediction points on the training data to obtain a posterior distribution. Another approach of obtaining the regression is the so-called weight-space view in which the GP regression is viewed as an optimization problem to identify optimal weights. To simplify the discussion suppose the mean is a known constant m and again.

The simple kriging estimator $f_*(\mathbf{x}_*)$ at site \mathbf{x}_* is a linear combination of k variables at sites \mathbf{x}_i for $i = 1, \dots, k$:

$$f_*(\mathbf{x}_*) = m + \sum_{i=1}^k \lambda_i (f(\mathbf{x}_i) - m) \quad (40)$$

where λ_i is a set of weights. The estimation variance at site \mathbf{x}_* is

$$\sigma^2(\mathbf{x}_*) = \text{Var}(f_*(\mathbf{x}_*) - f(\mathbf{x}_*)) \quad (41)$$

$$= \text{Var}\left(m + \sum_{i=1}^k \lambda_i (f(\mathbf{x}_i) - m) - f(\mathbf{x}_*)\right) \quad (42)$$

$$= \text{Var}\left(\sum_{i=1}^k \lambda_i (f(\mathbf{x}_i) - m) - [f(\mathbf{x}_*) - m]\right) \quad (43)$$

$$= \text{Var}\left(\sum_{i=1}^k \lambda_i \epsilon(\mathbf{x}_i) - \epsilon(\mathbf{x}_*)\right). \quad (44)$$

where $\epsilon(\mathbf{x}) = f(\mathbf{x}) - E[f(\mathbf{x})] = f(\mathbf{x}) - m$ is the residual. Note that

$$E[\epsilon(\mathbf{x})] = E[f(\mathbf{x}) - m] = 0 \quad (45)$$

Changing the index from $i = 1$ to $i = 0$ and defining $\lambda_0 = -1$:

$$\sigma^2(\mathbf{x}_*) = \text{Var}\left(\sum_{i=1}^k \lambda_i \epsilon(\mathbf{x}_i) + \lambda_0 \epsilon(\mathbf{x}_*)\right) \quad (46)$$

$$= \text{Var}\left(\sum_{i=0}^k \lambda_i \epsilon(\mathbf{x}_i)\right) \quad (47)$$

Using the definition of variance the expression (??) can be expanded and factored as

$$\sigma^2(\mathbf{x}_*) = E \left[\left(\sum_{i=0}^k \lambda_i \epsilon(\mathbf{x}_i) \right)^2 \right] - \left(E \left[\sum_{i=0}^k \lambda_i \epsilon(\mathbf{x}_i) \right] \right)^2 \quad (48)$$

$$= E \left[\left(\sum_{i=0}^k \lambda_i \epsilon(\mathbf{x}_i) \right) \left(\sum_{j=0}^k \lambda_j \epsilon(\mathbf{x}_j) \right) \right] - \left(\left[\sum_{i=0}^k \lambda_i E[\epsilon(\mathbf{x}_i)] \right] \right)^2 \quad (49)$$

$$= E \left[\left(\sum_{i=0}^k \lambda_i \epsilon(\mathbf{x}_i) \right) \left(\sum_{j=0}^k \lambda_j \epsilon(\mathbf{x}_j) \right) \right] - \left(\left[\sum_{i=0}^k \lambda_i E[\epsilon(\mathbf{x}_i)] \right] \right)^2 \quad (50)$$

$$= E \left[\sum_{i=0}^k \sum_{j=0}^k \lambda_i \lambda_j \epsilon(\mathbf{x}_i) \epsilon(\mathbf{x}_j) \right] - \left(\left[\sum_{i=0}^k \lambda_i \underbrace{E[\epsilon(\mathbf{x}_i)]}_{=0} \right] \right)^2 \quad (51)$$

$$= \sum_{i=0}^k \sum_{j=0}^k \lambda_i \lambda_j E[\epsilon(\mathbf{x}_i) \epsilon(\mathbf{x}_j)] \quad (52)$$

$$= \sum_{i=0}^k \sum_{j=0}^k \lambda_i \lambda_j E[(f(\mathbf{x}_i) - m)(f(\mathbf{x}_j) - m)] \quad (53)$$

$$= \sum_{i=0}^k \sum_{j=0}^k \lambda_i \lambda_j \left\{ E[f(\mathbf{x}_i) f(\mathbf{x}_j)] - \underbrace{m}_{=m} E[f(\mathbf{x}_i)] - \underbrace{m}_{=m} E[f(\mathbf{x}_j)] + m^2 \right\} \quad (54)$$

$$= \sum_{i=0}^k \sum_{j=0}^k \lambda_i \lambda_j \left\{ \underbrace{E[f(\mathbf{x}_i) f(\mathbf{x}_j)] - m^2}_{\text{Cov}(f(\mathbf{x}_i), f(\mathbf{x}_j))} \right\} \quad (55)$$

$$= \sum_{i=0}^k \sum_{j=0}^k \lambda_i \lambda_j k(\mathbf{x}_i, \mathbf{x}_j) \quad (56)$$

Change the limits of the summation back to $i = 1$ and $j = 1$, recall that k is symmetric:

$$\sigma^2(\mathbf{x}_*) = \sum_{i=0}^k \sum_{j=0}^k \lambda_i \lambda_j k(\mathbf{x}_i, \mathbf{x}_j) \quad (57)$$

$$\begin{aligned} &= \lambda_0 \lambda_0 k(\mathbf{x}_*, \mathbf{x}_*) + \lambda_0 \lambda_1 k(\mathbf{x}_*, \mathbf{x}_1) + \cdots + \lambda_0 \lambda_k k(\mathbf{x}_*, \mathbf{x}_k) + \\ &\quad \lambda_1 \lambda_0 k(\mathbf{x}_1, \mathbf{x}_*) + \lambda_1 \lambda_1 k(\mathbf{x}_1, \mathbf{x}_1) + \cdots + \lambda_1 \lambda_k k(\mathbf{x}_1, \mathbf{x}_k) + \\ &\quad \vdots \\ &\quad \lambda_k \lambda_0 k(\mathbf{x}_k, \mathbf{x}_*) + \lambda_k \lambda_1 k(\mathbf{x}_k, \mathbf{x}_1) + \cdots + \lambda_k \lambda_k k(\mathbf{x}_k, \mathbf{x}_k) \end{aligned} \quad (58)$$

$$= \sum_{i=1}^k \sum_{j=1}^k \lambda_i \lambda_j k(\mathbf{x}_i, \mathbf{x}_j) + 2 \underbrace{\lambda_0}_{=-1} \sum_{i=1}^k \lambda_i k(\mathbf{x}_i, \mathbf{x}_*) + \lambda_0^2 k(\mathbf{x}_*, \mathbf{x}_*) \quad (59)$$

$$= \sum_{i=1}^k \sum_{j=1}^k \lambda_i \lambda_j k(\mathbf{x}_i, \mathbf{x}_j) - 2 \sum_{i=1}^k \lambda_i k(\mathbf{x}_i, \mathbf{x}_*) + k(\mathbf{x}_*, \mathbf{x}_*) \quad (60)$$

Define the matrix

$$\mathbf{K} = \mathbf{K}(\mathbf{x}_1, \dots, \mathbf{x}_k) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_1) & \cdots & k(\mathbf{x}_k, \mathbf{x}_1) \\ k(\mathbf{x}_1, \mathbf{x}_2) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_k, \mathbf{x}_2) \\ \cdots & \cdots & \cdots & \cdots \\ k(\mathbf{x}_1, \mathbf{x}_k) & k(\mathbf{x}_2, \mathbf{x}_k) & \cdots & k(\mathbf{x}_k, \mathbf{x}_k) \end{bmatrix}, \quad (61)$$

the vector of weights $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \dots, \lambda_k]^T$, and the vector of covariance of sample sites with the estimation site

$$\boldsymbol{\nu} = [k(\mathbf{x}_1, \mathbf{x}_*) \quad k(\mathbf{x}_2, \mathbf{x}_*) \quad \cdots \quad k(\mathbf{x}_k, \mathbf{x}_*)]^T, \quad (62)$$

then

$$\sigma^2(\mathbf{x}_*; \mathbf{x}_1, \dots, \mathbf{x}_k, \boldsymbol{\lambda}) = \boldsymbol{\lambda}^T \mathbf{K} \boldsymbol{\lambda} - 2\boldsymbol{\lambda}^T \boldsymbol{\nu} + k(\mathbf{x}_*, \mathbf{x}_*). \quad (63)$$

To find the optimal weights $\boldsymbol{\lambda}^*$ that minimize $\sigma^2(\mathbf{x}_*)$ we can formulate the following optimization problem:

$$\text{minimize } \sigma^2(\boldsymbol{\lambda}; \mathbf{x}_*, \mathbf{x}_1, \dots, \mathbf{x}_k) \quad (64)$$

where we have rewritten $\sigma^2(\boldsymbol{\lambda}; \mathbf{x}_*, \mathbf{x}_1, \dots, \mathbf{x}_k)$ to emphasize that $\boldsymbol{\lambda}$ is the optimization variable and $\mathbf{x}_*, \mathbf{x}_1, \dots, \mathbf{x}_k$ are fixed parameters defining a problem instance.

The optimal weights $\boldsymbol{\lambda}^*$ are found by minimizing the variance (??). Since (??) is a linear equation it can only has one critical point that occurs where the first derivative is zero

$$\frac{d\sigma^2}{d\boldsymbol{\lambda}} = 2\mathbf{K}\boldsymbol{\lambda} - 2\boldsymbol{\nu}. \quad (65)$$

Setting the above derivative to zero and solving for the optimal weights

$$\boldsymbol{\lambda}^* = \mathbf{K}^{-1}\boldsymbol{\nu} \quad (66)$$

Moreover, this point is a minimum if the Hessian

$$\frac{d^2\sigma^2}{d\boldsymbol{\lambda}^2} = 2\mathbf{K} \quad (67)$$

is positive definite. The matrix \mathbf{K} is positive definite if for all vectors $\mathbf{x} \in \mathbb{R}^k$ the scalar $\mathbf{x}^T \mathbf{K} \mathbf{x} = 0$. The simple kriging estimator (??) can now be rewritten with the optimal weights as

$$\begin{aligned} f_*(\mathbf{x}_*) &= m + \sum_{i=1}^k \lambda_i^* (f(\mathbf{x}_i) - m) \\ &= m + (\boldsymbol{\lambda}^*)^T \bar{\mathbf{f}} \end{aligned} \quad (68)$$

where $\bar{\mathbf{f}} = [f(\mathbf{x}_1) - m \quad f(\mathbf{x}_2) - m \quad \cdots \quad f(\mathbf{x}_k) - m]^T$. Note that $\boldsymbol{\nu}$ is analagous to the matrix $\mathbf{K}(\mathbf{X}, \mathbf{X}_*)$ in (??)–(??) and with with the substitution of $m = 0$ the expression for the predicted values is equivalent. Also, using (??).

and the corresponding variance is

$$\sigma^2(\mathbf{x}_*) = (\boldsymbol{\lambda}^*)^T \mathbf{K} \boldsymbol{\lambda}^* - 2(\boldsymbol{\lambda}^*)^T \boldsymbol{\nu} + k(\mathbf{x}_*, \mathbf{x}_*) \quad (69)$$

$$= (\mathbf{K}^{-1} \boldsymbol{\nu})^T \mathbf{K} \mathbf{K}^{-1} \boldsymbol{\nu} - 2(\mathbf{K}^{-1} \boldsymbol{\nu})^T \boldsymbol{\nu} + k(\mathbf{x}_*, \mathbf{x}_*) \quad (70)$$

$$= (\mathbf{K}^{-1} \boldsymbol{\nu})^T \boldsymbol{\nu} - 2(\mathbf{K}^{-1} \boldsymbol{\nu})^T \boldsymbol{\nu} + k(\mathbf{x}_*, \mathbf{x}_*) \quad (71)$$

$$= \boldsymbol{\nu}^T (\mathbf{K}^{-1})^T \boldsymbol{\nu} - 2\boldsymbol{\nu}^T (\mathbf{K}^{-1})^T \boldsymbol{\nu} + k(\mathbf{x}_*, \mathbf{x}_*) \quad (72)$$

$$= k(\mathbf{x}_*, \mathbf{x}_*) - \boldsymbol{\nu}^T \mathbf{K}^{-1} \boldsymbol{\nu} \quad (73)$$

$$\mu_{\mathbf{f}_*|\mathbf{f}} = \mathbf{K}(\mathbf{X}_*, \mathbf{X}) \mathbf{K}(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f} \quad (74)$$

$$\mathbf{P}_{\mathbf{f}_*|\mathbf{f}} = \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) - \mathbf{K}(\mathbf{X}_*, \mathbf{X}) \mathbf{K}(\mathbf{X}, \mathbf{X})^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \quad (75)$$

References

- [1] C. M. Bishop and N. M. Nasrabadi, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [2] C. K. Williams and C. E. Rasmussen, *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [3] R. A. Olea, *Geostatistics for Engineering and Earth Scientists*. Springer Science+Buisness Media, 1999.