

AutoML: Gaussian Processes

Gaussian Process Prediction

Bernd Bischl Frank Hutter Lars Kotthoff
Marius Lindauer Joaquin Vanschoren

Motivation

- So far, we have learned how to **sample** from a Gaussian process prior.
- However, most of the time, we are not interested in drawing random functions from the prior. Instead, we usually like to implement the knowledge provided by the training data to predict values of f at a new test point \mathbf{x}_* .
- In what follows, we will investigate how to update Gaussian process prior (\rightarrow posterior process) and how to make predictions.

Gaussian Posterior Process and Prediction

Posterior Process I

- Let us distinguish between **observed training** inputs (also denoted by a design matrix \mathbf{X}), their corresponding values

$$\mathbf{f} = \left[f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(n)}) \right],$$

and one single **unobserved test** point $f_* = f(\mathbf{x}_*)$.

- Assuming a zero-mean GP prior $\mathcal{G}(\mathbf{0}, k(\mathbf{x}, \mathbf{x}'))$, we can assert that

$$\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{k}_* \\ \mathbf{k}_*^T & k_{**} \end{bmatrix}\right),$$

where, $\mathbf{K} = (k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}))_{i,j}$, $\mathbf{k}_* = [k(\mathbf{x}_*, \mathbf{x}^{(1)}), \dots, k(\mathbf{x}_*, \mathbf{x}^{(n)})]$ and $k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$.

Posterior Process II

() A General Rule of Conditioning for Gaussian Random Variables*

Let $\mathbf{z}_1 \in \mathbb{R}^{m_1}$, $\mathbf{z}_2 \in \mathbb{R}^{m_2}$, and $\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2)$. If the m -dimensional Gaussian vector \mathbf{z} can be partitioned as

$$\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2), \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix},$$

then the conditional distribution $\mathbf{a} = \mathbf{z}_2 \mid \mathbf{z}_1$ will be a multivariate normal distribution:

$$\mathcal{N}(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\mathbf{a} - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12})$$

Posterior Process III

- Given that \mathbf{f} is observed, we can exploit the general rule (*) to obtain the following formula:

$$f_* \mid \mathbf{x}_*, \mathbf{X}, \mathbf{f} \sim \mathcal{N}(\mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{f}, \mathbf{k}_{**} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*).$$

- 💡 As the posterior is Gaussian, the maximum a-posteriori estimate (i.e., the mode of the posterior distribution) is:

$$\mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{f}$$

.

GP Prediction: Two Points I

To visualize the above idea, assume that we have observed a single training point $\mathbf{x} = -0.5$. Based on this point, we intend to make a prediction at the test point $\mathbf{x}_* = 0.5$.

- Under a zero-mean \mathcal{G} with $k(\mathbf{x}, \mathbf{x}') = \exp(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2)$, we assume (the covariance matrix has been computed):

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} 1 & 0.61 \\ 0.61 & 1 \end{bmatrix}\right).$$

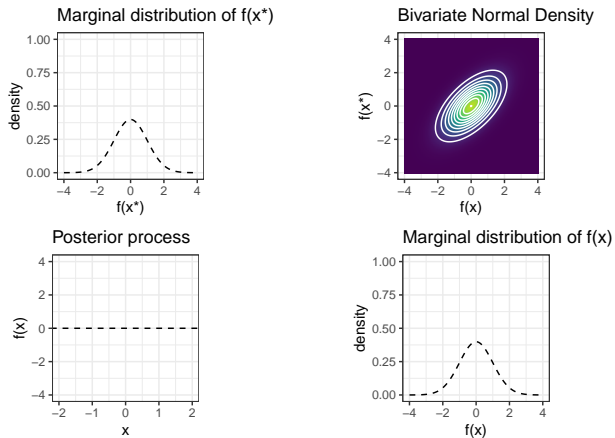
- Let us assume that we observe the point $f(\mathbf{x}) = 1$. We can compute the posterior distribution:

$$\begin{aligned} f_* \mid \mathbf{x}_*, \mathbf{x}, f &\sim \mathcal{N}(\mathbf{k}_*^T \mathbf{K}^{-1} f, k_{**} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*) \\ &\sim \mathcal{N}(0.61 \cdot 1 \cdot 1, 1 - 0.61 \cdot 1 \cdot 0.61) \\ &\sim \mathcal{N}(0.61, 0.6279) \end{aligned}$$

- The MAP-estimate for \mathbf{x}_* is $f(\mathbf{x}_*) = 0.61$, and the uncertainty estimate is 0.6279.

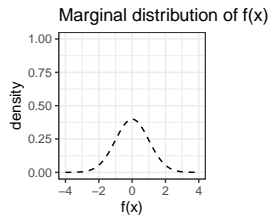
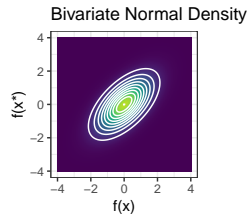
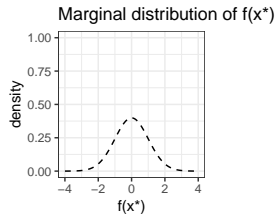
GP Prediction: Two Points II

The figures show the bivariate normal density as well as the corresponding marginals.



GP Prediction: Two Points III

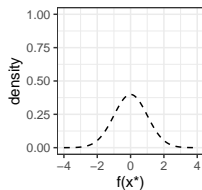
Now, assume that we observe the value of $f(\mathbf{x}) = 1$ corresponding to the training point $\mathbf{x} = -0.5$.



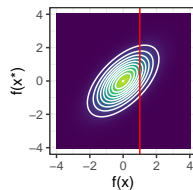
GP Prediction: Two Points IV

We condition the Gaussian on $f(\mathbf{x}) = 1$.

Marginal distribution of $f(x^*)$



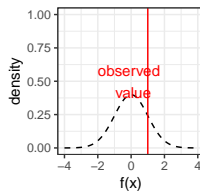
Bivariate Normal Density



Posterior process



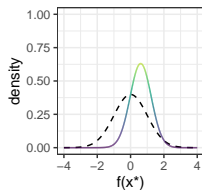
Marginal distribution of $f(x)$



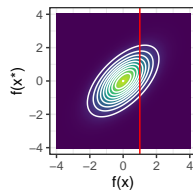
GP Prediction: Two Points V

We then compute the posterior distribution of $f(\mathbf{x}_*)$ given that $f(\mathbf{x}) = 1$.

Marginal distribution of $f(\mathbf{x}^*)$



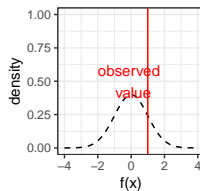
Bivariate Normal Density



Posterior process

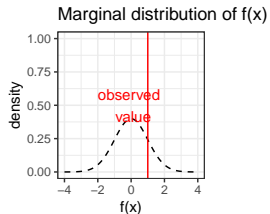
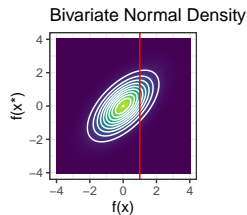
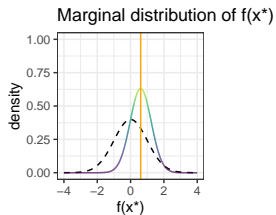


Marginal distribution of $f(\mathbf{x})$



GP Prediction: Two Points VI

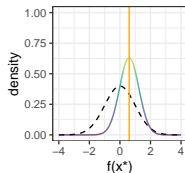
A possible predictor for f at \mathbf{x}_* is the MAP of the posterior distribution.



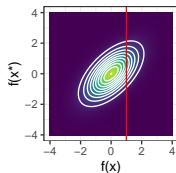
GP Prediction: Two Points VII

We can repeat this process for different \mathbf{x}_* and find the respective mean (grey line) and standard deviation (grey area). Note that the grey area is mean $\pm 2 \times$ standard deviation.

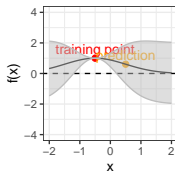
Marginal distribution of $f(x^*)$



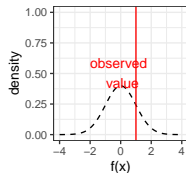
Bivariate Normal Density



Posterior process



Marginal distribution of $f(x)$



Posterior Process I

- The previous discussion was restricted to a single test point. However, one can generalize it to posterior processes with multiple unobserved test points:

$$\mathbf{f}_* = \left[f\left(\mathbf{x}_*^{(1)}\right), \dots, f\left(\mathbf{x}_*^{(n)}\right) \right].$$

- Under a zero-mean Gaussian process, we have:

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

where $\mathbf{K}_* = \left(k\left(\mathbf{x}^{(i)}, \mathbf{x}_*^{(j)}\right) \right)_{i,j}$ and $\mathbf{K}_{**} = k(\mathbf{x}_*^{(i)}, \mathbf{x}_*^{(j)})$.

Posterior Process II

- 💡 Similar to the single test point situation, to get the posterior distribution, we exploit the general rule of conditioning for Gaussians:

$$\mathbf{f}_* \mid \mathbf{X}_*, \mathbf{X}, \mathbf{f} \sim \mathcal{N}(\mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{f}, \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*).$$

- 💡 This formula enables us to talk about correlations among different test points and sample functions from the posterior process.

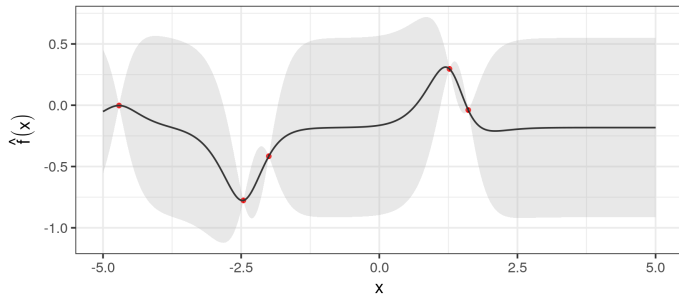
Properties of a Gaussian Process

GP as an Interpolator

- The “prediction” for a training point $\mathbf{x}^{(i)}$ is the exact function value $f(\mathbf{x}^{(i)})$. That is,

$$\mathbf{f} \mid \mathbf{X}, \mathbf{f} \sim \mathcal{N}(\mathbf{K}\mathbf{K}^{-1}\mathbf{f}, \mathbf{K} - \mathbf{K}^T\mathbf{K}^{-1}\mathbf{K}) = \mathcal{N}(\mathbf{f}, \mathbf{0}).$$

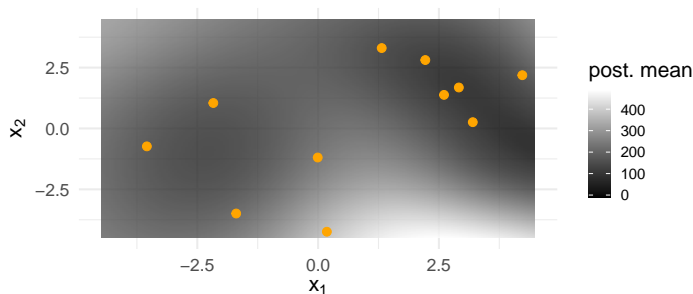
- Thus, a Gaussian process is a function **interpolator**.



After observing the training points (red), the posterior process (black) interpolates the training points.
($k(x, x')$ is Matérn with $\nu = 2.5$, the default for `DiceKriging::km`)

GP as a Spatial Model I

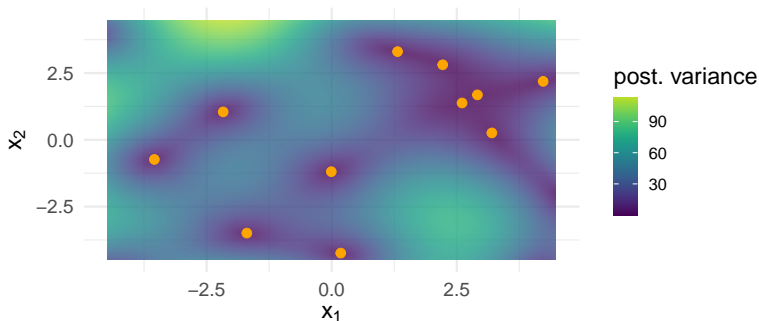
- The correlation among two outputs depends on the distance of the corresponding input points \mathbf{x} and \mathbf{x}' . For instance, the Gaussian covariance kernel is $k(\mathbf{x}, \mathbf{x}') = \exp\left(\frac{-\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right)$.
- Hence, close data points with high spatial similarity $k(\mathbf{x}, \mathbf{x}')$ enter into more strongly correlated predictions: $\mathbf{k}_*^\top \mathbf{K}^{-1} \mathbf{f}$ ($\mathbf{k}_* := (k(\mathbf{x}, \mathbf{x}^{(1)}), \dots, k(\mathbf{x}, \mathbf{x}^{(n)}))$).



Example: the posterior mean of a GP that is fitted with the Gaussian covariance kernel with $\ell = 1$.

GP as a Spatial Model II

- 💡 Posterior uncertainty increases if the new data points are far from the design points.
- 💡 The uncertainty is minimal at the design points, since the posterior variance is zero at these points.

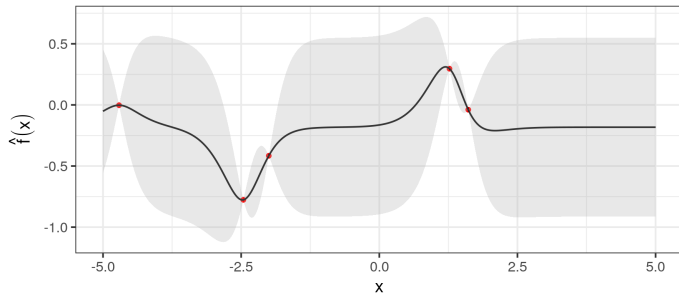


Example (continued): posterior variance

Noisy Gaussian Process

Noisy Gaussian Process I

- So far, we have implicitly assumed that we access the true function values $f(\mathbf{x})$.
- For the squared exponential kernel, for example, we had $\text{cov}(f(\mathbf{x}^{(i)}), f(\mathbf{x}^{(j)})) = 1$.
- Consequently, the posterior Gaussian process was an interpolator.



After observing the training points (red), the posterior process (black) interpolates the training points.
($k(x, x')$ is Matérn with $\nu = 2.5$, the default for `DiceKriging::km`)

Noisy Gaussian Process II

- However, in reality that is not often the case. Rather, we often only have access to a noisy version of the true function values:

$$y = f(\mathbf{x}) + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \sigma^2).$$

- Let us assume that $f(\mathbf{x})$ is still a Gaussian process. Then, we would have the following:

$$\begin{aligned} \text{cov}(y^{(i)}, y^{(j)}) &= \text{cov}\left(f(\mathbf{x}^{(i)}) + \epsilon^{(i)}, f(\mathbf{x}^{(j)}) + \epsilon^{(j)}\right) \\ &= \text{cov}\left(f(\mathbf{x}^{(i)}), f(\mathbf{x}^{(j)})\right) + 2 \cdot \text{cov}\left(f(\mathbf{x}^{(i)}), \epsilon^{(j)}\right) + \text{cov}\left(\epsilon^{(i)}, \epsilon^{(j)}\right) \\ &= k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) + \sigma^2 \delta_{ij}. \end{aligned}$$

💡 σ^2 is called **nugget**.

Noisy Gaussian Process III

- We can now derive the predictive distribution for the case of noisy observations.
- Assuming that f is modeled by a Gaussian process, the prior distribution of y is

$$\mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{pmatrix} \sim \mathcal{N}(\mathbf{m}, \mathbf{K} + \sigma^2 \delta_{ij}),$$

with

$$\mathbf{m} := \left(m(\mathbf{x}^{(i)}) \right)_i, \quad \mathbf{K} := \left(k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \right)_{i,j}.$$

Noisy Gaussian Process IV

We distinguish again between:

- Observed training points and their corresponding values, i.e, \mathbf{X} and y .
- Unobserved test points and their corresponding values, i.e, \mathbf{X}_* and \mathbf{f}_* .

and get:

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

Noisy Gaussian Process V

- Similar to the noise-free case, we condition according to the rule of conditioning for Gaussians to get the posterior distribution for the test outputs \mathbf{f}_* at \mathbf{X}_* :

$$\mathbf{f}_* \mid \mathbf{X}_*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mathbf{m}_{\text{post}}, \mathbf{K}_{\text{post}}),$$

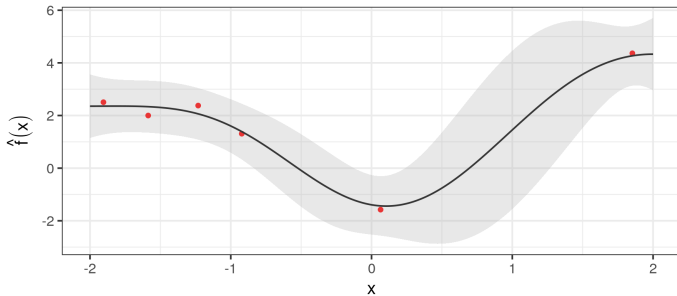
with

$$\begin{aligned}\mathbf{m}_{\text{post}} &= \mathbf{K}_*^T (\mathbf{K} + \sigma^2 \cdot \mathbf{I})^{-1} \mathbf{y} \\ \mathbf{K}_{\text{post}} &= \mathbf{K}_{**} - \mathbf{K}_*^T (\mathbf{K}^{-1} + \sigma^2 \cdot \mathbf{I}) \mathbf{K}_*.\end{aligned}$$

💡 This converts back to the noise-free formula if $\sigma^2 = 0$.

Noisy Gaussian Process VI

- The noisy Gaussian process is not an interpolator any more.
- A larger nugget term leads to a wider “band” around the observed training points.
- In general, the effect of the nugget term is estimated during training (see the next section).



After observing the training points (red), we have a nugget-band around the observed points.
($k(x, x')$ is the squared exponential)

Decision Theory for Gaussian Processes

Risk Minimization for Gaussian Processes I

In machine learning, we usually choose a loss function and try to minimize the empirical risk:

$$\mathcal{R}_{\text{emp}}(f) := \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)})),$$

as an approximation to the theoretical risk:

$$\mathcal{R}(f) := \mathbb{E}_{xy}[L(y, f(\mathbf{x}))] = \int L(y, f(\mathbf{x})) d\mathbb{P}_{xy}.$$

- How does the theory of Gaussian processes fit into this scenario?
- What if we were looking for predictions that are optimal w.r.t. a certain loss function?

Risk Minimization for Gaussian Processes II

- The theory of Gaussian process provides us with a posterior distribution, i.e, $p(y \mid \mathcal{D})$.
- To make a prediction at a test point \mathbf{x}_* , we can approximate the theoretical risk by exploiting the posterior distribution:

$$\mathcal{R}(y_* \mid \mathbf{x}_*) \approx \int L(\tilde{y}_*, y_*) p(\tilde{y}_* \mid \mathbf{x}_*, \mathcal{D}) d\tilde{y}_*.$$

- The optimal prediction w.r.t the loss function is then:

$$\hat{y}_* \mid \mathbf{x}_* = \arg \min_{y_*} \mathcal{R}(y_* \mid \mathbf{x}_*)$$