

GAUSSIAN POSTERIOR PROCESS AND PREDICTION

Introduction to Machine Learning

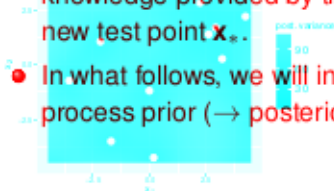


Gaussian Process Prediction

- So far, we have learned how to sample from a GP prior.
- However, most of the time, we are not interested in drawing random functions from the prior. Instead, we usually like to use 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 the Gaussian process prior (\rightarrow posterior process) and how to make predictions.

Learning goals

- Know how to derive the posterior process
- GPs are interpolating and spatial models
- Model noise via a nugget term



GAUSSIAN POSTERIOR PROCESS AND PREDICTION



Gaussian Posterior Process and Prediction

- However, most of the time, we are not interested in drawing random functions from the prior. Instead, we usually like to use 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 the Gaussian process prior (\rightarrow posterior process) and how to make predictions.

POSTERIOR PROCESS

- Let us now distinguish between observed training inputs, also denote by a design matrix \mathbf{X} , and the corresponding observed values

Gaussian Posterior Process and Prediction

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

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

- We now want to infer the distribution of $f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}$.

$$f_* = f(\mathbf{x}_*)$$

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

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

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



POSTERIOR PROCESS / 2

- Given that \mathbf{f} is observed, we can apply the general rule for condition (*) of Gaussian random variables and obtain the following formula:

$$\mathbf{f} = \begin{bmatrix} f(\mathbf{x}^{(1)}) \\ \vdots \\ f(\mathbf{x}^{(n)}) \end{bmatrix}$$

and one single observed point f_* at \mathbf{x}_* with $\mathbf{f}_* = k_*^T \mathbf{K}^{-1} \mathbf{f}$.

- We now want to infer the distribution of $f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}$.
- As the posterior is a Gaussian, the maximum a-posteriori estimate, i.e. the mode of the posterior distribution, is $k_*^T \mathbf{K}^{-1} \mathbf{f}$.

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

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

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



POSTERIOR PROCESS / 3

- (*) General rule for condition of Gaussian random variables: for condition (*) of Gaussian random variables and obtain the

If the m -dimensional Gaussian vector $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ can be partitioned with $\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2)$ where \mathbf{z}_1 is m_1 -dimensional and \mathbf{z}_2 is m_2 -dimensional, and:

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

- As the posterior is a Gaussian, the maximum a-posteriori estimate, i.e. the mode of the posterior distribution, is $\mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{f}$.

then the conditioned distribution of $\mathbf{z}_2 | \mathbf{z}_1 = \mathbf{a}$ is a multivariate normal

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



GP PREDICTION: TWO POINTS

Let us visualize this by a simple example:

- Assume we observed a single training point $\mathbf{x} = -0.5$, and want to make a prediction at a test point $\mathbf{x}_* = 0.5$.
- Under a zero-mean GP with $k(\mathbf{x}, \mathbf{x}') = \exp(-\frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|^2)$, we compute the cov-matrix:

$$\begin{pmatrix} f \\ f_* \end{pmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$$

then the conditioned distribution of $\mathbf{z}_2 \mid \mathbf{z}_1 = \mathbf{a}$ is a multivariate normal

- Assume that we observe the point $f(\mathbf{x}) = 1$.
- We 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.



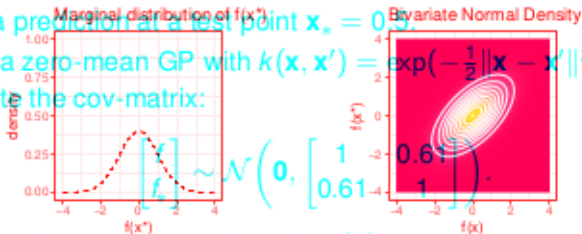
GP PREDICTION: TWO POINTS

Shows is the bivariate normal density, and the respective marginals.

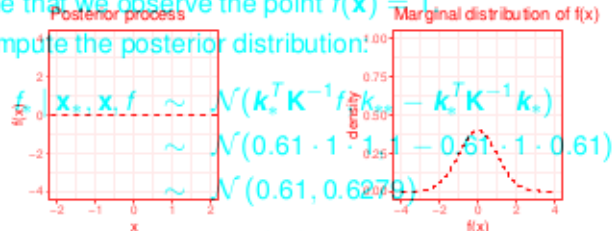


- Assume we observed a single training point $\mathbf{x} = -0.5$, and want to make a prediction at a test point $\mathbf{x}_* = 0.5$.

- Under a zero-mean GP with $k(\mathbf{x}, \mathbf{x}') = \exp(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2)$, we compute the cov-matrix:



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



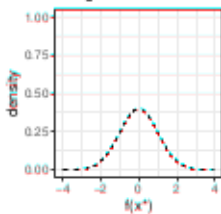
- The MAP-estimate for \mathbf{x}_* is $f(\mathbf{x}_*) = 0.61$, and the uncertainty estimate is 0.6279.

GP PREDICTION: TWO POINTS

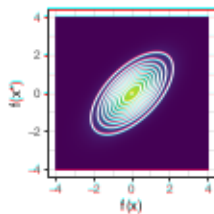
Assume we observed $f(x) = 1$ for the training point $x = -0.5$.



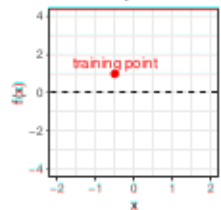
Marginal distribution of $f(x^*)$



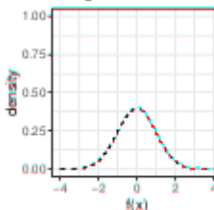
Bivariate Normal Density



Posterior process



Marginal distribution of $f(x)$

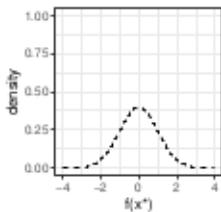


GP PREDICTION: TWO POINTS

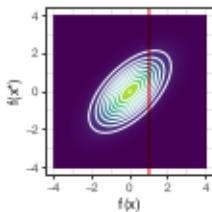
We condition the Gaussian on $f(x) = 1$ training point $x = -0.5$.



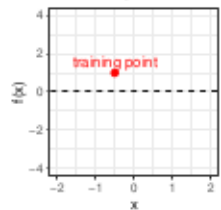
Marginal distribution of $f(x^*)$



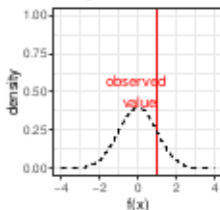
Bivariate Normal Density



Posterior process



Marginal distribution of $f(x)$

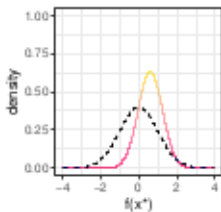


GP PREDICTION: TWO POINTS

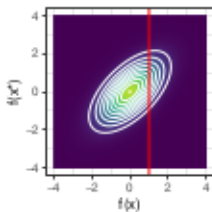
We compute the posterior distribution of $f(x_*)$ given that $f(x) = 1$.



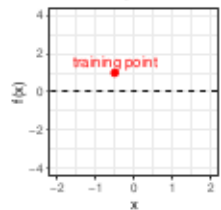
Marginal distribution of $f(x^*)$



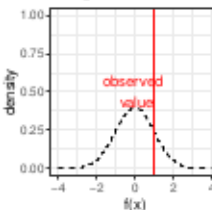
Bivariate Normal Density



Posterior process

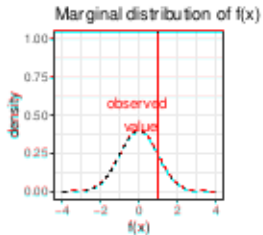
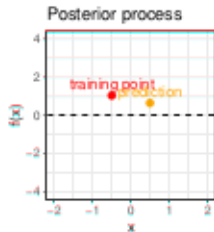
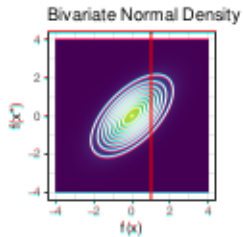
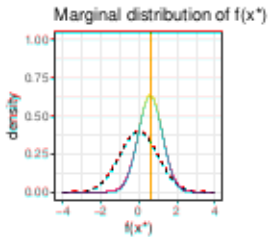


Marginal distribution of $f(x)$



GP PREDICTION: TWO POINTS

A possible predictor for f at x_* is the MAP (of the posterior distribution).

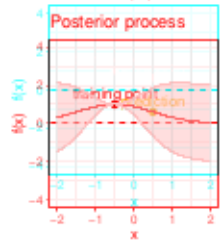
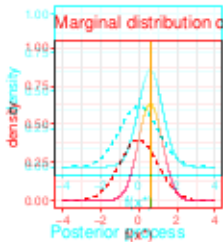


GP PREDICTION: TWO POINTS

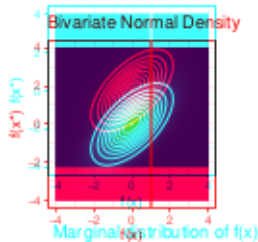
We can do this for different values x_* and show the respective mean (grey line) and standard deviations (grey area is mean $\pm 2 \cdot$ posterior standard deviation).



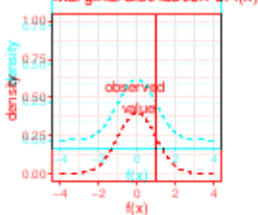
Marginal distribution of $f(x^*)$



Bivariate Normal Density



Marginal distribution of $f(x)$

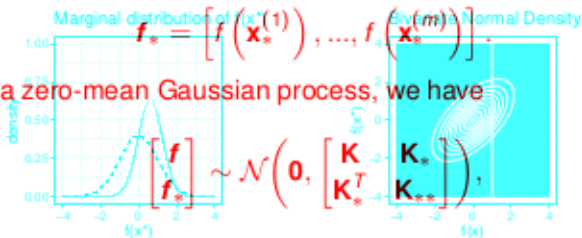


POSTERIOR PROCESS POINTS

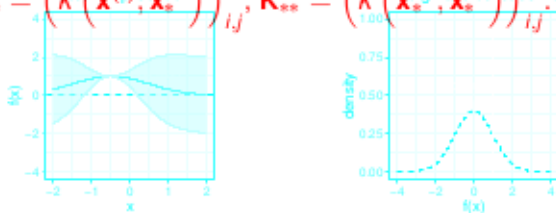
We can generalize the formula for the posterior process for multiple unobserved test points.



- Under a zero-mean Gaussian process, we have



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



POSTERIOR PROCESS / 2

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

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

- Under a zero-mean Gaussian process, we have
- This formula enables us to talk about correlations among different test points and sample functions from the posterior process.

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

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



POSTERIOR PROCESS

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

Properties of a Gaussian Process

$$f_* | \mathcal{X}_*, \mathcal{X}, Y \sim \mathcal{N}(K_* K^{-1} Y, K_{**} - K_* K^{-1} K_*).$$

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



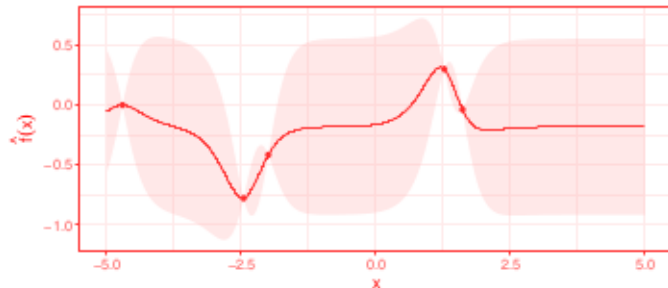
GP AS INTERPOLATOR

The "prediction" for a training point $\mathbf{x}^{(i)}$ is the exact function value $f(\mathbf{x}^{(i)})$

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

Properties of a Gaussian Process

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`)

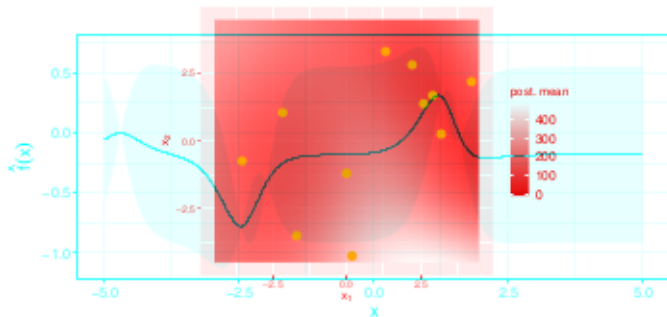


GP AS A SPATIAL MODEL

The "prediction" for a training point $\mathbf{x}^{(i)}$ is the exact function value $f(\mathbf{x}^{(i)})$

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(\frac{-\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right)$$

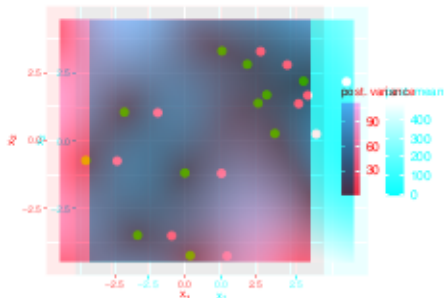
Hence, close data points with high spatial similarity $k(\mathbf{x}, \mathbf{x}')$ enter into more strongly correlated predictions: $\mathbf{k}_*^\top \mathbf{K}^{-1} \mathbf{f}$ ($\mathbf{k}_* := (k(\mathbf{x}, \mathbf{x}^{(1)}), \dots, k(\mathbf{x}, \mathbf{x}^{(n)}))$).
Thus, a Gaussian process is a function interpolator.



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

GP AS A SPATIAL MODEL / 2

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

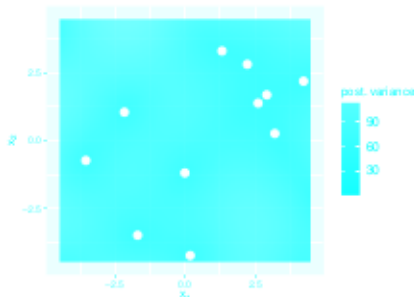


Example (continued): Posterior mean of a GP that was fitted with the Gaussian covariance kernel with $\ell = 1$.

GP AS A SPATIAL MODEL

- 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 the Gaussian

Noisy Gaussian Process



Example (continued): Posterior variance.



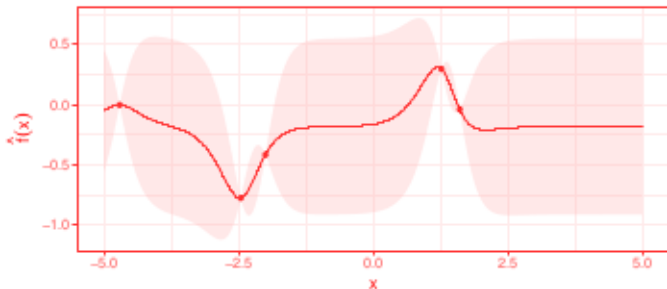
NOISY GAUSSIAN PROCESS

- So far, we implicitly assumed that we had access to the true function value $f(\mathbf{x})$.
- For the squared exponential kernel, for example, we have

Noisy Gaussian Process

$$\text{Cov}(f(\mathbf{x}^{(i)}), f(\mathbf{x}^{(i)})) = 1.$$

- As a result, the posterior Gaussian process is an interpolator:



After observing the training points (red), the posterior process (black) interpolates the training points.
($k(x,x')$ is Mat rn with $\nu = 2.5$, the default for DiceKriging: 3m)



NOISY GAUSSIAN PROCESS / 2

- In reality, however, this is often not the case.
- So far, we implicitly assumed that we had access to the true function $f(\mathbf{x})$.
- We often only have access to a noisy version of the true function value.
- For the squared exponential kernel, for example, we have

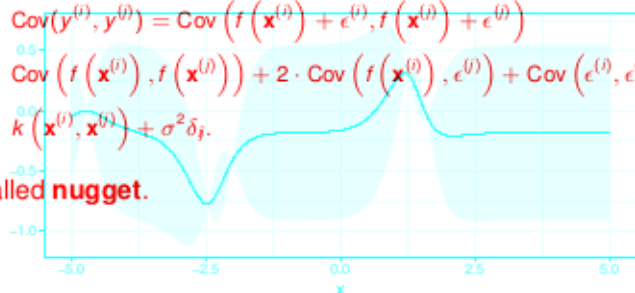
$$y = f(\mathbf{x}) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2).$$

$$\text{Cov}(f(\mathbf{x}^{(i)}), f(\mathbf{x}^{(j)})) = 1.$$

- Let us still assume that $f(\mathbf{x})$ is a Gaussian process.
- Then,
- As a result, the posterior Gaussian process is an interpolator:

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

- σ^2 is called **nugget**.



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



NOISY GAUSSIAN PROCESS / 3

- Let us now derive the predictive distribution for the case of noisy observations.
- We often only have access to a noisy version of the true function
- The prior distribution of y , assuming that f is modeled by a Gaussian process is then $f(\mathbf{x}) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$.

- Let us still assume that $f(\mathbf{x})$ is a Gaussian process.
- Then,

$$\begin{aligned} \text{Cov}(y^{(i)}, y^{(j)}) &= \text{Cov}\left(f(\mathbf{x}^{(i)}) + \epsilon^{(i)}, f(\mathbf{x}^{(j)}) + \epsilon^{(j)}\right) \\ &= \text{Cov}\left(f(\mathbf{x}^{(i)}), f(\mathbf{x}^{(j)})\right) + 2 \cdot \text{Cov}\left(f(\mathbf{x}^{(i)}), \epsilon^{(j)}\right) + \text{Cov}\left(\epsilon^{(i)}, \epsilon^{(j)}\right) \end{aligned}$$

with $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) + \sigma^2 \delta_j$.

- σ^2 is called **nugget** $m(\mathbf{x}^{(i)})$, $\mathbf{K} := \left(k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})\right)_{i,j}$.



NOISY GAUSSIAN PROCESS / 4

- We distinguish again between
 - observed training points \mathbf{X} , \mathbf{y} , and
 - unobserved test inputs \mathbf{X}_* with unobserved values \mathbf{f}_*
- The joint Gaussian process is then

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma^2 \mathbf{I}_n & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{bmatrix} \right).$$
$$\mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{pmatrix} \sim \mathcal{N}(\mathbf{m}, \mathbf{K} + \sigma^2 \mathbf{I}_n),$$

with

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



NOISY GAUSSIAN PROCESS / 5

- Similarly to the noise-free case, we condition according to the rule of conditioning for Gaussians to get the posterior distribution for the test outputs f_* at X_* :
- observed training points X, y , and
- unobserved test inputs X_* with unobserved values f_* .

and get

$$f_* \mid X_*, X, y \sim \mathcal{N}\left(0, \begin{bmatrix} K + \sigma^2 I_n & K_* \\ K_*^T & K_{**} \end{bmatrix}\right).$$

with

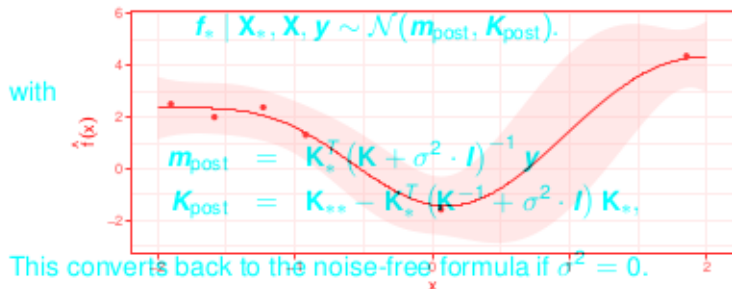
$$m_{\text{post}} = K_*^T (K + \sigma^2 \cdot I)^{-1} y$$
$$K_{\text{post}} = K_{**} - K_*^T (K^{-1} + \sigma^2 \cdot I) K_*$$

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



NOISY GAUSSIAN PROCESS / 6

- The noisy Gaussian process is not an interpolator any more.
- A larger nugget term leads to a wider "band" around the observed training points.
- The nugget term is estimated during training.



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

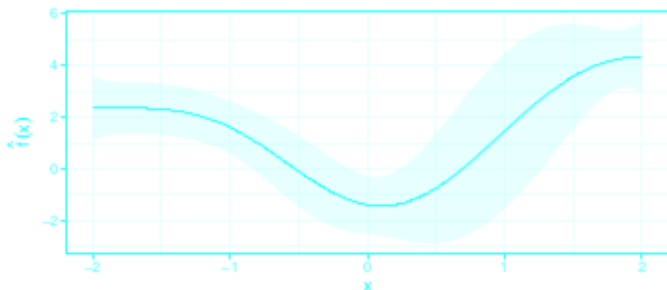
After observing the training points (red), we have a nugget-band around the observed points.
($k(x, x')$ is the squared exponential)

NOISY GAUSSIAN PROCESS

- The noisy Gaussian process is not an interpolator any more.
- A larger nugget term leads to a wider “band” around the observed training points.

Decision Theory for Gaussian Processes

- The nugget term is estimated during testing.



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



RISK MINIMIZATION FOR GAUSSIAN PROCESSES

/ 2



In the theory of Gaussian process, we can choose a loss function and minimize the empirical risk

$$p(y | \mathcal{D})$$

If we now want to make a prediction at a test point \mathbf{x}_* , we approximate the theoretical risk in a different way, by using the posterior distribution as a theoretical risk

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

- The optimal predictor w.r.t the loss function is then:
- What if we want to make a prediction which is optimal w.r.t. a certain loss function? $\hat{y}_* | \mathbf{x}_* = \arg \min_{y_*} \mathcal{R}(y_* | \mathbf{x}_*)$.

RISK MINIMIZATION FOR GAUSSIAN PROCESSES

- The theory of Gaussian process gives us a posterior distribution

$$p(y | \mathcal{D})$$

- If we now want to make a prediction at a test point \mathbf{x}_* , we approximate the theoretical risk in a different way, by using the posterior distribution:

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

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

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