

# Gaussian Process - Definition

- **Definition:** A *Gaussian process* is a (potentially infinite) collection of random variables such that the joint distribution of any finite number of them is multivariate Gaussian:

$$f \sim GP(\mu, K)$$

where  $\mu(x)$  and  $K(x, x')$  are the mean and covariance function.

- Gaussian processes take a *nonparameteric* approach to regression. We select a *prior distribution* over the function  $f$  (a GP prior:  $f \sim GP(\mu, K)$ ) and condition this distribution on our observations, using the posterior distribution to make predictions.
- Gaussian processes are very *powerful* and leverage the many *convenient properties* of the Gaussian distribution to enable tractable inference.

# Regression

Consider the general *regression* problem. Here we have:

- an input domain  $\mathbf{X}$  (for example,  $\mathbb{R}^n$ , but in general anything),
- an unknown function  $f: \mathbf{X} \rightarrow \mathbb{R}$ , and
- and (perhaps noisy) observations of the function:  
 $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$ , where  $y_i = f(\mathbf{x}_i) + \varepsilon_i$ .

Our goal is to *predict* the value of the function  $f(\mathbf{X}_*)$  at some test locations  $\mathbf{X}_*$ .

A Gaussian process distribution on  $f$  is written

$$p(f) = GP(f; \mu, K),$$

and just like the multivariate Gaussian distribution, is parameterized by its first two functions:

- $\mathbb{E}[f] = \mu: \mathbf{X} \rightarrow \mathbb{R}$ , the *mean function*, and
- $\mathbb{E}[(f(x) - \mu(x))(f(x') - \mu(x')))] = K: \mathbf{X} \times \mathbf{X} \rightarrow \mathbb{R}$ , a positive semidefinite *covariance function* or *kernel*.

# GPs: Mean and covariance functions

- The mean function encodes the *central tendency* of the function, and is often assumed to be a constant (usually zero).
- The covariance function encodes information about the *shape* and *structure* we expect the function to have. A simple and very common example is the *squared exponential* covariance:

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|_2^2\right),$$

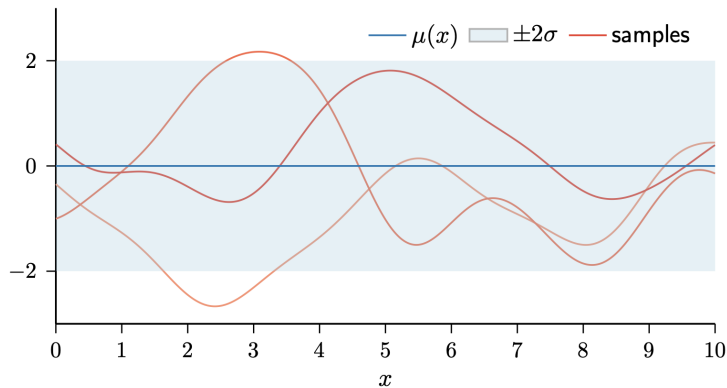
which encodes the notation that “nearby points should have similar function values.”

Suppose we have selected a GP prior  $GP(f; \mu, K)$  for the function  $f$ . Consider a finite set of points  $\mathbf{X} \subseteq \mathcal{X}$ . The GP prior on  $f$ , by definition, *implies* the following joint distribution on the associated function values  $\mathbf{f} = f(\mathbf{X})$ :

$$p(\mathbf{f} \mid \mathbf{X}) = \mathcal{N}(\mathbf{f}; \mu(\mathbf{X}), K(\mathbf{X}, \mathbf{X})).$$

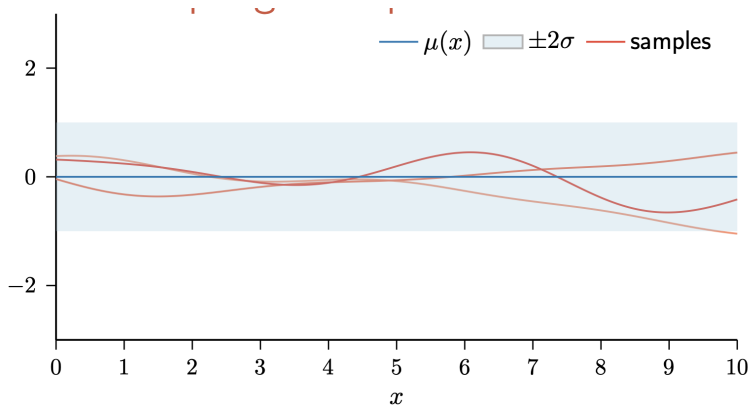
That is, we simply evaluate the mean and covariance functions at  $\mathbf{X}$  and take the associated multivariate Gaussian distribution.

# Prior: Sampling examples



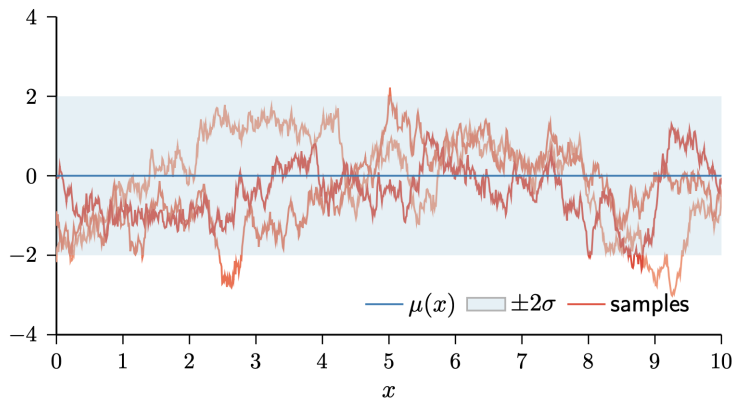
$$K = \exp\left(-\frac{1}{2}\|x - x'\|^2\right)$$

# Prior: Sampling examples



$$K = \lambda^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) \quad \lambda = \frac{1}{2}, \ell = 2$$

# Prior: Sampling examples



$$K = \exp(-\|x - x'\|)$$



# From the prior to the posterior

So far, we have constructed *prior* distributions over the function  $f$ .  
How do we *condition* our prior on some observations  $\mathcal{D} = (\mathbf{X}, \mathbf{f})$   
to *make predictions* about the value of  $f$  at some points  $\mathbf{X}_*$ ?

# From the prior to the posterior

We begin by writing the *joint distribution* between the training function values  $f(\mathbf{X}) = \mathbf{f}$  and the test function values  $f(\mathbf{X}_*) = \mathbf{f}_*$ :

$$p(\mathbf{f}, \mathbf{f}_*) = \mathcal{N} \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} ; \begin{bmatrix} \mu(\mathbf{X}) \\ \mu(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) & K(\mathbf{X}, \mathbf{X}_*) \\ K(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right) \cdots$$

# From the prior to the posterior

... we then *condition* this multivariate Gaussian on the known training values  $\mathbf{f}$ . We already know how to do that!

$$p(\mathbf{f}_* \mid \mathbf{X}_*, \mathcal{D}) = \mathcal{N}(\mathbf{f}_*; \mu_{f|\mathcal{D}}(\mathbf{X}_*), K_{f|\mathcal{D}}(\mathbf{X}_*, \mathbf{X}_*)),$$

where

$$\begin{aligned}\mu_{f|\mathcal{D}}(\mathbf{x}) &= \mu(\mathbf{x}) + K(\mathbf{x}, \mathbf{X})\mathbf{K}^{-1}(\mathbf{f} - \mu(\mathbf{X})) \\ K_{f|\mathcal{D}}(\mathbf{x}, \mathbf{x}') &= K(\mathbf{x}, \mathbf{x}') - K(\mathbf{x}, \mathbf{X})\mathbf{K}^{-1}K(\mathbf{X}, \mathbf{x}').\end{aligned}$$

The posterior distribution over  $\mathbf{f}$  is a Gaussian process!

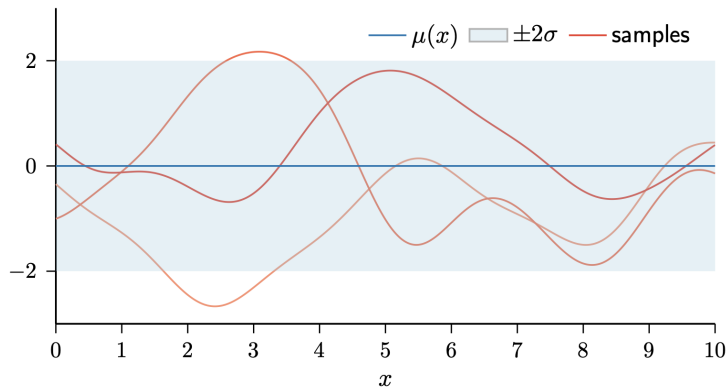
# The posterior mean

One way to understand the posterior mean function  $\mu_{f|\mathcal{D}}$  is as a *correction to the prior mean* consisting of a *weighted combination* of kernel functions, one for each training data point:

$$\begin{aligned}\mu_{f|\mathcal{D}}(\mathbf{x}) &= \mu(\mathbf{x}) + K(\mathbf{x}, \mathbf{X}) (K(\mathbf{X}, \mathbf{X}))^{-1} (\mathbf{f} - \mu(\mathbf{X})) \\ &= \mu(\mathbf{x}) + \sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \mathbf{x}),\end{aligned}$$

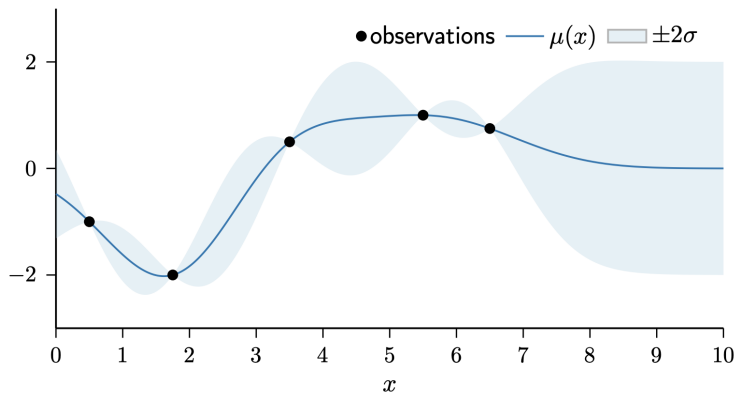
where  $\alpha_i = K(\mathbf{X}, \mathbf{X})^{-1} (f(\mathbf{x}_i) - \mu(\mathbf{x}_i))$ .

# Prior

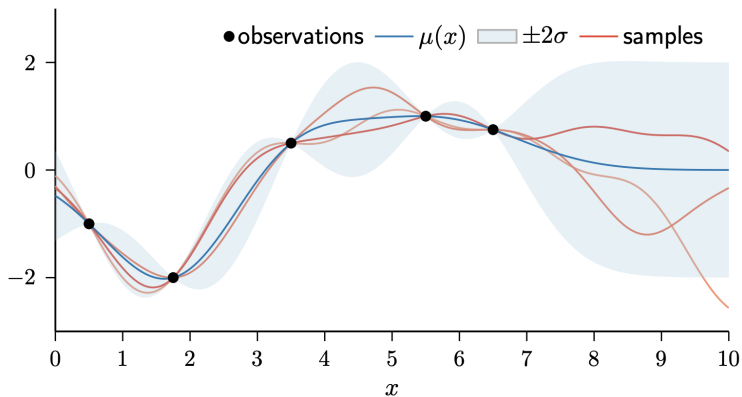


$$K = \exp\left(-\frac{1}{2}\|x - x'\|^2\right)$$

# Posterior example



# Posterior: Sampling



So far, we have assumed we can sample the function  $f$  *exactly*, which is uncommon in regression settings. How do we deal with *observation noise*?



We must create a *model* for our observations given the latent function. To begin, we will choose the simple iid, zero-mean additive Gaussian noise model:

$$y(\mathbf{x}) = f(\mathbf{x}) + \varepsilon,$$
$$p(\varepsilon \mid \mathbf{x}) = \mathcal{N}(\varepsilon; \mathbf{0}, \sigma^2);$$

combined we have

$$p(\mathbf{y} \mid \mathbf{f}) = \mathcal{N}(\mathbf{y}; \mathbf{f}, \sigma^2 \mathbf{I}).$$

To derive the posterior given *noisy observations*  $\mathcal{D}$ , we again write the joint distribution between the training function values  $\mathbf{y}$  and the test function values  $\mathbf{f}_*$ :

$$p(\mathbf{y}, \mathbf{f}_*) = \mathcal{N} \left( \begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix}; \begin{bmatrix} \mu(\mathbf{X}) \\ \mu(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & K(\mathbf{X}, \mathbf{X}_*) \\ K(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right) \cdots \quad (1)$$

... and *condition* as before.

$$p(\mathbf{f}_* \mid \mathbf{X}_*, \mathcal{D}) = \mathcal{N}(\mathbf{f}_*; \mu_{f|\mathcal{D}}(\mathbf{X}_*), K_{f|\mathcal{D}}(\mathbf{X}_*, \mathbf{X}_*)),$$

where

$$\begin{aligned}\mu_{f|\mathcal{D}}(\mathbf{x}) &= \mu(\mathbf{x}) + K(\mathbf{x}, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1}(\mathbf{y} - \mu(\mathbf{X})) \\ K_{f|\mathcal{D}}(\mathbf{x}, \mathbf{x}') &= K(\mathbf{x}, \mathbf{x}') - K(\mathbf{x}, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1}K(\mathbf{X}, \mathbf{x}').\end{aligned}$$

# Hyperparameters

- So far, we have assumed that the prior distribution on  $f$  follows a Gaussian Process.
- But this prior distribution *itself* has parameters, for example the length scale  $\ell$ , the output scale  $\lambda$ , and the noise variance  $\sigma^2$ . As parameters of a prior distribution, we call these *hyperparameters*.
- For convenience, we will write  $\theta$  to denote the vector of all hyperparameters of the model (including of  $\mu$  and  $K$ ).
- How do we *learn*  $\theta$ ?

# Marginal likelihood

Assume we have chosen a parameterized prior

$$p(\mathbf{f} \mid \theta) = GP(\mathbf{f}; \mu(\mathbf{x}; \theta), K(\mathbf{x}, \mathbf{x}'; \theta)).$$

We will measure the quality of the fit to our training data  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$  with the *marginal likelihood*, the probability of *observing the given data* under our prior:

$$p(\mathbf{y} \mid \mathbf{X}, \theta) = \int p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f} \mid \mathbf{X}, \theta) d\mathbf{f},$$

where we have *marginalized* the unknown function values  $\mathbf{f}$  (hence, marginal likelihood).

# Marginal likelihood: Evaluating

Thankfully, this is an integral we can do *analytically* under the Gaussian noise assumption!

$$\begin{aligned} p(\mathbf{y} \mid \mathbf{X}, \theta) &= \int p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f} \mid \mathbf{X}, \theta) d\mathbf{f}, \\ &= \int \overset{\text{iid noise}}{\mathcal{N}(\mathbf{y}; \mathbf{f}, \sigma^2 \mathbf{I})} \overset{\text{GP prior}}{\mathcal{N}(\mathbf{f}; \mu(\mathbf{X}; \theta), K(\mathbf{X}, \mathbf{X}; \theta))} d\mathbf{f} \\ &= \mathcal{N}(\mathbf{y}; \mu(\mathbf{X}; \theta), K(\mathbf{X}, \mathbf{X}; \theta) + \sigma^2 \mathbf{I}). \end{aligned}$$

(Convolutions of two Gaussians are Gaussian.)

# Marginal likelihood: Evaluating

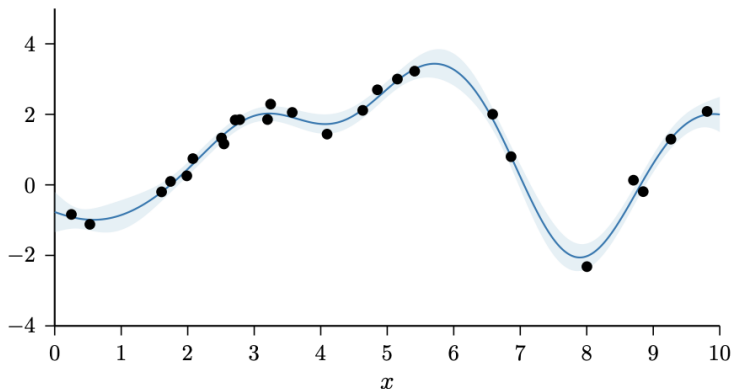
The log-likelihood of our data under the chosen prior are then (writing  $\mathbf{V} = (K(\mathbf{X}, \mathbf{X}; \theta) + \sigma^2 \mathbf{I})$ ):

$$\log p(\mathbf{y} \mid \mathbf{X}, \theta) =$$

$$-\frac{(\mathbf{y} - \overset{\text{data fit}}{\boldsymbol{\mu}})^\top \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})}{2} - \frac{\log \det \mathbf{V}}{2} - \frac{N \log 2\pi}{2}$$

The first term is large when the *data fit the model well*, and the second term is large when the *volume of the prior covariance is small*; that is, when the model is *simpler*.

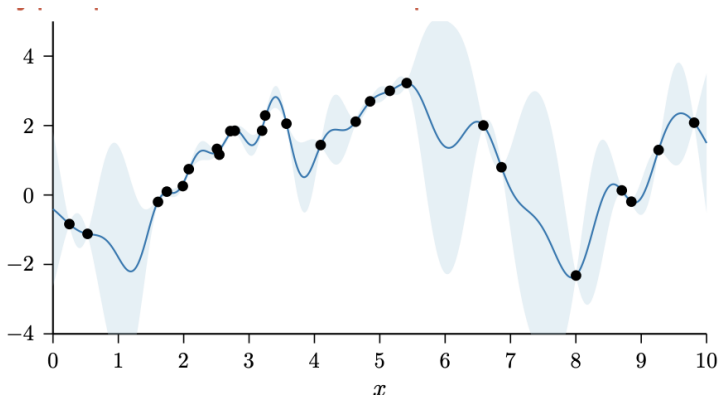
# Hyperparameters: Example



$$\theta = (\lambda, \ell, \sigma) = (1, 1, \frac{1}{5}), \quad \log p(\mathbf{y} \mid \mathbf{X}, \theta) = -27.6$$



# Hyperparameters: Example



$$\theta = (\lambda, \ell, \sigma) = (2, \frac{1}{3}, \frac{1}{20}), \quad \log p(\mathbf{y} \mid \mathbf{X}, \theta) = -46.5$$

# Hyperparameters are important

Comparing the marginal likelihoods, we see that the observed data are *over 100 million times more likely* to have been generated by the first model rather than from the second model! Clearly hyperparameters can be *quite important*.