

Gaussian Processes

Marc Deisenroth
Centre for Artificial Intelligence
Department of Computer Science
University College London

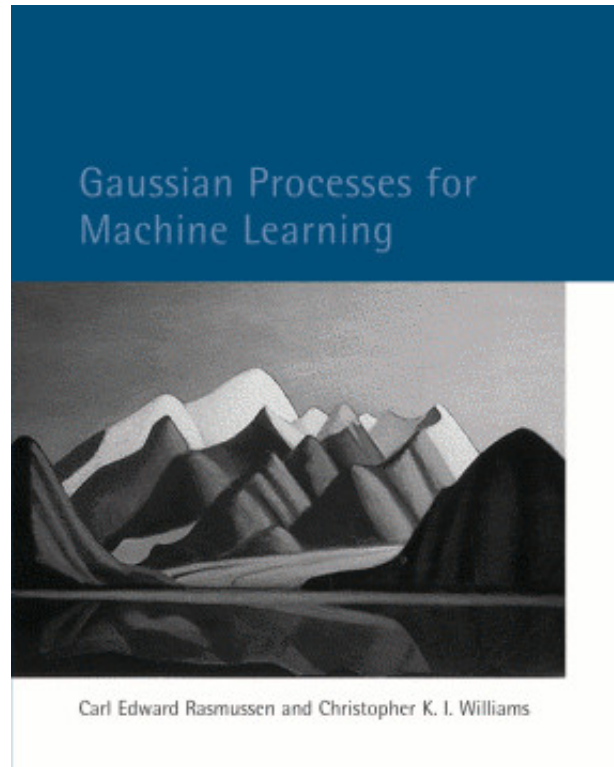
 @mpd37

`m.deisenroth@ucl.ac.uk`

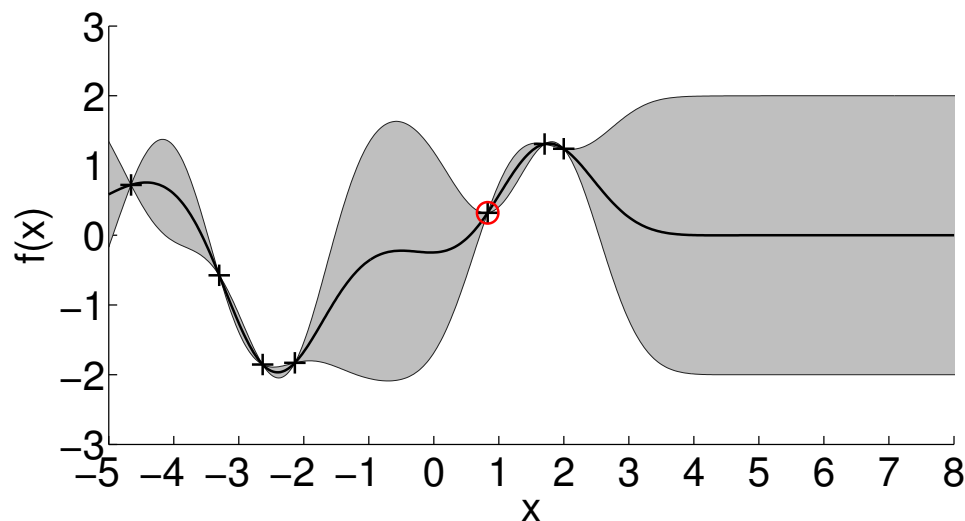
`https://deisenroth.cc`

AIMS Rwanda and AIMS Ghana

March/April 2020



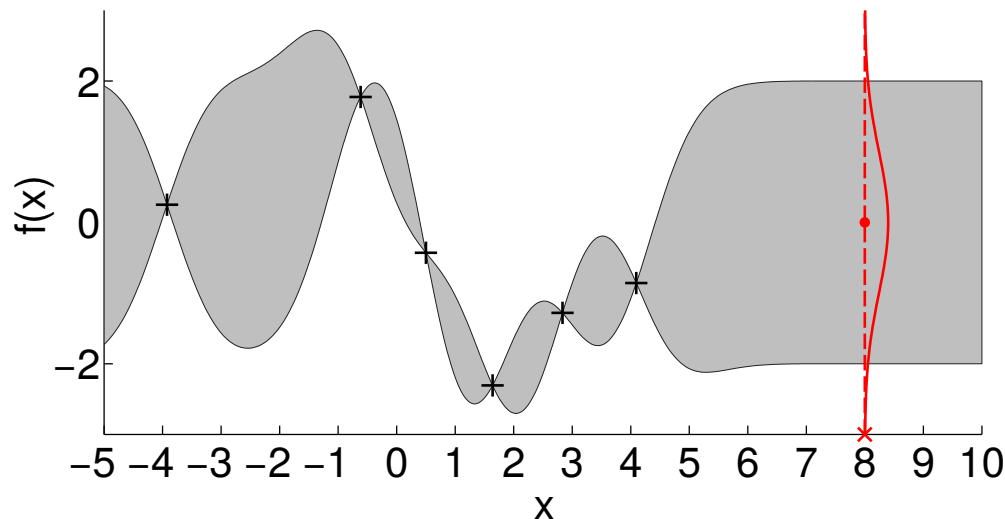
<http://www.gaussianprocess.org/>



Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$, find a **distribution over functions** $p(f)$ that explains the data

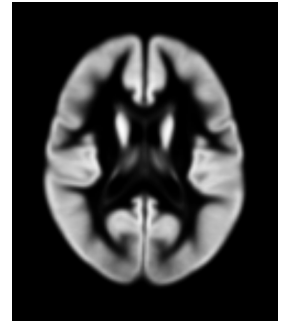
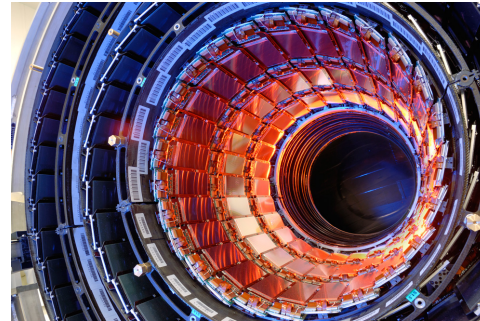
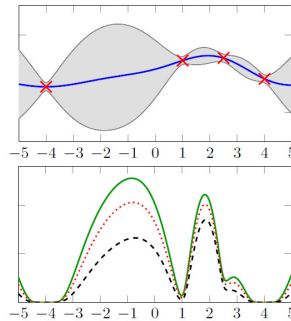
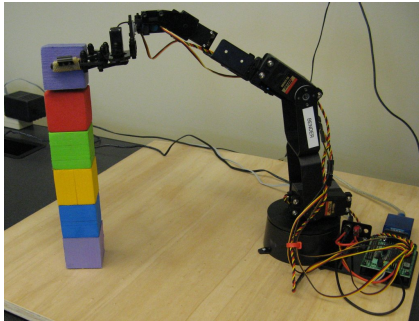
▶▶ Probabilistic regression problem



Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$, find a **distribution over functions** $p(f)$ that explains the data

▶▶ Probabilistic regression problem

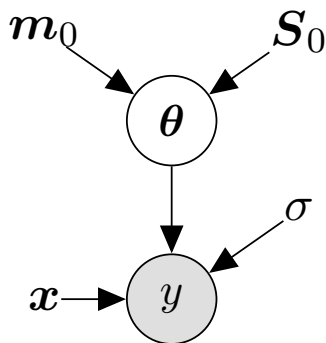


- Reinforcement learning and robotics
- Bayesian optimization (experimental design)
- Geostatistics
- Sensor networks
- Time-series modeling and forecasting
- High-energy physics
- Medical applications

Prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$

Likelihood $p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta}, \sigma^2)$

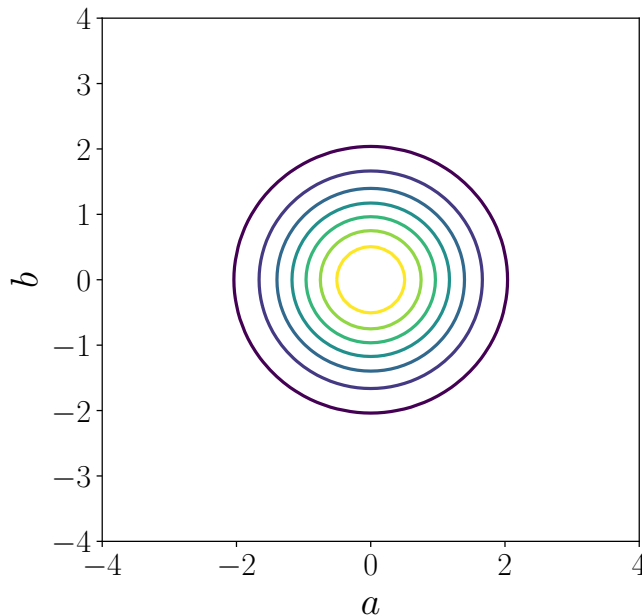
$$\implies y = \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$



- Parameter $\boldsymbol{\theta}$ becomes a latent (random) variable
- Distribution $p(\boldsymbol{\theta})$ induces a **distribution over plausible functions**
- Choose a conjugate Gaussian prior
 - Gaussian posterior $p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_N, \mathbf{S}_N)$
 - Closed-form computations (e.g., predictions, marginal likelihood)

Consider a linear regression setting

$$y = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

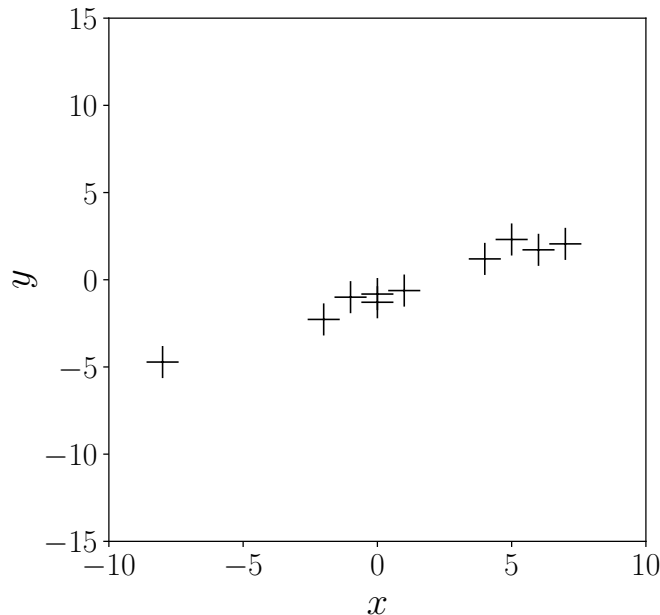
$$f_i(x) = a_i + b_i x, \quad [a_i, b_i] \sim p(a, b)$$

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{X} = [x_1, \dots, x_N], \quad \mathbf{y} = [y_1, \dots, y_N] \quad \text{Training data}$$

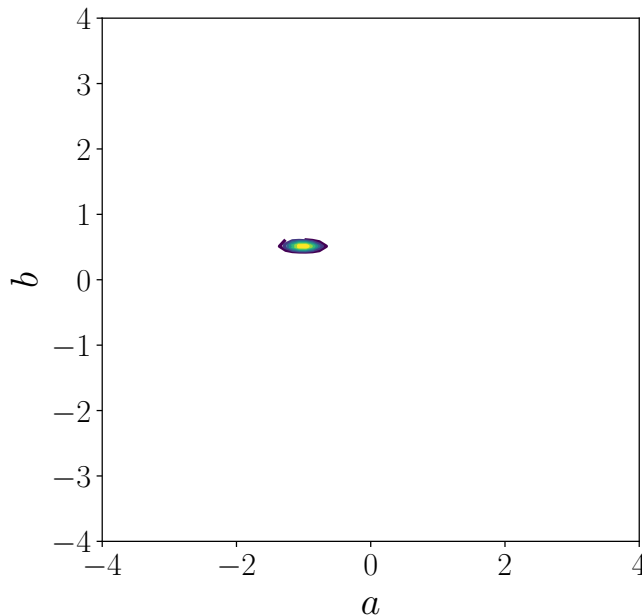


Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p(a, b | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N) \quad \text{Posterior}$$

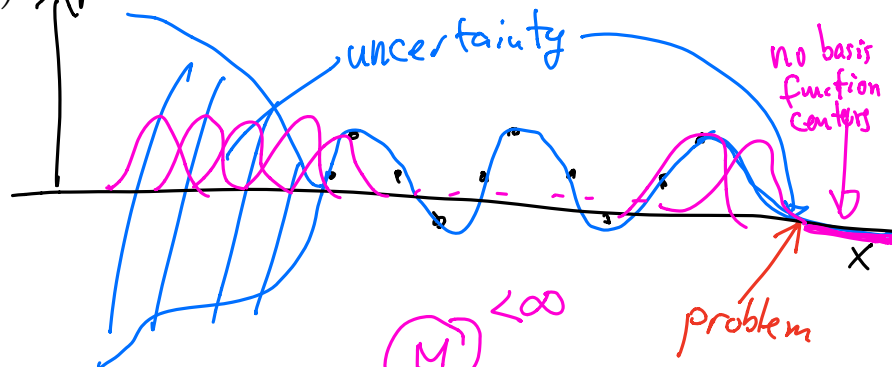


Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$[a_i, b_i] \sim p(a, b | \mathbf{X}, \mathbf{y}) \quad \uparrow$$

$$f_i = a_i + b_i x$$



$$f(x, \theta) = \sum_{m=1}^M \theta_m \phi_m(x) \quad < \infty$$

$$p(\theta) = \mathcal{N}(0, \mathbf{I})$$

- Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
 - ▶▶ Place a prior on functions
 - ▶▶ Make assumptions on the distribution of functions

- Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
 - ▶▶ Place a prior on functions
 - ▶▶ Make assumptions on the distribution of functions
- Intuition: function = infinitely long vector of function values
 - ▶▶ Make assumptions on the distribution of function values

- Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
 - ▶▶ Place a prior on functions
 - ▶▶ Make assumptions on the distribution of functions
- Intuition: function = infinitely long vector of function values
 - ▶▶ Make assumptions on the distribution of function values
 - ▶▶ **Gaussian process**

1 Gaussian Process: Definition

2 Regression as Inference

- GP Prior
- Likelihood
- Marginal Likelihood
- Posterior
- Predictions

3 Model Selection

- GP Training
- ~~Training~~ kernel etc.

4 Limitations and Guidelines

5 Application Areas

BLR:

$$p(y_* | x_*) = \int p(y_* | x_*, \theta) p(\theta) d\theta$$

parameters

consider all plausible (∞) values/settings of θ

GP:

$$p(y_* | x_*) = \int p(y_* | x_*, f) p(f) df$$

function

consider all plausible (∞) values/settings of f

$$f(\cdot) = \sin(\cdot)$$
$$f: \mathcal{X} \rightarrow \mathcal{Y}; f: \mathbb{R}^D \rightarrow \mathbb{R}$$

Gaussian Process: Definition

- We will place a distribution $p(f)$ on functions f
- Informally, a function can be considered an infinitely long vector of function values $f = [f_1, f_2, f_3, \dots]$

- We will place a distribution $p(f)$ on functions f
- Informally, a function can be considered an infinitely long vector of function values $f = [f_1, f_2, f_3, \dots]$
- A Gaussian process is a generalization of a multivariate Gaussian distribution to infinitely many variables.

- We will place a distribution $p(f)$ on functions f
- Informally, a function can be considered an infinitely long vector of function values $f = [f_1, f_2, f_3, \dots]$
- A Gaussian process is a generalization of a multivariate Gaussian distribution to infinitely many variables.

any finite subset is jointly Gaussian distributed

Definition (Rasmussen & Williams, 2006)

A **Gaussian process** (GP) is a collection of random variables f_1, f_2, \dots , any finite number of which is Gaussian distributed.

- training data is finite
- test data is finite
 - ↳ locations at which we want to evaluate $f(\cdot)$

```
def f(x):  
    return np.sin(x)  
xx = np.linspace(-10, 10, 50)  
f(xx)
```

test points

- We will place a distribution $p(f)$ on functions f
- Informally, a function can be considered an infinitely long vector of function values $f = [f_1, f_2, f_3, \dots]$
- A Gaussian process is a generalization of a multivariate Gaussian distribution to infinitely many variables.

Definition (Rasmussen & Williams, 2006)

A **Gaussian process** (GP) is a collection of random variables f_1, f_2, \dots , any finite number of which is Gaussian distributed.

- A Gaussian distribution is specified by a mean vector μ and a covariance matrix Σ
- A Gaussian process is specified by a **mean function** $m(\cdot)$ and a **covariance function (kernel)** $k(\cdot, \cdot)$ ► More on this later

Regression as Inference

Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$, find a (posterior) **distribution over functions** $p(f(\cdot) | \mathbf{X}, \mathbf{y})$ that explains the data. Here: \mathbf{X} training inputs, \mathbf{y} training targets

Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$, find a (posterior) **distribution over functions** $p(f(\cdot)|\mathbf{X}, \mathbf{y})$ that explains the data. Here: \mathbf{X} training inputs, \mathbf{y} training targets

Training data: \mathbf{X}, \mathbf{y} . Bayes' theorem yields

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$, find a (posterior) **distribution over functions** $p(f(\cdot) | \mathbf{X}, \mathbf{y})$ that explains the data. Here: \mathbf{X} training inputs, \mathbf{y} training targets

Training data: \mathbf{X}, \mathbf{y} . Bayes' theorem yields

$$p(f(\cdot) | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y} | \mathbf{X})}$$

Prior: $p(f(\cdot)) = GP(m, k)$ \blacktriangleright Specify mean m function and kernel k .

Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$, find a (posterior) **distribution over functions** $p(f(\cdot)|\mathbf{X}, \mathbf{y})$ that explains the data. Here: \mathbf{X} training inputs, \mathbf{y} training targets

Training data: \mathbf{X}, \mathbf{y} . Bayes' theorem yields

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Prior: $p(f(\cdot)) = GP(m, k)$ ► Specify mean m function and kernel k .

Likelihood (noise model): $p(\mathbf{y}|f(\cdot), \mathbf{X}) = \mathcal{N}(f(\mathbf{X}), \sigma_n^2 \mathbf{I})$

Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$, find a (posterior) **distribution over functions** $p(f(\cdot)|\mathbf{X}, \mathbf{y})$ that explains the data. Here: \mathbf{X} training inputs, \mathbf{y} training targets

Training data: \mathbf{X}, \mathbf{y} . Bayes' theorem yields

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Prior: $p(f(\cdot)) = GP(m, k)$ \blacktriangleright Specify mean m function and kernel k .

Likelihood (noise model): $p(\mathbf{y}|f(\cdot), \mathbf{X}) = \mathcal{N}(f(\mathbf{X}), \sigma_n^2 \mathbf{I})$

Marginal likelihood (evidence): $p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot)|\mathbf{X}) df$

Objective

For a set of observations $y_i = f(\mathbf{x}_i) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$, find a (posterior) **distribution over functions** $p(f(\cdot)|\mathbf{X}, \mathbf{y})$ that explains the data. Here: \mathbf{X} training inputs, \mathbf{y} training targets

Training data: \mathbf{X}, \mathbf{y} . Bayes' theorem yields

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Prior: $p(f(\cdot)) = GP(m, k)$ \blacktriangleright Specify mean m function and kernel k .

Likelihood (noise model): $p(\mathbf{y}|f(\cdot), \mathbf{X}) = \mathcal{N}(f(\mathbf{X}), \sigma_n^2 \mathbf{I})$

Marginal likelihood (evidence): $p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot)|\mathbf{X}) df$

Posterior: $p(f(\cdot)|\mathbf{y}, \mathbf{X}) = GP(m_{\text{post}}, k_{\text{post}})$

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Bayesian linear regression:

- Prior $p(\boldsymbol{\theta})$ on the parameters $\boldsymbol{\theta}$ allows us to encode some properties of the parameters (e.g., range, reasonable values, ...)
- Every sample $\boldsymbol{\theta}_i \sim p(\boldsymbol{\theta})$ induces a function $f_i(\cdot) := \boldsymbol{\theta}_i^\top \boldsymbol{\phi}(\cdot)$

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Bayesian linear regression:

- Prior $p(\boldsymbol{\theta})$ on the parameters $\boldsymbol{\theta}$ allows us to encode some properties of the parameters (e.g., range, reasonable values, ...)
- Every sample $\boldsymbol{\theta}_i \sim p(\boldsymbol{\theta})$ induces a function $f_i(\cdot) := \boldsymbol{\theta}_i^\top \boldsymbol{\phi}(\cdot)$

Gaussian process:

- GP prior: $p(f(\cdot))$
- Function plays the role of the parameters
 - ▶▶ Every sample $f_i(\cdot) \sim GP$ is a function

- continuity
- differentiability
- function is positive
- function varies slowly
(no rapid change in curve)

- strictly monotonic (invertible)
 - symmetry
 - periodicity
 - bounded function values
-

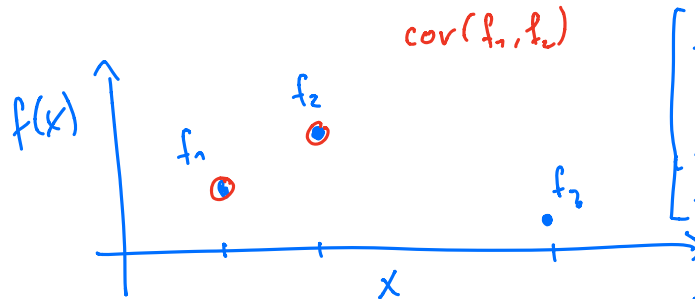
- Bayesian prior specifies assumptions on the quantity of interest
- What assumptions could we make on the underlying function?
- What characterizes the function we want to model?

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \in \mathbb{R}^2$$

$$f_1 \sim \text{GP}: \mathbb{R}^D \rightarrow \mathbb{R}$$

$$f_2 \sim \text{GP}: \mathbb{R}^D \rightarrow \mathbb{R}$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = W \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$$



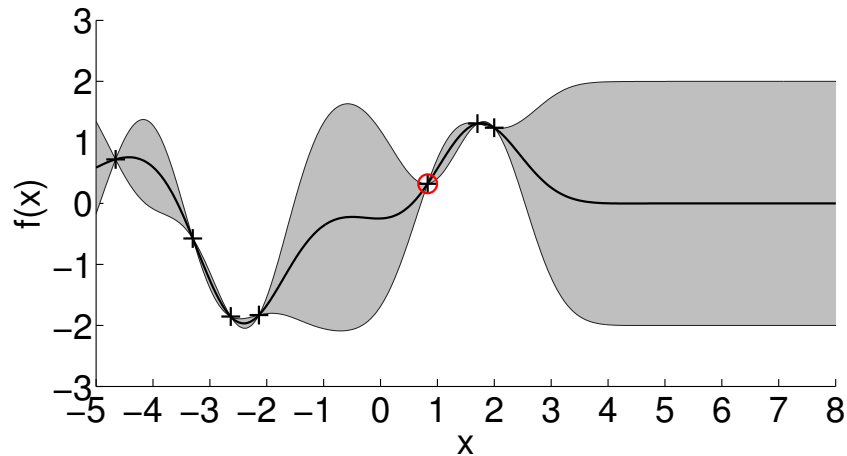
$$\begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \end{bmatrix} \sim \mathcal{N}(m, S)$$

$\in \mathbb{R}^{4 \times 4}$

→ multi-output Gaussian process

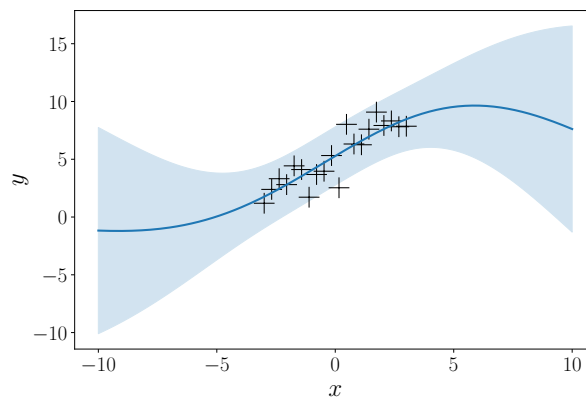
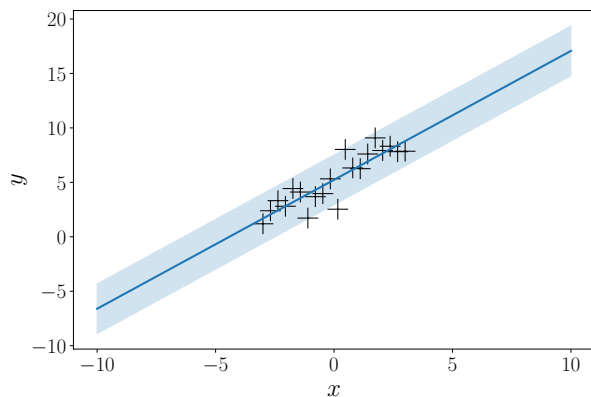
$$S = \begin{bmatrix} \text{var}[f_1] & \text{cov}[f_1, f_2] & \dots \\ \vdots & \text{var}[f_2] & \dots \\ & & \dots & \text{var}[f_4] \end{bmatrix}$$

- Bayesian prior specifies assumptions on the quantity of interest
- What assumptions could we make on the underlying function?
- What characterizes the function we want to model?
 - Mean function
 - Covariance function

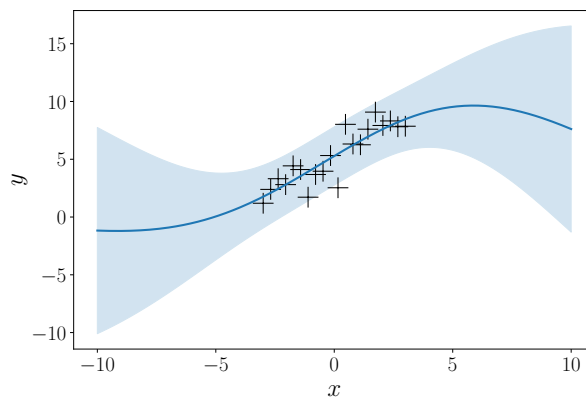
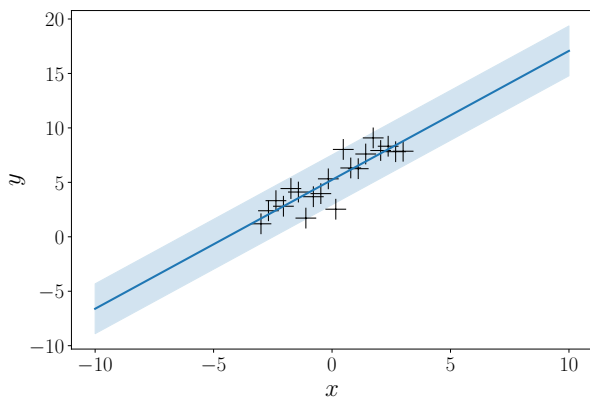


$$m(\mathbf{x}) = \mathbb{E}_f[f(\mathbf{x})], \quad f \sim GP$$

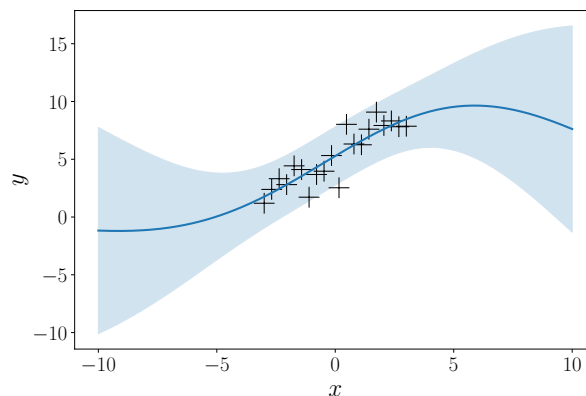
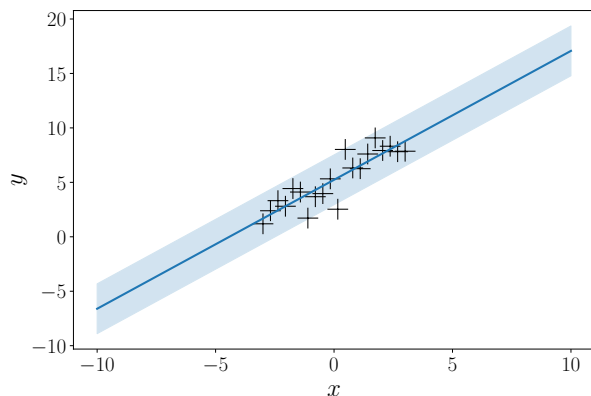
- The **average function** of the distribution over functions
- Allows us to **bias the model** (can make sense in application-specific settings)



- Can be a parametrized function, e.g., linear, exponential, or neural network. Example: $m_{\theta}(\mathbf{x}) = \boldsymbol{\theta}^{\top} \boldsymbol{\phi}(\mathbf{x})$



- Can be a parametrized function, e.g., linear, exponential, or neural network. Example: $m_{\theta}(\mathbf{x}) = \theta^{\top} \phi(\mathbf{x})$
- Prior mean function m_{θ} can incorporate **problem-specific prior knowledge** (e.g., in robotics, natural sciences)
- Can simplify the learning problem



- Can be a parametrized function, e.g., linear, exponential, or neural network. Example: $m_{\theta}(\mathbf{x}) = \theta^{\top} \phi(\mathbf{x})$
- Prior mean function m_{θ} can incorporate **problem-specific prior knowledge** (e.g., in robotics, natural sciences)
- Can simplify the learning problem
- Often: “Agnostic” mean function in the absence of data or prior knowledge: $m(\cdot) \equiv 0$ everywhere (for symmetry reasons)

- Covariance function (kernel) is symmetric and positive semi-definite

- Covariance function (kernel) is symmetric and positive semi-definite
- Compute covariances/correlations between (unknown) function values by just looking at the corresponding inputs:

$$\text{Cov}[f(\mathbf{x}_i), f(\mathbf{x}_j)] = k(\mathbf{x}_i, \mathbf{x}_j)$$

▶▶ Kernel trick (Schölkopf & Smola, 2002)

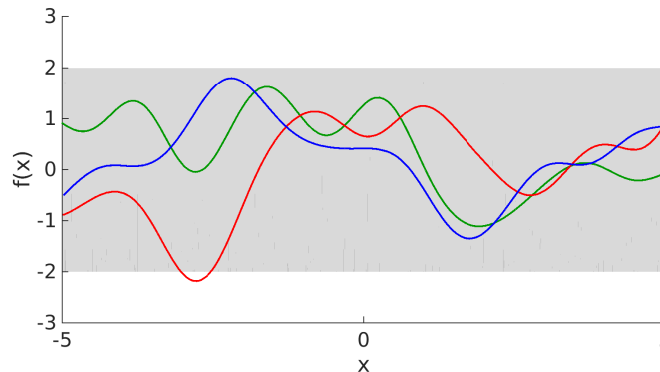
- Covariance function (kernel) is symmetric and positive semi-definite
- Compute covariances/correlations between (unknown) function values by just looking at the corresponding inputs:

$$\text{Cov}[f(\mathbf{x}_i), f(\mathbf{x}_j)] = k(\mathbf{x}_i, \mathbf{x}_j)$$

- ▶▶ Kernel trick (Schölkopf & Smola, 2002)
- Encodes high-level structural assumptions (e.g., smoothness, periodicity) of the function we want to model

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$

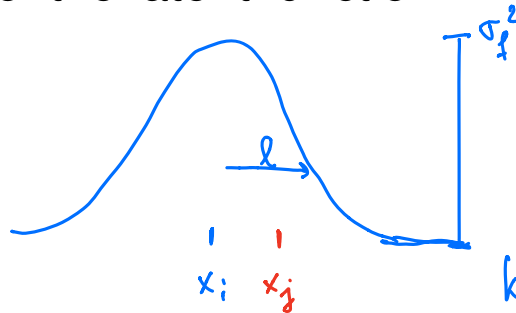
- Assumption on latent function: **Smooth** (∞ differentiable)



$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\ell^2}\right)$$

Handwritten notes: $\in [0, 1]$ (bracketed over the exponential term), $\|\frac{\mathbf{x}_i - \mathbf{x}_j}{\ell}\|^2$ (pointing to the exponent)

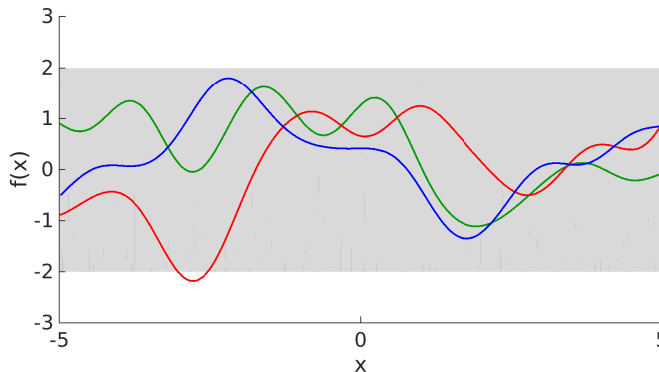
- Assumption on latent function: **Smooth** (∞ differentiable)
- σ_f : **Amplitude** of the latent function



$\ell \approx$ standard deviation in a Gaussian distribution

$k(x_i, x_j) = \text{cov}(f(x_i), f(x_j))$
 ≈ 0 if x_i and x_j are "far" from each other
 \rightarrow depends on ℓ

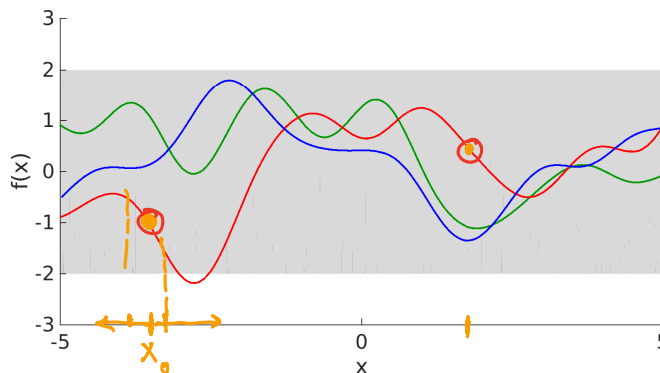
Handwritten note: x_j (pointing to the second x_j in the equation)



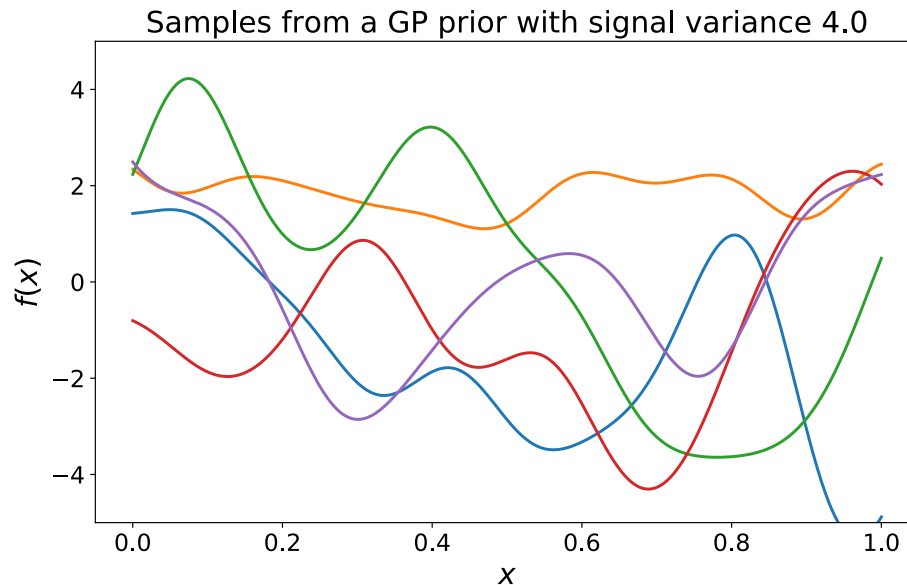
$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$

- Assumption on latent function: **Smooth** (∞ differentiable)
- σ_f : **Amplitude** of the latent function
- ℓ : **Length-scale**. How far do we have to move in input space before the function value changes significantly, i.e., when do function values become uncorrelated?

►► Smoothness parameter

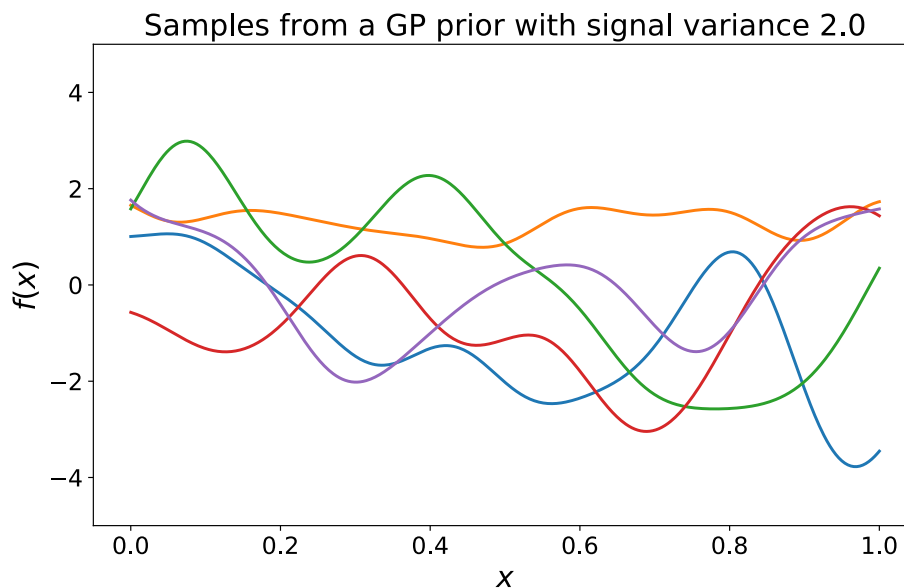


$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



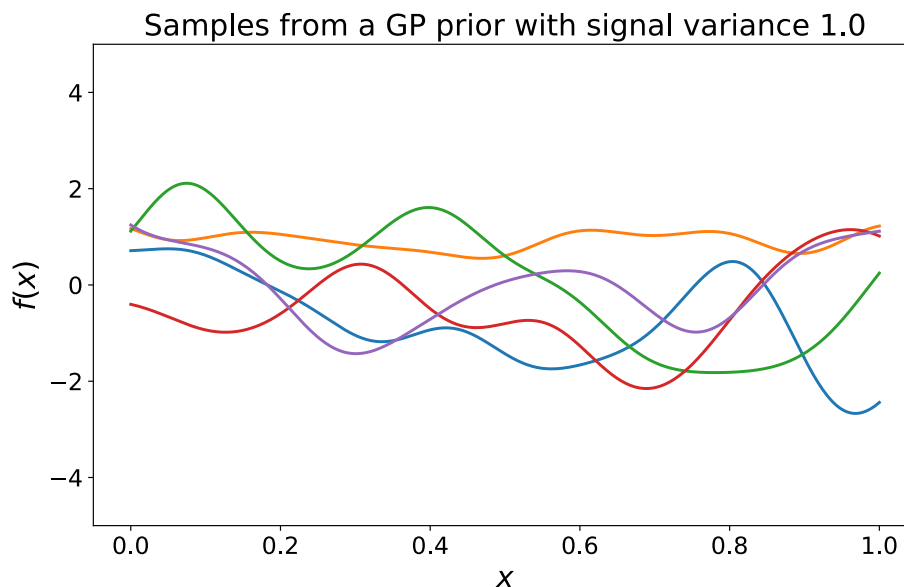
- Controls the amplitude (vertical magnitude) of the function we wish to model

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



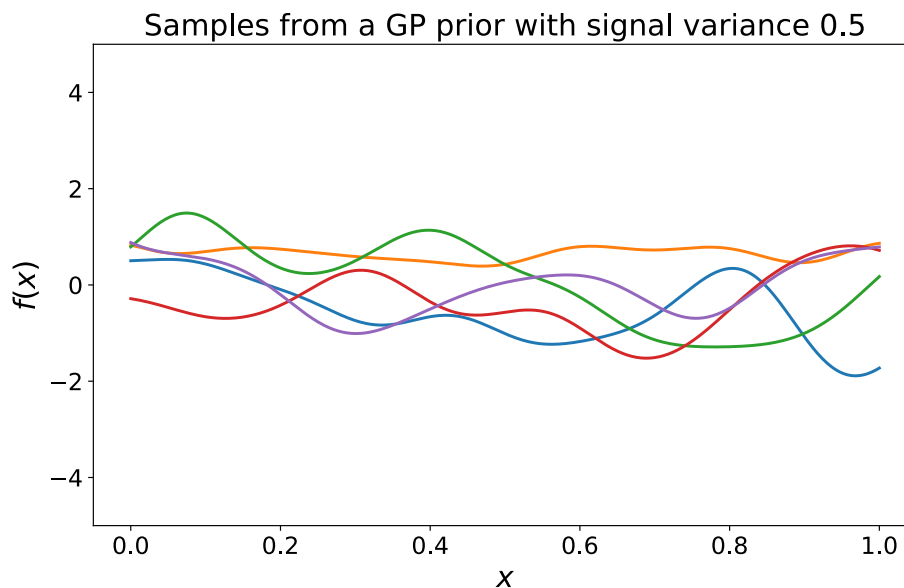
- Controls the amplitude (vertical magnitude) of the function we wish to model

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



- Controls the amplitude (vertical magnitude) of the function we wish to model

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$

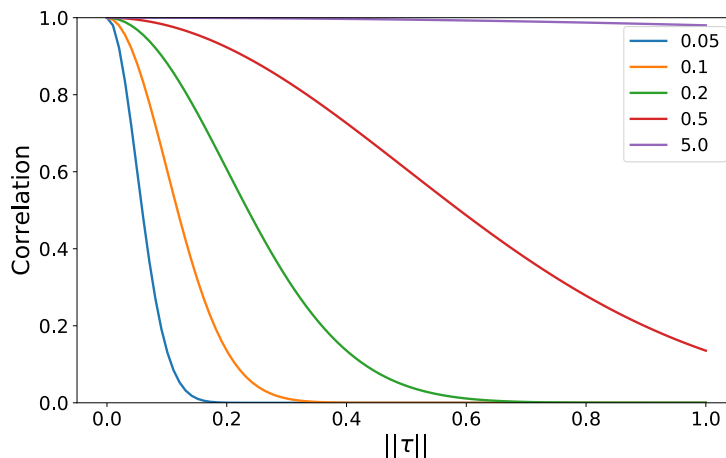


- Controls the amplitude (vertical magnitude) of the function we wish to model

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$

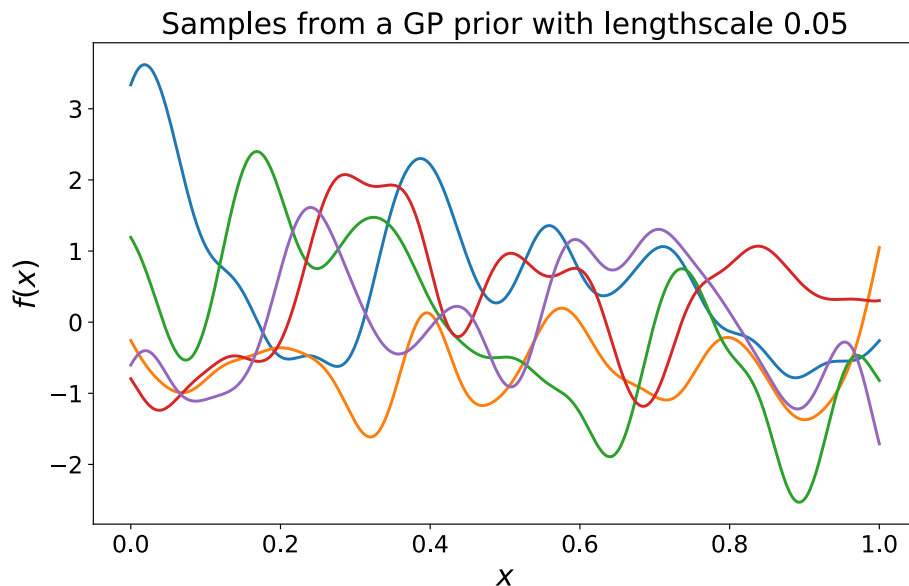
- How “wiggly” is the function?
- How much information we can transfer to other function values?
 - ▶▶ Correlation between function values
- How far do we have to move in input space from \mathbf{x} to \mathbf{x}' to make $f(\mathbf{x})$ and $f(\mathbf{x}')$ uncorrelated?

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



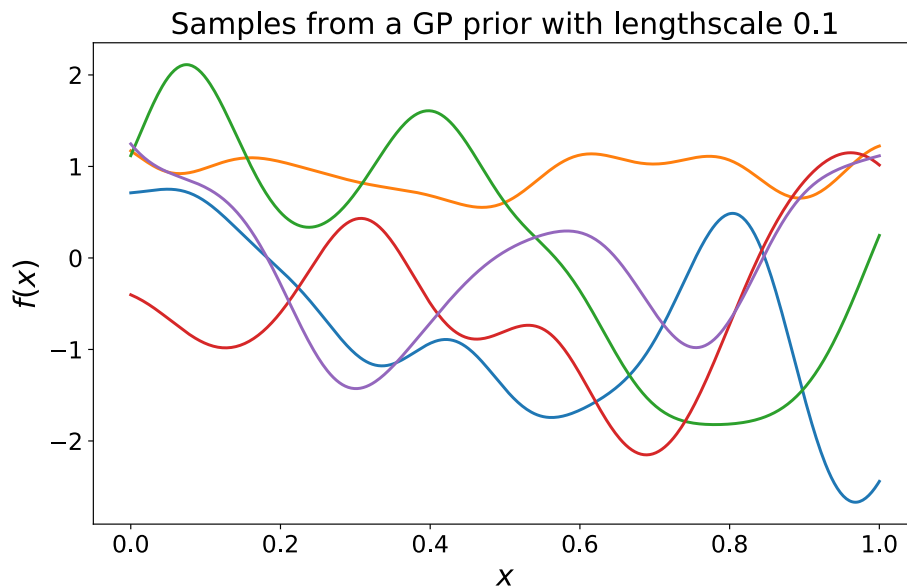
- Correlation between function values $f(\mathbf{x})$ and $f(\mathbf{x}')$ depends on the (scaled) distance $\|\boldsymbol{\tau}\|/\ell = \|\mathbf{x} - \mathbf{x}'\|/\ell$ of the corresponding inputs.
- What does a short/long length-scale ℓ imply?

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



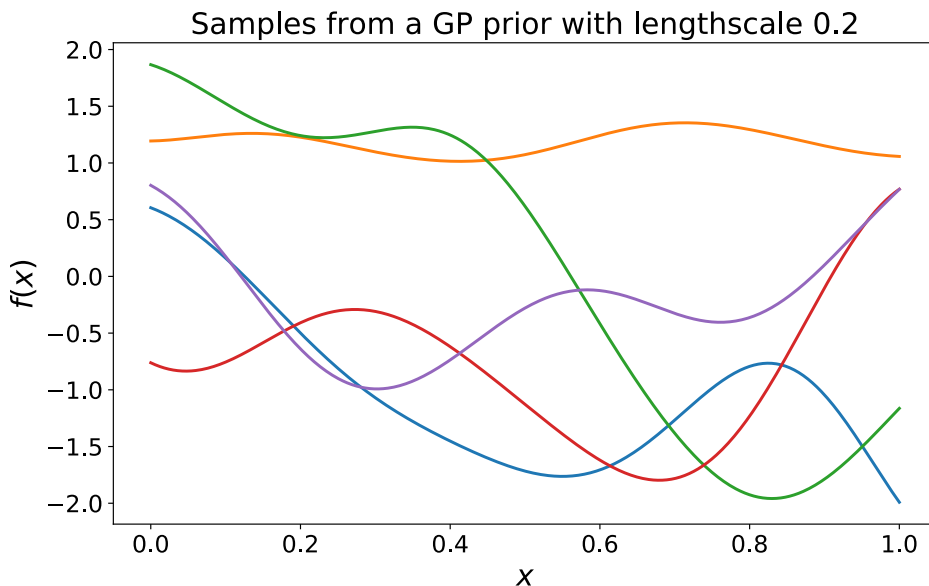
►► Explore interactive diagrams at
<https://drafts.distill.pub/gp/>

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



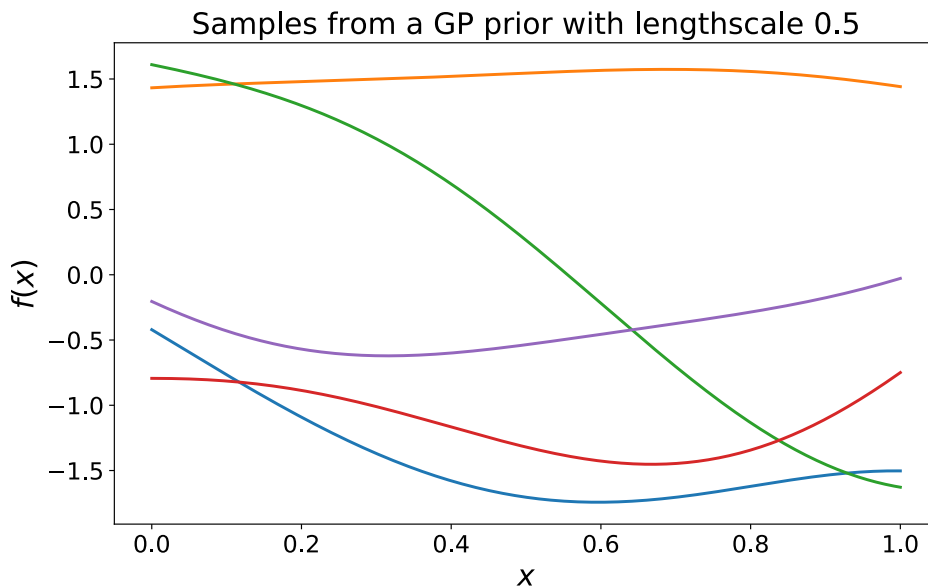
►► Explore interactive diagrams at
<https://drafts.distill.pub/gp/>

$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



►► Explore interactive diagrams at
<https://drafts.distill.pub/gp/>

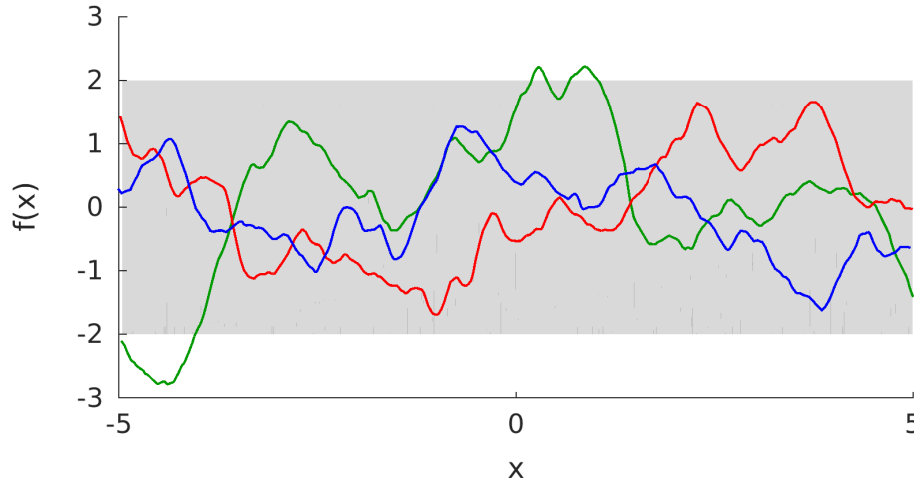
$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \ell^2\right)$$



►► Explore interactive diagrams at
<https://drafts.distill.pub/gp/>

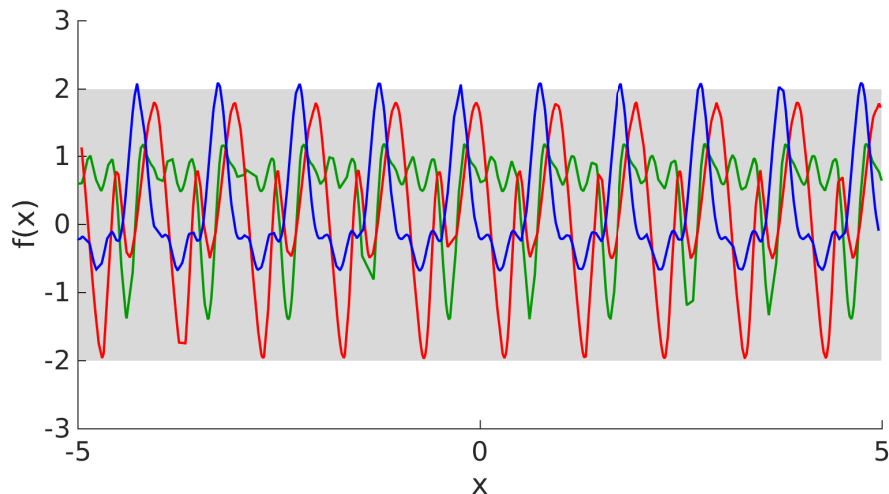
$$k_{Mat,3/2}(x_i, x_j) = \sigma_f^2 \left(1 + \frac{\sqrt{3}\|x_i - x_j\|}{\ell} \right) \exp \left(- \frac{\sqrt{3}\|x_i - x_j\|}{\ell} \right)$$

- Assumption on latent function: **1-times differentiable**
- σ_f : **Amplitude** of the latent function
- ℓ : **Length-scale**. How far do we have to move in input space before the function value changes significantly?



$$k_{per}(x_i, x_j) = \sigma_f^2 \exp\left(-\frac{2 \sin^2\left(\frac{\kappa(x_i - x_j)}{2\pi}\right)}{\ell^2}\right)$$
$$= k_{Gauss}(\mathbf{u}(x_i), \mathbf{u}(x_j)), \quad \mathbf{u}(x) = \begin{bmatrix} \cos(\kappa x) \\ \sin(\kappa x) \end{bmatrix}$$

- Assumption on latent function: **periodic**
- **Periodicity parameter** κ



Assume k_1 and k_2 are valid covariance functions and $u(\cdot)$ is a (nonlinear) transformation of the input space. Then

- $k_1 + k_2$ is a valid covariance function

Assume k_1 and k_2 are valid covariance functions and $u(\cdot)$ is a (nonlinear) transformation of the input space. Then

- $k_1 + k_2$ is a valid covariance function
- $k_1 k_2$ is a valid covariance function

Assume k_1 and k_2 are valid covariance functions and $u(\cdot)$ is a (nonlinear) transformation of the input space. Then

- $k_1 + k_2$ is a valid covariance function
- $k_1 k_2$ is a valid covariance function
- $k(u(\mathbf{x}), u(\mathbf{x}'))$ is a valid covariance function (MacKay, 1998)
 - ▶▶ Periodic covariance function
 - ▶▶ [Manifold Gaussian process](#) (Calandra et al., 2016)
 - ▶▶ [Deep kernel learning](#) (Wilson et al., 2016)

Assume k_1 and k_2 are valid covariance functions and $u(\cdot)$ is a (nonlinear) transformation of the input space. Then

- $k_1 + k_2$ is a valid covariance function
- $k_1 k_2$ is a valid covariance function
- $k(u(\mathbf{x}), u(\mathbf{x}'))$ is a valid covariance function (MacKay, 1998)
 - ▶▶ Periodic covariance function
 - ▶▶ **Manifold Gaussian process** (Calandra et al., 2016)
 - ▶▶ **Deep kernel learning** (Wilson et al., 2016)
- ▶▶ **Automatic Statistician** (Lloyd et al., 2014)

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Gaussian likelihood in linear regression:

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \boldsymbol{\theta}^\top \mathbf{x}, \sigma^2)$$

- **Function** (not a distribution) **of the parameters**
- Describes how parameters and observed data are connected
- Tells us **how to transform parameters into (noisy) data**

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Gaussian likelihood in linear regression:

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \boldsymbol{\theta}^\top \mathbf{x}, \sigma^2)$$

- **Function** (not a distribution) **of the parameters**
- Describes how parameters and observed data are connected
- Tells us **how to transform parameters into (noisy) data**

Gaussian likelihood in Gaussian processes:

$$p(y|\mathbf{x}, f(\cdot)) = \mathcal{N}(y | f(\mathbf{x}), \sigma^2)$$

- Parameters are the function f itself

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Bayesian linear regression with a Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$:

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

- Normalizes the posterior distribution

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Bayesian linear regression with a Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$:

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{0}, \boldsymbol{\Phi}\boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I}) \end{aligned}$$

- Normalizes the posterior distribution
- Can be computed analytically

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Bayesian linear regression with a Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$:

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{0}, \boldsymbol{\Phi}\boldsymbol{\Phi}^\top + \sigma^2\mathbf{I}) \\ &= \mathbb{E}_{\boldsymbol{\theta}}[p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})] \end{aligned}$$

- Normalizes the posterior distribution
- Can be computed analytically
- Expected likelihood (under the parameter prior)
- Expected predictive distribution of the training targets \mathbf{y} (under the parameter prior)

Gaussian process marginal likelihood

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

- Normalizes the posterior distribution

Gaussian process marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(f(\cdot))d\mathbf{f} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I}) \end{aligned}$$

- Normalizes the posterior distribution
- Can be computed analytically

Gaussian process marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(f(\cdot))d\boldsymbol{\theta} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I}) \\ &= \mathbb{E}_f[p(\mathbf{y}|\mathbf{X}, f(\cdot))] \end{aligned}$$

- Normalizes the posterior distribution
- Can be computed analytically
- Expected likelihood (under the GP prior)
- Expected predictive distribution of the training targets \mathbf{y} (under the GP prior)

Gaussian process marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(f(\cdot))d\boldsymbol{\theta} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I}) \\ &= \mathbb{E}_f[p(\mathbf{y}|\mathbf{X}, f(\cdot))] \end{aligned}$$

- Normalizes the posterior distribution
- Can be computed analytically
- Expected likelihood (under the GP prior)
- Expected predictive distribution of the training targets \mathbf{y} (under the GP prior)

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

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, N$$

Posterior over functions (with training data \mathbf{X}, \mathbf{y}):

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y}|\mathbf{X})}$$

Posterior over functions (with training data \mathbf{X}, \mathbf{y}):

$$p(f(\cdot) | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y} | \mathbf{X})}$$

Using the properties of Gaussians, we obtain (with $\mathbf{K} := k(\mathbf{X}, \mathbf{X})$)

$$p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot)) = \mathcal{N}(\mathbf{y} | f(\mathbf{X}), \sigma_n^2 \mathbf{I}) GP(m(\cdot), k(\cdot, \cdot))$$

Posterior over functions (with training data \mathbf{X}, \mathbf{y}):

$$p(f(\cdot) | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y} | \mathbf{X})}$$

Using the properties of Gaussians, we obtain (with $\mathbf{K} := k(\mathbf{X}, \mathbf{X})$)

$$p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot)) = \mathcal{N}(\mathbf{y} | f(\mathbf{X}), \sigma_n^2 \mathbf{I}) GP(m(\cdot), k(\cdot, \cdot))$$

$$= Z \times GP(m_{\text{post}}(\cdot), k_{\text{post}}(\cdot, \cdot))$$

$$m_{\text{post}}(\cdot) = m(\cdot) + k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))$$

$$k_{\text{post}}(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \cdot)$$

Posterior over functions (with training data \mathbf{X}, \mathbf{y}):

$$p(f(\cdot) | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot))}{p(\mathbf{y} | \mathbf{X})}$$

Using the properties of Gaussians, we obtain (with $\mathbf{K} := k(\mathbf{X}, \mathbf{X})$)

$$p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot)) = \mathcal{N}(\mathbf{y} | f(\mathbf{X}), \sigma_n^2 \mathbf{I}) GP(m(\cdot), k(\cdot, \cdot))$$

$$= Z \times GP(m_{\text{post}}(\cdot), k_{\text{post}}(\cdot, \cdot))$$

$$m_{\text{post}}(\cdot) = m(\cdot) + k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))$$

$$k_{\text{post}}(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \cdot)$$

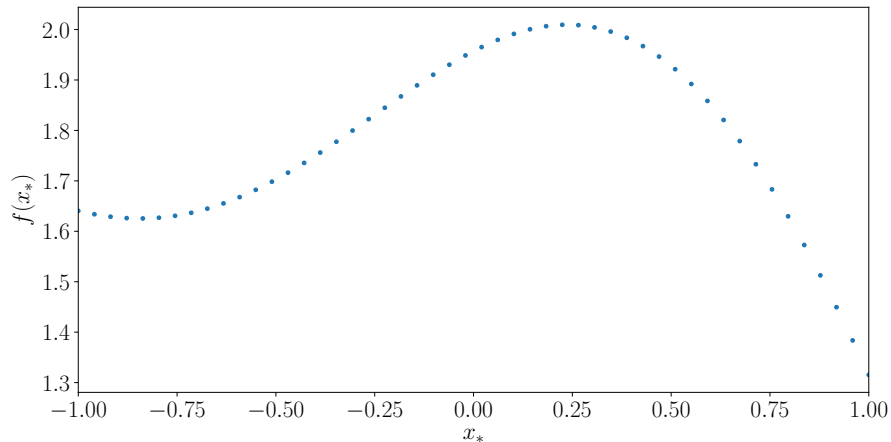
Marginal likelihood:

$$Z = p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | f(\cdot), \mathbf{X}) p(f(\cdot)) df = \mathcal{N}(\mathbf{y} | m(\mathbf{X}), \mathbf{K} + \sigma_n^2 \mathbf{I})$$

- GP is a distribution over functions
 - ▶▶ A sample from a GP will be an entire function

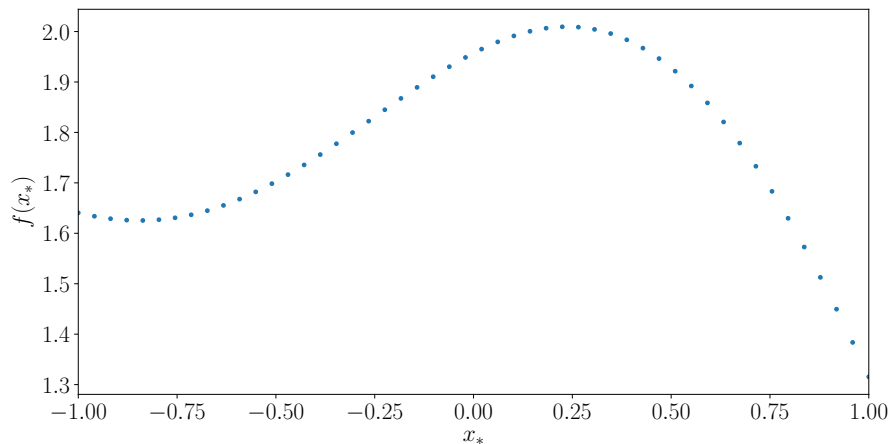
- GP is a distribution over functions
 - ▶▶ A sample from a GP will be an entire function
- In practice, we cannot sample functions directly

- GP is a distribution over functions
 - ▶▶ A sample from a GP will be an entire function
- In practice, we cannot sample functions directly
- Instead: function = collection of function values



- GP is a distribution over functions
 - ▶▶ A sample from a GP will be an entire function
- In practice, we cannot sample functions directly
- Instead: function = collection of function values
- Determine function values at a finite set of input locations

$$\mathbf{X}_* = [\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_*^{(K)}]$$



- Without any training data, the predictive distribution at test points \mathbf{X}_* is

$$\begin{aligned} p(\mathbf{f}(\mathbf{X}_*)|\mathbf{X}_*) &= \mathcal{N}(\mathbb{E}_f[f(\mathbf{X}_*)], \mathbb{V}_f[f(\mathbf{X}_*)]) \\ &= \mathcal{N}(m_{\text{prior}}(\mathbf{X}_*), k_{\text{prior}}(\mathbf{X}_*, \mathbf{X}_*)) \end{aligned}$$

- Without any training data, the predictive distribution at test points \mathbf{X}_* is

$$\begin{aligned} p(\mathbf{f}(\mathbf{X}_*)|\mathbf{X}_*) &= \mathcal{N}(\mathbb{E}_f[f(\mathbf{X}_*)], \mathbb{V}_f[f(\mathbf{X}_*)]) \\ &= \mathcal{N}(m_{\text{prior}}(\mathbf{X}_*), k_{\text{prior}}(\mathbf{X}_*, \mathbf{X}_*)) \end{aligned}$$

- Exploited: Definition of GP that **all function values are jointly Gaussian distributed**

- Without any training data, the predictive distribution at test points \mathbf{X}_* is

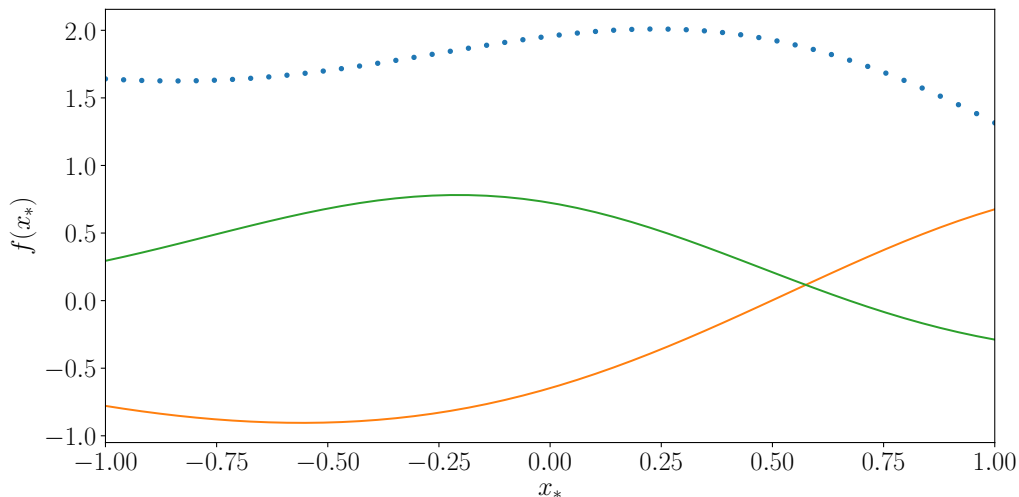
$$\begin{aligned} p(\mathbf{f}(\mathbf{X}_*)|\mathbf{X}_*) &= \mathcal{N}(\mathbb{E}_f[f(\mathbf{X}_*)], \mathbb{V}_f[f(\mathbf{X}_*)]) \\ &= \mathcal{N}(m_{\text{prior}}(\mathbf{X}_*), k_{\text{prior}}(\mathbf{X}_*, \mathbf{X}_*)) \end{aligned}$$

- Exploited: Definition of GP that **all function values are jointly Gaussian distributed**
- Generate “function draws” (samples from the GP prior)

$$f_k(\mathbf{X}_*) \sim \mathcal{N}(m_{\text{prior}}(\mathbf{X}_*), k_{\text{prior}}(\mathbf{X}_*, \mathbf{X}_*))$$

- Goal: Generate random functions f_k , so that

$$f_k(\mathbf{X}_*) \sim \mathcal{N}(m_{\text{prior}}(\mathbf{X}_*), k_{\text{prior}}(\mathbf{X}_*, \mathbf{X}_*))$$



- Goal: Generate random functions f_k , so that

$$f_k(\mathbf{X}_*) \sim \mathcal{N}(m_{\text{prior}}(\mathbf{X}_*), k_{\text{prior}}(\mathbf{X}_*, \mathbf{X}_*))$$

- Define $\mathbf{m}_* := m_{\text{prior}}(\mathbf{X}_*)$ and $\mathbf{K}_{**} := k_{\text{prior}}(\mathbf{X}_*, \mathbf{X}_*)$. Then

$$f_k(\mathbf{X}_*) \sim \mathcal{N}(\mathbf{m}_*, \mathbf{K}_{**})$$

►► Sample from a multivariate Gaussian

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

- **Objective:** Find $p(f(\mathbf{X}_*)|\mathbf{X}, \mathbf{y}, \mathbf{X}_*)$ for training data \mathbf{X}, \mathbf{y} and test inputs \mathbf{X}_* .
- GP prior at training inputs: $p(f|\mathbf{X}) = \mathcal{N}(m(\mathbf{X}), \mathbf{K})$
- Gaussian Likelihood: $p(\mathbf{y}|f, \mathbf{X}) = \mathcal{N}(f(\mathbf{X}), \sigma_n^2 \mathbf{I})$

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

- **Objective:** Find $p(f(\mathbf{X}_*) | \mathbf{X}, \mathbf{y}, \mathbf{X}_*)$ for training data \mathbf{X}, \mathbf{y} and test inputs \mathbf{X}_* .

- GP prior at training inputs: $p(f | \mathbf{X}) = \mathcal{N}(m(\mathbf{X}), \mathbf{K})$ cov($f(x), f(x_*)$)

- Gaussian Likelihood: $p(\mathbf{y} | f, \mathbf{X}) = \mathcal{N}(f(\mathbf{X}), \sigma_n^2 \mathbf{I})$

- With $f \sim GP$ it follows that \mathbf{f}, \mathbf{f}_* are jointly Gaussian distributed:

$$p(\mathbf{f}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*) = \mathcal{N} \left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

f_1, f_2, \dots
 $p(f_1, f_2) = \mathcal{N}(m, S)$
 $\mathbf{f} := [f_1, \dots, f_N] = [f(x_1), \dots, f(x_N)] \in \mathbb{R}^N$
 $\mathbf{f}_* := [f_1, \dots, f_K] = [f(x_1), \dots, f(x_K)] \in \mathbb{R}^K$

cov($f(x_*), f(x)$)
 $p(\mathbf{f}, \mathbf{f}_*) = \mathcal{N}(\dots)$
 $\mathbb{E}_{f(x)} [f] = m(x)$
 $\mathbb{E}_{f(x_*)} [f_*] = m(x_*)$
var $f(x) = K$
var $f(x_*) = K(x_*, x_*)$

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

- **Objective:** Find $p(f(\mathbf{X}_*)|\mathbf{X}, \mathbf{y}, \mathbf{X}_*)$ for training data \mathbf{X}, \mathbf{y} and test inputs \mathbf{X}_* .
- GP prior at training inputs: $p(f|\mathbf{X}) = \mathcal{N}(m(\mathbf{X}), \mathbf{K})$
- Gaussian Likelihood: $p(\mathbf{y}|f, \mathbf{X}) = \mathcal{N}(f(\mathbf{X}), \sigma_n^2 \mathbf{I})$
- With $f \sim GP$ it follows that \mathbf{f}, \mathbf{f}_* are jointly Gaussian distributed:

$$p(\mathbf{f}, \mathbf{f}_*|\mathbf{X}, \mathbf{X}_*) = \mathcal{N} \left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

- Due to the Gaussian likelihood, we also get (\mathbf{f} is unobserved)

$$p(\mathbf{y}, \mathbf{f}_*|\mathbf{X}, \mathbf{X}_*) = \mathcal{N} \left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

Prior evaluated at \mathbf{X}, \mathbf{X}_* :

$$p(\mathbf{y}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*) = \mathcal{N} \left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

Posterior **predictive distribution** $p(\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*)$ at test inputs \mathbf{X}_*

Prior evaluated at \mathbf{X}, \mathbf{X}_* :

$$p(\mathbf{y}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*) = \mathcal{N} \left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

Posterior **predictive distribution** $p(\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*)$ at test inputs \mathbf{X}_* obtained by **Gaussian conditioning**:

$$p(\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*) = \mathcal{N} \left(\mathbb{E}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*], \mathbb{V}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] \right)$$

$$\mathbb{E}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] = \underbrace{m(\mathbf{X}_*)}_{\text{prior mean}} + \underbrace{k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}}_{\text{"Kalman gain"}} \underbrace{(\mathbf{y} - m(\mathbf{X}))}_{\text{error}}$$

Prior evaluated at \mathbf{X}, \mathbf{X}_* :

$$p(\mathbf{y}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*) = \mathcal{N} \left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

Posterior **predictive distribution** $p(\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*)$ at test inputs \mathbf{X}_* obtained by **Gaussian conditioning**:

$$p(\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*) = \mathcal{N}(\mathbb{E}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*], \mathbb{V}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*])$$

$$\mathbb{E}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] = \underbrace{m(\mathbf{X}_*)}_{\text{prior mean}} + \underbrace{k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}}_{\text{“Kalman gain”}} \underbrace{(\mathbf{y} - m(\mathbf{X}))}_{\text{error}}$$

$$\mathbb{V}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] = \underbrace{k(\mathbf{X}_*, \mathbf{X}_*)}_{\text{prior variance}} - \underbrace{k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \mathbf{X}_*)}_{\geq 0}$$

- GP posterior (from earlier):

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = GP(m_{\text{post}}(\cdot), k_{\text{post}}(\cdot, \cdot))$$

$$m_{\text{post}}(\cdot) = m(\cdot) + k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))$$

$$k_{\text{post}}(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \cdot)$$

- GP posterior (from earlier):

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = GP(m_{\text{post}}(\cdot), k_{\text{post}}(\cdot, \cdot))$$

$$m_{\text{post}}(\cdot) = m(\cdot) + k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))$$

$$k_{\text{post}}(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \cdot)$$

- GP posterior predictions at \mathbf{X}_* :

$$p(\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*) = \mathcal{N}(\mathbb{E}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*], \mathbb{V}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*])$$

$$\mathbb{E}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*] = m(\mathbf{X}_*) + k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))$$

$$\mathbb{V}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*] = k(\mathbf{X}_*, \mathbf{X}_*) - k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \mathbf{X}_*)$$

- GP posterior (from earlier):

$$p(f(\cdot)|\mathbf{X}, \mathbf{y}) = GP(m_{\text{post}}(\cdot), k_{\text{post}}(\cdot, \cdot))$$

$$m_{\text{post}}(\cdot) = m(\cdot) + k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))$$

$$k_{\text{post}}(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \cdot)$$

- GP posterior predictions at \mathbf{X}_* :

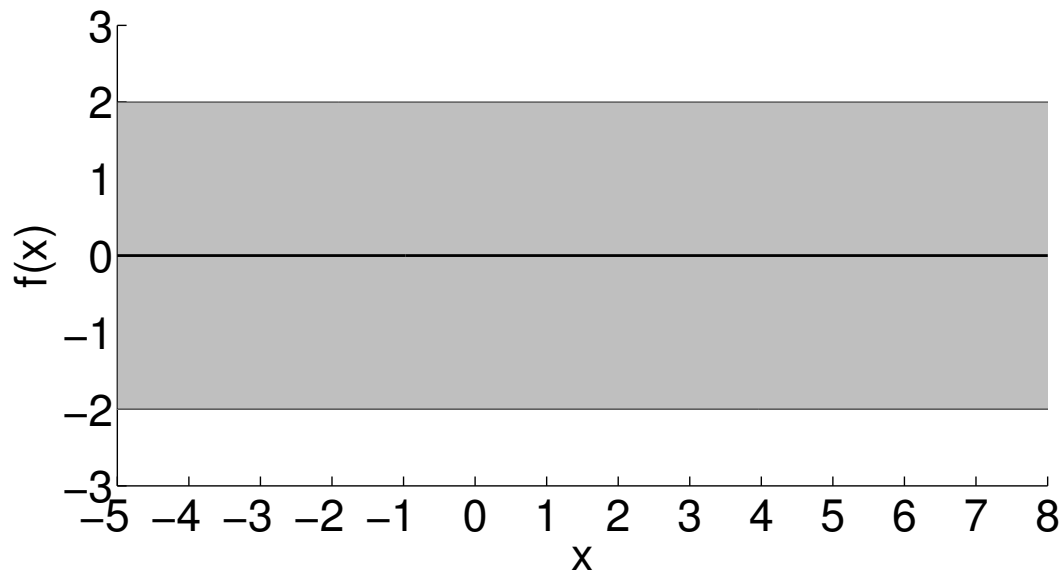
$$p(\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*) = \mathcal{N}(\mathbb{E}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*], \mathbb{V}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*])$$

$$\mathbb{E}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*] = m(\mathbf{X}_*) + k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))$$

$$\mathbb{V}[\mathbf{f}_*|\mathbf{X}, \mathbf{y}, \mathbf{X}_*] = k(\mathbf{X}_*, \mathbf{X}_*) - k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \mathbf{X}_*)$$

Predictions

Make predictions by evaluating the GP posterior mean and covariance function at a finite number of inputs \mathbf{X}_*

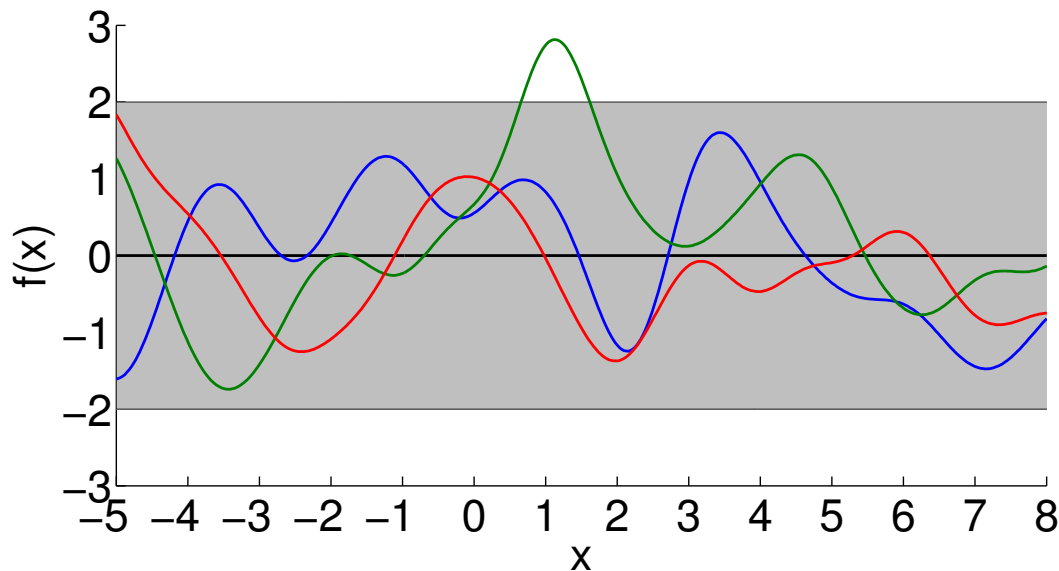


Prior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \emptyset] = m(\mathbf{x}_*) = 0$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \emptyset] = \sigma^2(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x}_*)$$

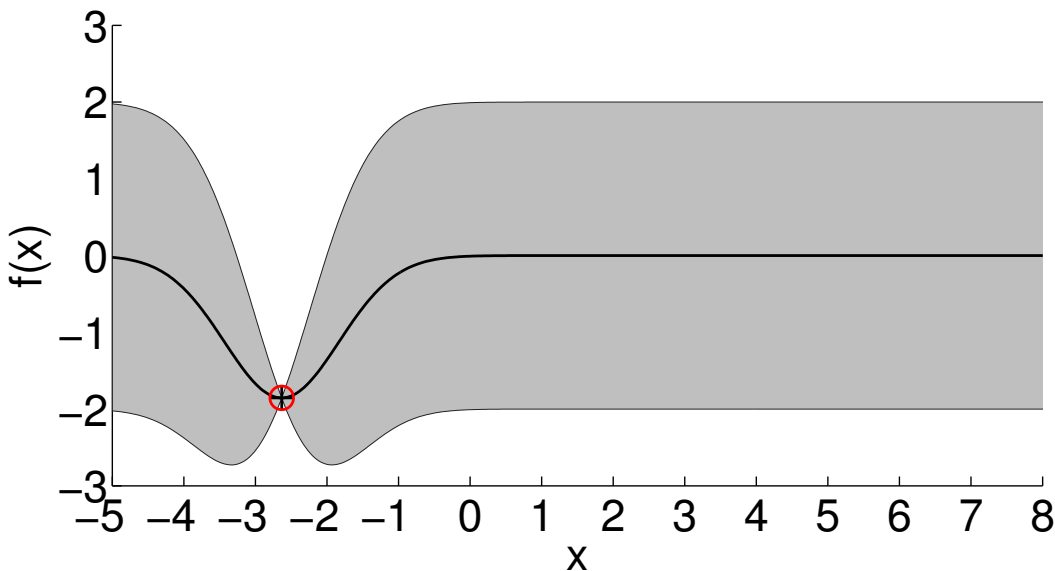


Prior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \emptyset] = m(\mathbf{x}_*) = 0$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \emptyset] = \sigma^2(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x}_*)$$

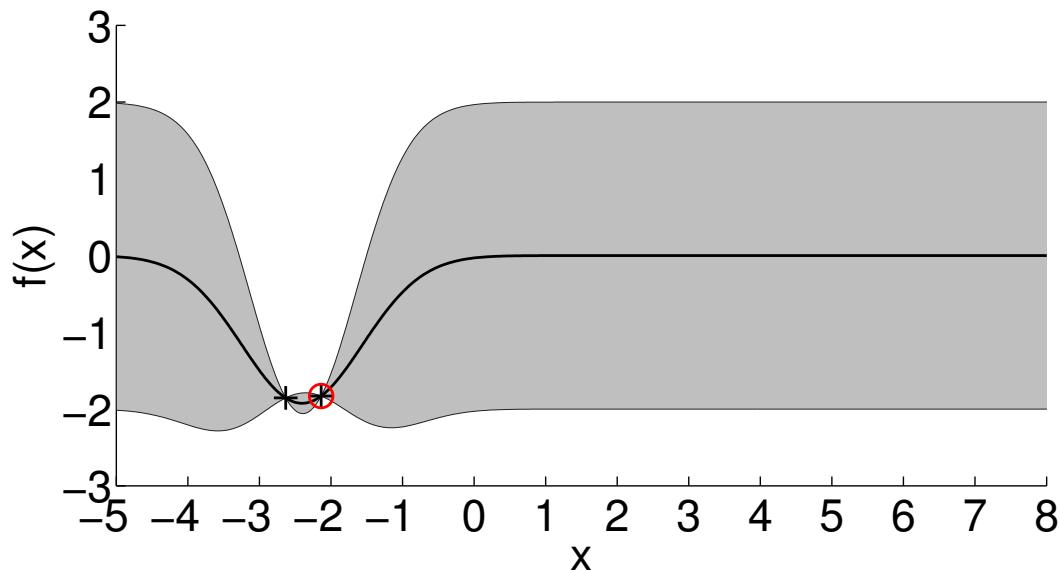


Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_*)$$

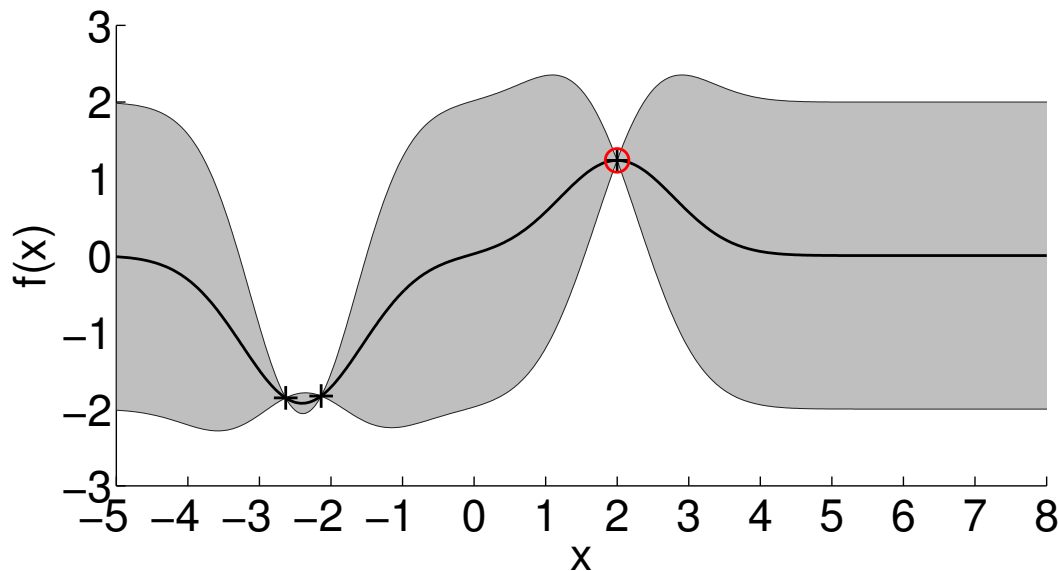


Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_*)$$

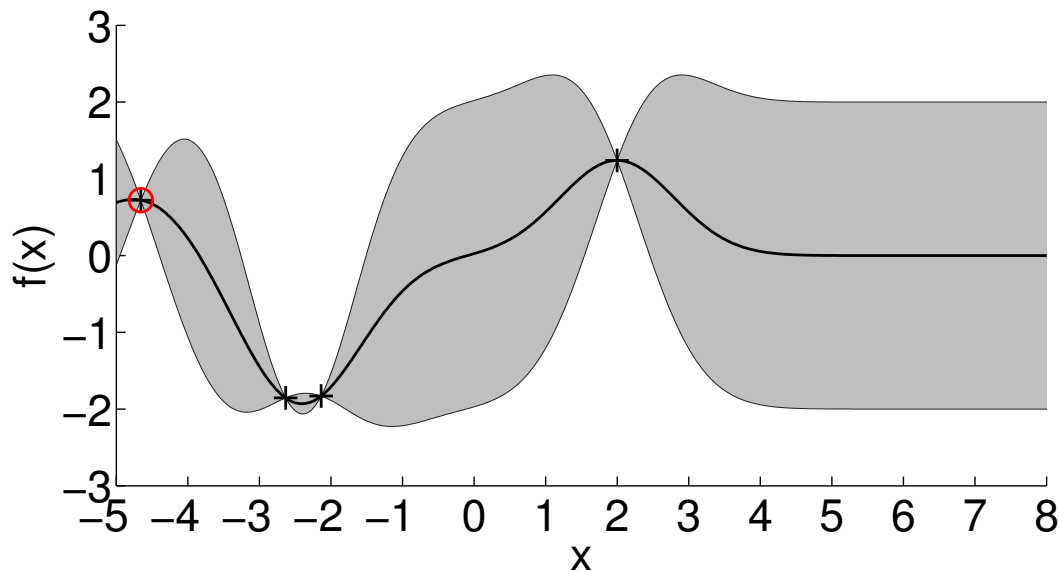


Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_*)$$

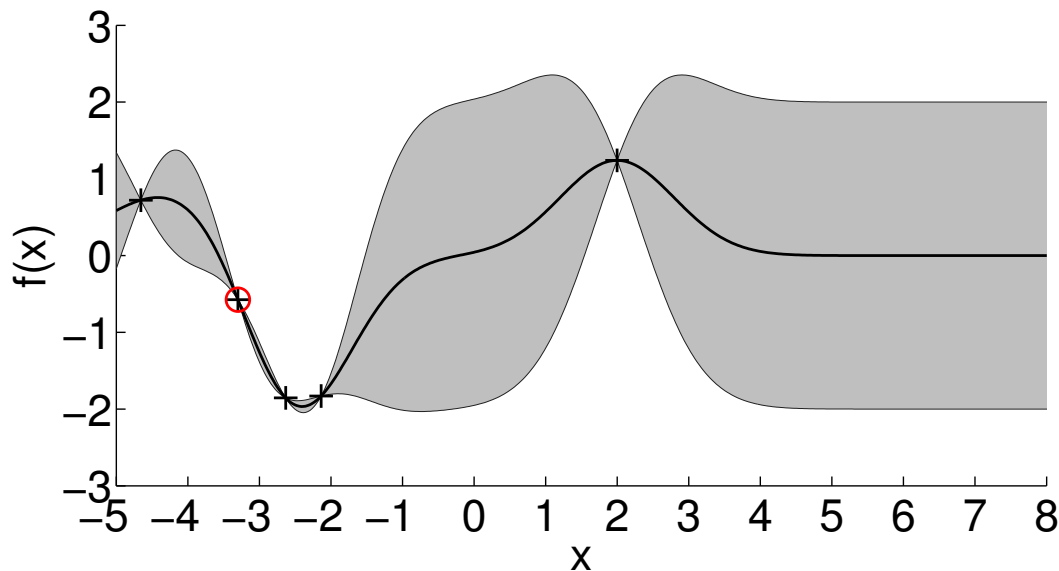


Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$

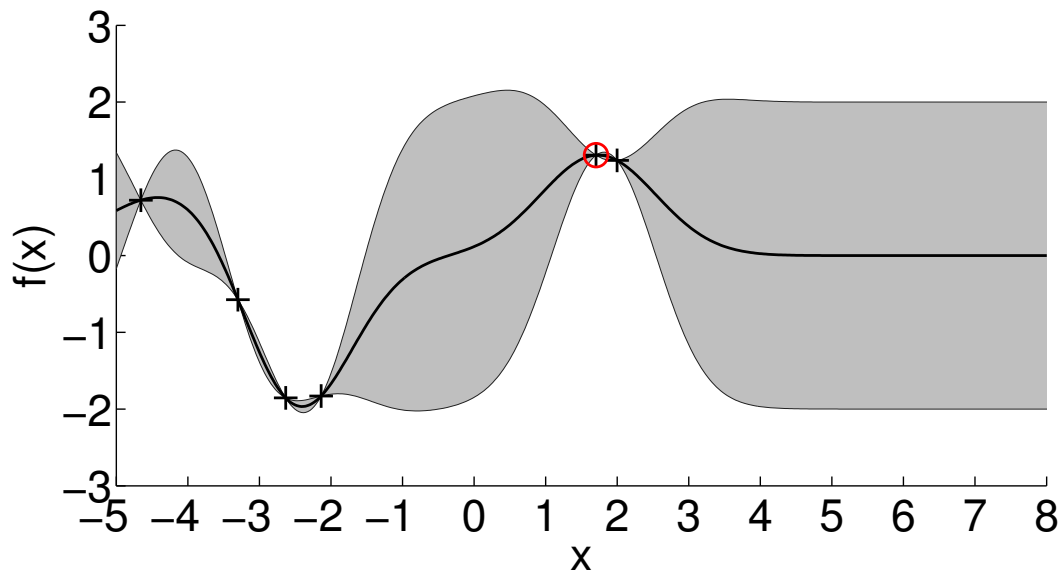


Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$

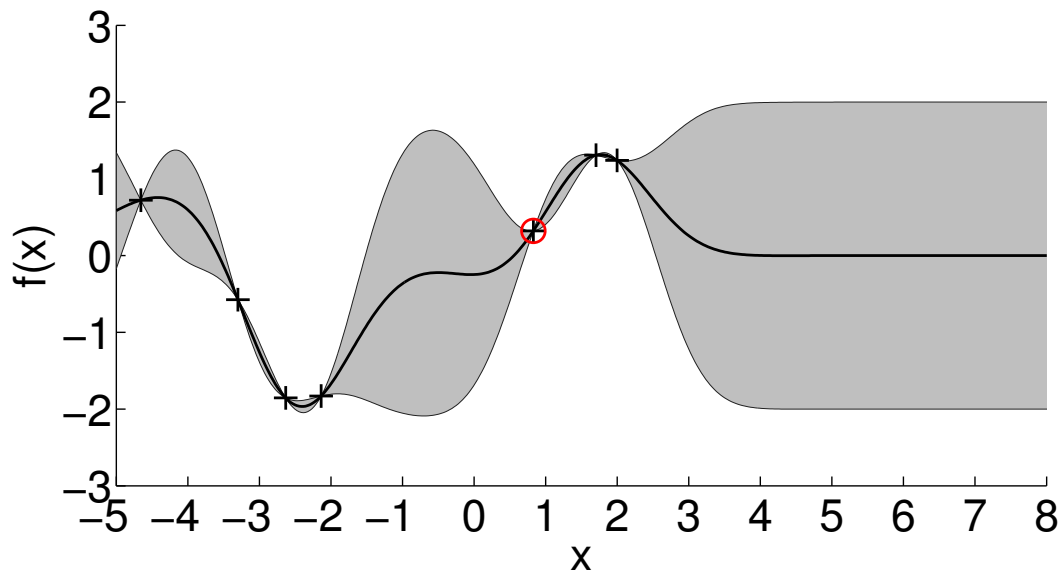


Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$

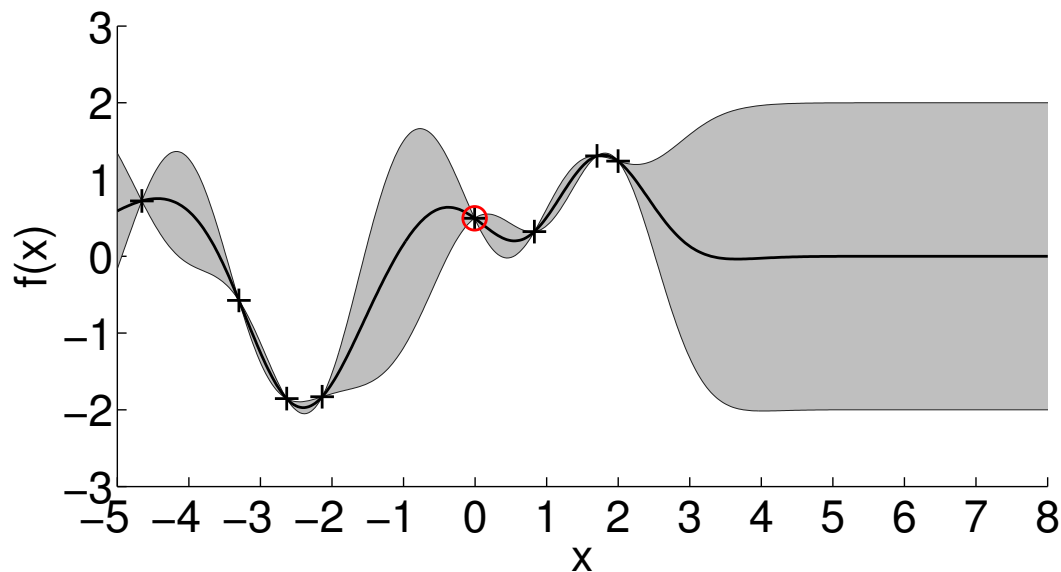


Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$



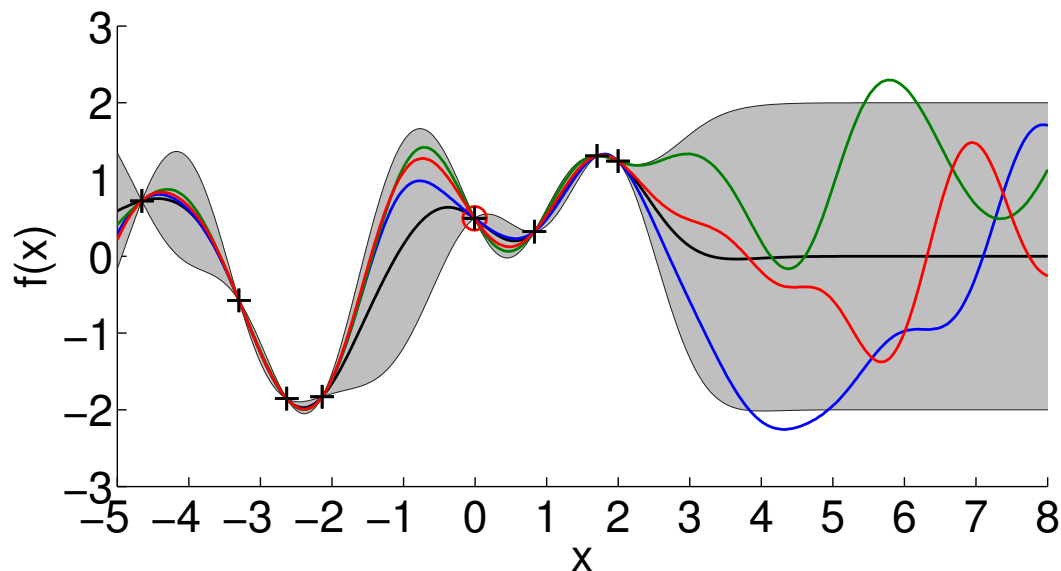
Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = \mathbf{k}(\mathbf{X}, \mathbf{x}_*)^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = \sigma^2(\mathbf{x}_*) =$$

$$\mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{X}, \mathbf{x}_*)^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = m(\mathbf{x}_*) = \mathbf{k}(\mathbf{X}, \mathbf{x}_*)^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}] = \sigma^2(\mathbf{x}_*) =$$

$$\mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{X}, \mathbf{x}_*)^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$

- [1] G. Bertone, M. P. Deisenroth, J. S. Kim, S. Liem, R. R. de Austri, and M. Welling. Accelerating the BSM Interpretation of LHC Data with Machine Learning. arXiv preprint arXiv:1611.02704, 2016.
- [2] R. Calandra, J. Peters, C. E. Rasmussen, and M. P. Deisenroth. Manifold Gaussian Processes for Regression. In *Proceedings of the International Joint Conference on Neural Networks*, 2016.
- [3] Y. Cao and D. J. Fleet. Generalized Product of Experts for Automatic and Principled Fusion of Gaussian Process Predictions. <http://arxiv.org/abs/1410.7827>, 2014.
- [4] N. A. C. Cressie. *Statistics for Spatial Data*. Wiley-Interscience, 1993.
- [5] M. Cutler and J. P. How. Efficient Reinforcement Learning for Robots using Informative Simulated Priors. In *Proceedings of the International Conference on Robotics and Automation*, 2015.
- [6] M. P. Deisenroth and J. W. Ng. Distributed Gaussian Processes. In *Proceedings of the International Conference on Machine Learning*, 2015.
- [7] M. P. Deisenroth, C. E. Rasmussen, and D. Fox. Learning to Control a Low-Cost Manipulator using Data-Efficient Reinforcement Learning. In *Proceedings of Robotics: Science and Systems*, 2011.
- [8] M. P. Deisenroth, C. E. Rasmussen, and J. Peters. Gaussian Process Dynamic Programming. *Neurocomputing*, 72(7–9):1508–1524, Mar. 2009.
- [9] M. P. Deisenroth, R. Turner, M. Huber, U. D. Hanebeck, and C. E. Rasmussen. Robust Filtering and Smoothing with Gaussian Processes. *IEEE Transactions on Automatic Control*, 57(7):1865–1871, 2012.
- [10] R. Frigola, F. Lindsten, T. B. Schön, and C. E. Rasmussen. Bayesian Inference and Learning in Gaussian Process State-Space Models with Particle MCMC. In *Advances in Neural Information Processing Systems*. 2013.
- [11] N. HajiGhassemi and M. P. Deisenroth. Approximate Inference for Long-Term Forecasting with Periodic Gaussian Processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2014.
- [12] J. Hensman, N. Durrande, and A. Solin. Variational Fourier Features for Gaussian Processes. pages 1–52.

- [13] J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian Processes for Big Data. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, 2013.
- [14] A. Krause, A. Singh, and C. Guestrin. Near-Optimal Sensor Placements in Gaussian Processes: Theory, Efficient Algorithms and Empirical Studies. *Journal of Machine Learning Research*, 9:235–284, Feb. 2008.
- [15] M. C. H. Lee, H. Salimbeni, M. P. Deisenroth, and B. Glocker. Patch Kernels for Gaussian Processes in High-Dimensional Imaging Problems. In *NIPS Workshop on Practical Bayesian Nonparametrics*, 2016.
- [16] J. R. Lloyd, D. Duvenaud, R. Grosse, J. B. Tenenbaum, and Z. Ghahramani. Automatic Construction and Natural-Language Description of Nonparametric Regression Models. In *AAAI Conference on Artificial Intelligence*, pages 1–11, 2014.
- [17] D. J. C. MacKay. Introduction to Gaussian Processes. In C. M. Bishop, editor, *Neural Networks and Machine Learning*, volume 168, pages 133–165. Springer, Berlin, Germany, 1998.
- [18] A. G. d. G. Matthews, J. Hensman, R. Turner, and Z. Ghahramani. On Sparse Variational Methods and the Kullback-Leibler Divergence between Stochastic Processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2016.
- [19] M. A. Osborne, S. J. Roberts, A. Rogers, S. D. Ramchurn, and N. R. Jennings. Towards Real-Time Information Processing of Sensor Network Data Using Computationally Efficient Multi-output Gaussian Processes. In *Proceedings of the International Conference on Information Processing in Sensor Networks*, pages 109–120. IEEE Computer Society, 2008.
- [20] J. Quiñonero-Candela and C. E. Rasmussen. A Unifying View of Sparse Approximate Gaussian Process Regression. *Journal of Machine Learning Research*, 6(2):1939–1960, 2005.
- [21] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- [22] S. Roberts, M. A. Osborne, M. Ebden, S. Reece, N. Gibson, and S. Aigrain. Gaussian Processes for Time Series Modelling. *Philosophical Transactions of the Royal Society (Part A)*, 371(1984), Feb. 2013.

- [23] B. Schölkopf and A. J. Smola. *Learning with Kernels—Support Vector Machines, Regularization, Optimization, and Beyond*. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA, 2002.
- [24] E. Snelson and Z. Ghahramani. Sparse Gaussian Processes using Pseudo-inputs. In Y. Weiss, B. Schölkopf, and J. C. Platt, editors, *Advances in Neural Information Processing Systems 18*, pages 1257–1264. The MIT Press, Cambridge, MA, USA, 2006.
- [25] M. K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2009.
- [26] V. Tresp. A Bayesian Committee Machine. *Neural Computation*, 12(11):2719–2741, 2000.
- [27] A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. Deep Kernel Learning. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2016.