

Gaussian Processes

Recommended reading:

Rasmussen/Williams: Chapters 1, 2, 4, 5

Deisenroth & Ng (2015)[3]

Marc Deisenroth

Department of Computing
Imperial College London

February 22, 2017

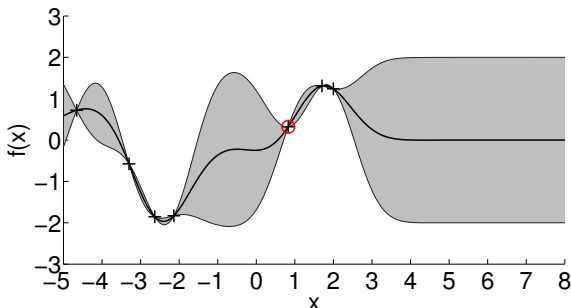
Gaussian Processes for Machine Learning



Carl Edward Rasmussen and Christopher K. I. Williams

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

Problem Setting

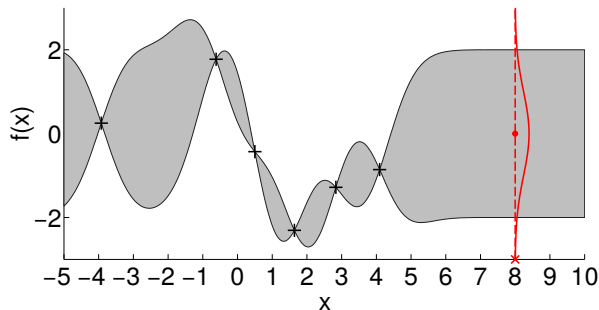


Objective

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

► Probabilistic regression problem

Problem Setting



Objective

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

► Probabilistic regression problem

Recap from CO-496: Bayesian Linear Regression

- ▶ Linear Regression Model:

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

Recap from CO-496: Bayesian Linear Regression

- ▶ Linear Regression Model:

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

- ▶ Integrating out the parameters when predicting leads to a distribution over functions:

$$p(f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$
$$= \mathcal{N}(\mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

Recap from CO-496: Bayesian Linear Regression

- ▶ Linear Regression Model:

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

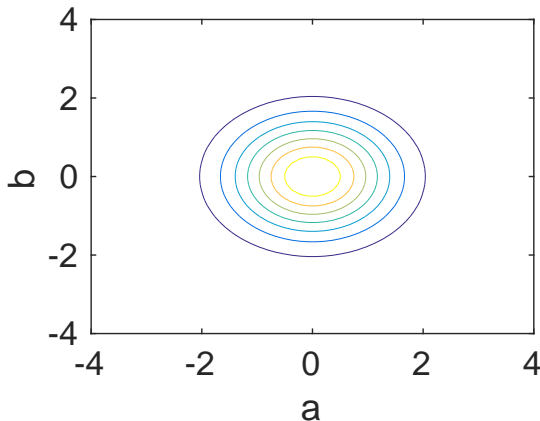
- ▶ Integrating out the parameters when predicting leads to a distribution over functions:

$$p(f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(f(\mathbf{x}_*) | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$
$$= \mathcal{N}(\mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$
$$\mu(\mathbf{x}_*) = \phi_*^\top \Sigma_p \Phi (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$
$$\sigma^2(\mathbf{x}_*) = \phi_*^\top \Sigma_p \phi_* - \phi_*^\top \Sigma_p \Phi (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \Phi^\top \Sigma_p \phi_*$$
$$\mathbf{K} = \Phi^\top \Sigma_p \Phi$$

Sampling from the Prior over Functions

Consider a linear regression setting

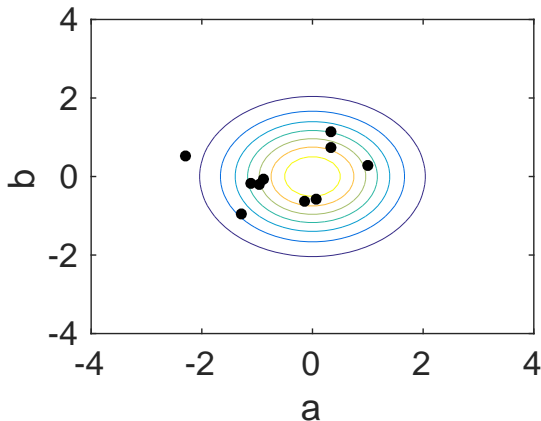
$$y = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Sampling from the Prior over Functions

Consider a linear regression setting

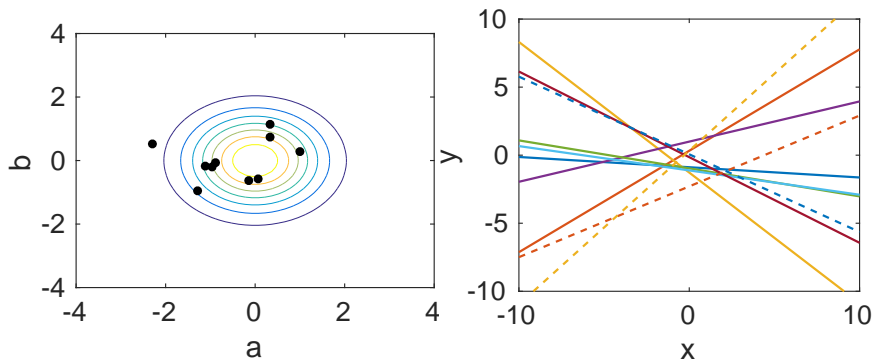
$$y = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Sampling from the Prior over Functions

Consider a linear regression setting

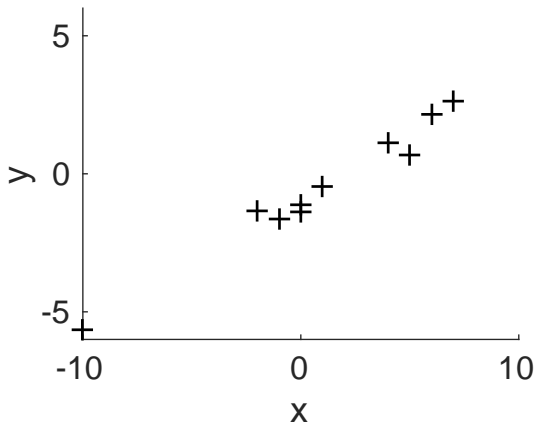
$$y = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Sampling from the Prior over Functions

Consider a linear regression setting

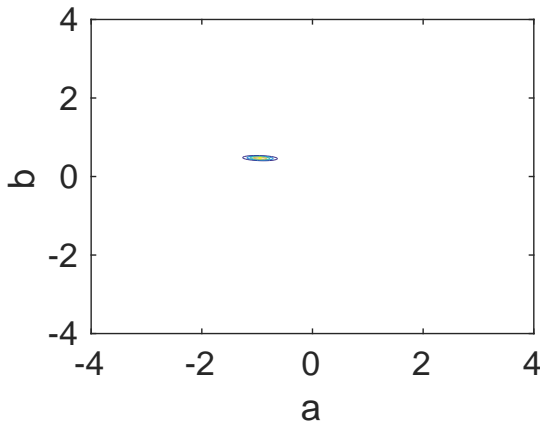
$$y = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Sampling from the Posterior over Functions

Consider a linear regression setting

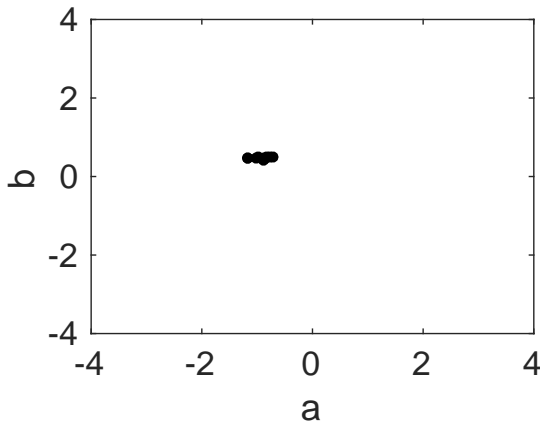
$$y = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Sampling from the Posterior over Functions

Consider a linear regression setting

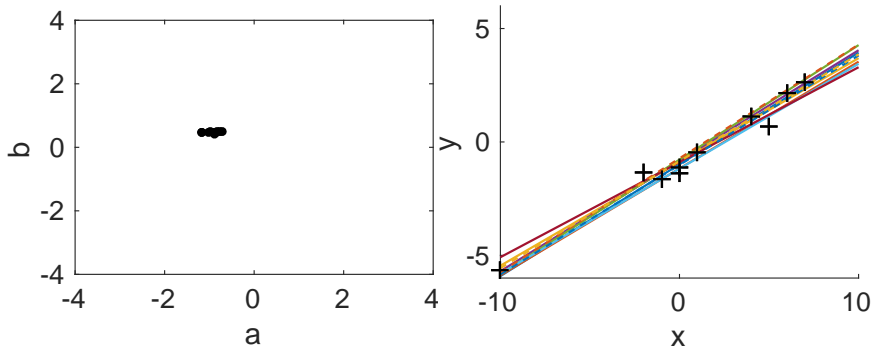
$$y = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Sampling from the Posterior over Functions

Consider a linear regression setting

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



Fitting Nonlinear Functions

- ▶ Fit nonlinear functions using (Bayesian) linear regression:
Linear combination of nonlinear features
- ▶ Example: Radial-basis-function (RBF) network

$$f(\mathbf{x}) = \sum_{i=1}^n w_i \phi_i(\mathbf{x}), \quad w_i \sim \mathcal{N}(0, \sigma_p^2)$$

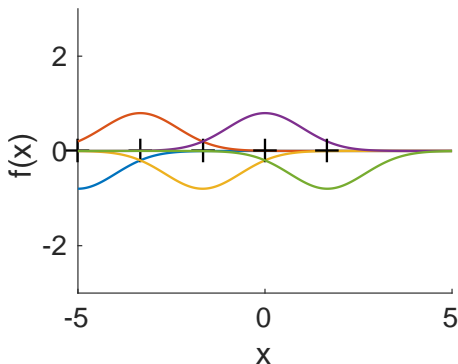
where

$$\phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^\top (\mathbf{x} - \boldsymbol{\mu}_i)\right)$$

for given “centers” $\boldsymbol{\mu}_i$

Illustration: Fitting a Radial Basis Function Network

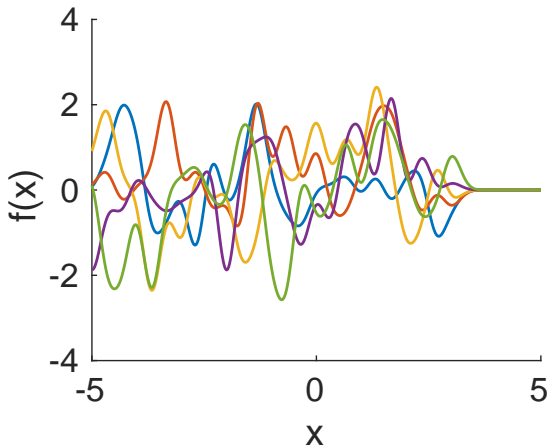
$$\phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^\top(\mathbf{x} - \boldsymbol{\mu}_i)\right)$$



- ▶ Place Gaussian-shaped basis functions ϕ_i at 25 input locations μ_i , linearly spaced in the interval $[-5, 3]$

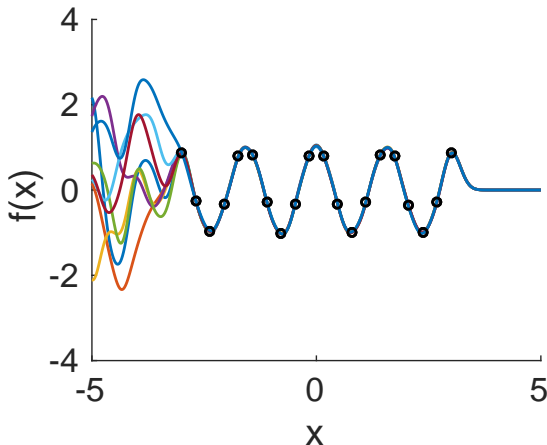
Samples from the RBF Prior

$$f(\mathbf{x}) = \sum_{i=1}^n w_i \phi_i(\mathbf{x}), \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

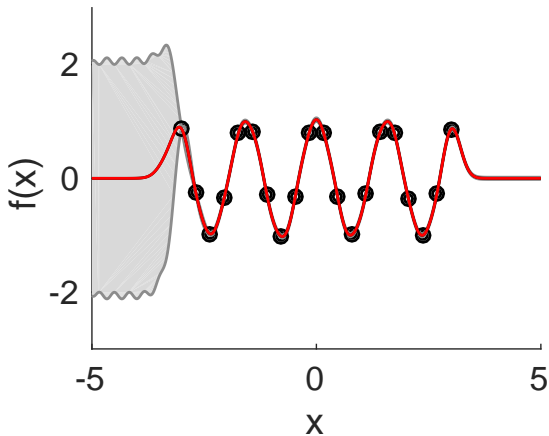


Samples from the RBF Posterior

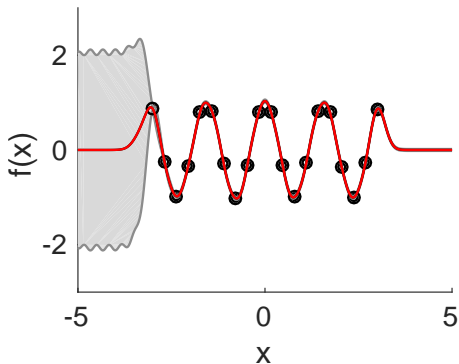
$$f(\mathbf{x}) = \sum_{i=1}^n w_i \phi_i(\mathbf{x}), \quad p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N)$$



RBF Posterior



Limitations



- ▶ Feature engineering
- ▶ Finite number of features:
 - ▶ Above: Without basis functions on the right, we cannot express any variability of the function
 - ▶ Ideally: Add more (infinitely many) basis functions

Approach

- ▶ Instead of sampling parameters, which induce a distribution over functions, sample functions directly
 - ▶▶ 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

- ▶ 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.

Gaussian Process

- ▶ 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

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

Gaussian Process

- ▶ 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

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)** $k(\cdot, \cdot)$

Covariance Function

- ▶ The covariance function (kernel) is symmetric and positive semi-definite
- ▶ It allows us to compute covariances between (unknown) function values by just looking at the corresponding inputs:

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

GP Regression as a Bayesian Inference Problem

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** $p(f|\mathbf{X}, \mathbf{y})$ that explains the data

GP Regression as a Bayesian Inference Problem

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** $p(f|\mathbf{X}, \mathbf{y})$ that explains the data

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

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

GP Regression as a Bayesian Inference Problem

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** $p(f|\mathbf{X}, \mathbf{y})$ that explains the data

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

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

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

GP Regression as a Bayesian Inference Problem

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** $p(f|\mathbf{X}, \mathbf{y})$ that explains the data

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

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

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

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

GP Regression as a Bayesian Inference Problem

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** $p(f|\mathbf{X}, \mathbf{y})$ that explains the data

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

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

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

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

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

GP Regression as a Bayesian Inference Problem

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** $p(f|\mathbf{X}, \mathbf{y})$ that explains the data

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

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

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

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

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

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

Prior over Functions

- ▶ Treat a function as a long vector of function values:

$$f = [f_1, f_2, \dots]$$

- ▶▶ Look at a **distribution over function values** $f_i = f(\mathbf{x}_i)$

Prior over Functions

- ▶ Treat a function as a long vector of function values:

$$f = [f_1, f_2, \dots]$$

- ▶ Look at a **distribution over function values** $f_i = f(\mathbf{x}_i)$
- ▶ Consider a finite number of N function values \mathbf{f} and all other (infinitely many) function values $\tilde{\mathbf{f}}$. Informally:

$$p(\mathbf{f}, \tilde{\mathbf{f}}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_f \\ \boldsymbol{\mu}_{\tilde{\mathbf{f}}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{ff} & \boldsymbol{\Sigma}_{f\tilde{\mathbf{f}}} \\ \boldsymbol{\Sigma}_{\tilde{\mathbf{f}}f} & \boldsymbol{\Sigma}_{\tilde{\mathbf{f}}\tilde{\mathbf{f}}} \end{bmatrix} \right)$$

where $\boldsymbol{\Sigma}_{\tilde{\mathbf{f}}\tilde{\mathbf{f}}} \in \mathbb{R}^{m \times m}$ and $\boldsymbol{\Sigma}_{f\tilde{\mathbf{f}}} \in \mathbb{R}^{N \times m}$, $m \rightarrow \infty$.

- ▶ $\Sigma_{ff}^{(i,j)} = \text{Cov}[f(\mathbf{x}_i), f(\mathbf{x}_j)] = k(\mathbf{x}_i, \mathbf{x}_j)$

Prior over Functions

- ▶ Treat a function as a long vector of function values:

$$f = [f_1, f_2, \dots]$$

- ▶ Look at a **distribution over function values** $f_i = f(\mathbf{x}_i)$
- ▶ Consider a finite number of N function values \mathbf{f} and all other (infinitely many) function values $\tilde{\mathbf{f}}$. Informally:

$$p(\mathbf{f}, \tilde{\mathbf{f}}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_f \\ \boldsymbol{\mu}_{\tilde{\mathbf{f}}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{ff} & \boldsymbol{\Sigma}_{f\tilde{\mathbf{f}}} \\ \boldsymbol{\Sigma}_{\tilde{\mathbf{f}}f} & \boldsymbol{\Sigma}_{\tilde{\mathbf{f}}\tilde{\mathbf{f}}} \end{bmatrix} \right)$$

where $\boldsymbol{\Sigma}_{\tilde{\mathbf{f}}\tilde{\mathbf{f}}} \in \mathbb{R}^{m \times m}$ and $\boldsymbol{\Sigma}_{f\tilde{\mathbf{f}}} \in \mathbb{R}^{N \times m}$, $m \rightarrow \infty$.

- ▶ $\Sigma_{ff}^{(i,j)} = \text{Cov}[f(\mathbf{x}_i), f(\mathbf{x}_j)] = k(\mathbf{x}_i, \mathbf{x}_j)$
- ▶ Key property: The **marginal remains finite**

$$p(\mathbf{f}) = \int p(\mathbf{f}, \tilde{\mathbf{f}}) d\tilde{\mathbf{f}} = \mathcal{N}(\boldsymbol{\mu}_f, \boldsymbol{\Sigma}_{ff})$$

Training and Test Marginal

- ▶ In practice, we always have **finite training and test inputs**
 $\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}$.
- ▶ Define $f_* := f_{\text{test}}, f := f_{\text{train}}$.

Training and Test Marginal

- ▶ In practice, we always have **finite training and test inputs** $\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}$.
- ▶ Define $\mathbf{f}_* := \mathbf{f}_{\text{test}}, \mathbf{f} := \mathbf{f}_{\text{train}}$.
- ▶ Then, we obtain the finite **marginal**

$$p(\mathbf{f}, \mathbf{f}_*) = \int p(\mathbf{f}, \mathbf{f}_*, \mathbf{f}_{\text{other}}) d\mathbf{f}_{\text{other}} = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_f \\ \boldsymbol{\mu}_* \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{ff} & \boldsymbol{\Sigma}_{f_*} \\ \boldsymbol{\Sigma}_{*f} & \boldsymbol{\Sigma}_{**} \end{bmatrix} \right)$$

GP Regression as a Bayesian Inference Problem (ctd.)

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

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

GP Regression as a Bayesian Inference Problem (ctd.)

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

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

Using the properties of Gaussians, we obtain

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

$$\mathbf{K} = k(\mathbf{X}, \mathbf{X})$$

GP Regression as a Bayesian Inference Problem (ctd.)

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

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

Using the properties of Gaussians, we obtain

$$\begin{aligned} p(\mathbf{y}|f, \mathbf{X}) p(f|\mathbf{X}) &= \mathcal{N}(\mathbf{y} | f(\mathbf{X}), \sigma_n^2 \mathbf{I}) \mathcal{N}(f(\mathbf{X}) | m(\mathbf{X}), \mathbf{K}) \\ &= Z \mathcal{N}(f(\mathbf{X}) | \underbrace{m(\mathbf{X}) + \mathbf{K}(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))}_{\text{posterior mean}}, \underbrace{\mathbf{K} - \mathbf{K}(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}\mathbf{K}}_{\text{posterior covariance}}) \end{aligned}$$

$$\mathbf{K} = k(\mathbf{X}, \mathbf{X})$$

GP Regression as a Bayesian Inference Problem (ctd.)

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

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

Using the properties of Gaussians, we obtain

$$\begin{aligned} p(\mathbf{y}|f, \mathbf{X}) p(f|\mathbf{X}) &= \mathcal{N}(\mathbf{y} | f(\mathbf{X}), \sigma_n^2 \mathbf{I}) \mathcal{N}(f(\mathbf{X}) | m(\mathbf{X}), \mathbf{K}) \\ &= Z \mathcal{N}(f(\mathbf{X}) | \underbrace{m(\mathbf{X}) + \mathbf{K}(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}(\mathbf{y} - m(\mathbf{X}))}_{\text{posterior mean}}, \underbrace{\mathbf{K} - \mathbf{K}(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}\mathbf{K}}_{\text{posterior covariance}}) \end{aligned}$$

$$\mathbf{K} = k(\mathbf{X}, \mathbf{X})$$

Marginal likelihood:

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

GP Predictions (1)

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

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

GP Predictions (1)

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

- ▶ **Objective:** Find $p(f(\mathbf{X}_*)|\mathbf{X}, \mathbf{y})$ for training data \mathbf{X}, \mathbf{y} and test inputs \mathbf{X}_* .
- ▶ GP prior: $p(f|\mathbf{X}) = \mathcal{N}(m(\mathbf{X}), \mathbf{K})$
- ▶ 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 f, f_* are jointly Gaussian distributed:

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

GP Predictions (1)

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

- ▶ **Objective:** Find $p(f(\mathbf{X}_*)|\mathbf{X}, \mathbf{y})$ for training data \mathbf{X}, \mathbf{y} and test inputs \mathbf{X}_* .
- ▶ GP prior: $p(f|\mathbf{X}) = \mathcal{N}(m(\mathbf{X}), \mathbf{K})$
- ▶ 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 f, f_* are jointly Gaussian distributed:

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

- ▶ Due to the Gaussian likelihood, we also get (f is unobserved)

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

GP Predictions (2)

Prior:

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

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

GP Predictions (2)

Prior:

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

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

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

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

$$\begin{aligned} \mathbb{V}[f_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] &= k_{\text{post}}(\mathbf{X}_*, \mathbf{X}_*) \\ &= \underbrace{k(\mathbf{X}_*, \mathbf{X}_*)}_{\text{prior variance}} - k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \mathbf{X}_*) \end{aligned}$$

GP Predictions (2)

Prior:

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

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

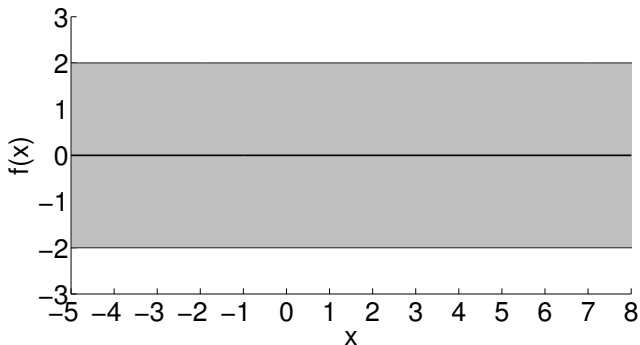
$$p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*) = \mathcal{N}(\mathbb{E}[f_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*], \mathbb{V}[f_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*])$$

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

$$\begin{aligned} \mathbb{V}[f_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] &= k_{\text{post}}(\mathbf{X}_*, \mathbf{X}_*) \\ &= \underbrace{k(\mathbf{X}_*, \mathbf{X}_*)}_{\text{prior variance}} - k(\mathbf{X}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}k(\mathbf{X}, \mathbf{X}_*) \end{aligned}$$

From now: Set prior mean function $m \equiv 0$

Illustration: Inference with Gaussian Processes



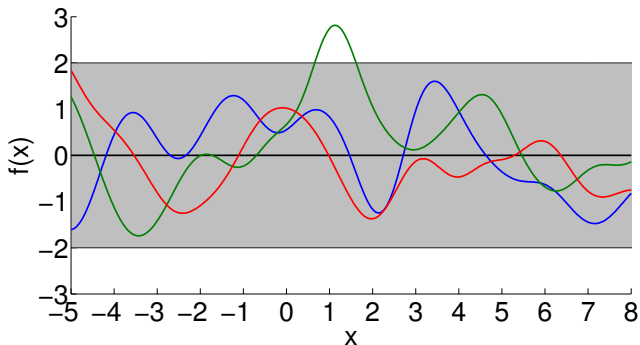
Prior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



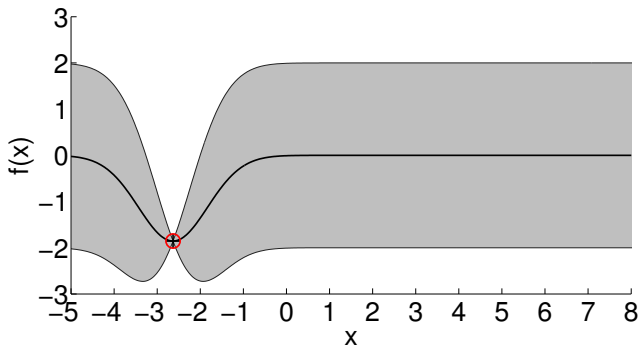
Prior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



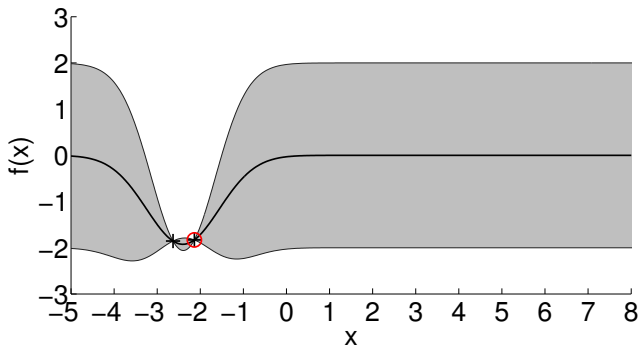
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



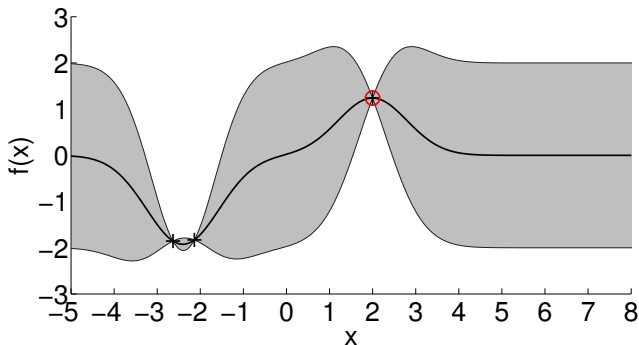
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



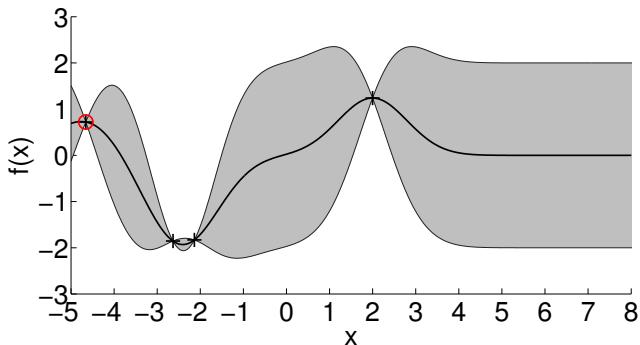
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



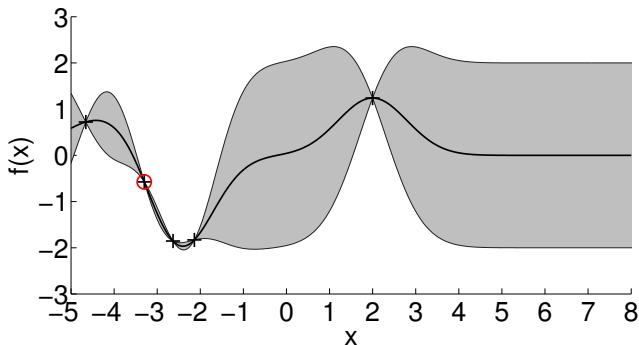
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



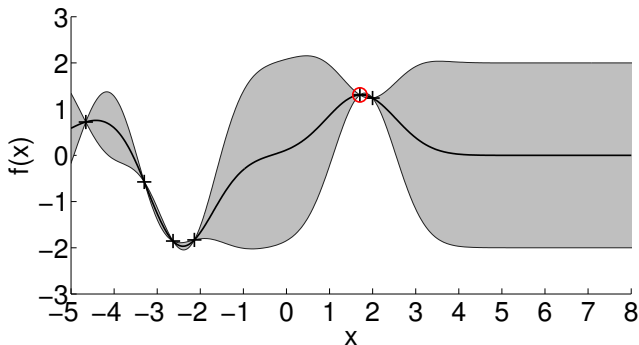
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



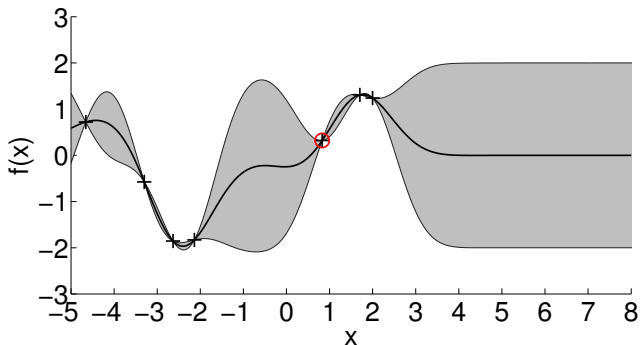
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



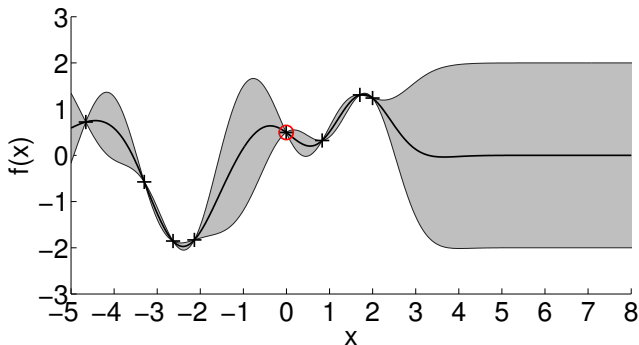
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



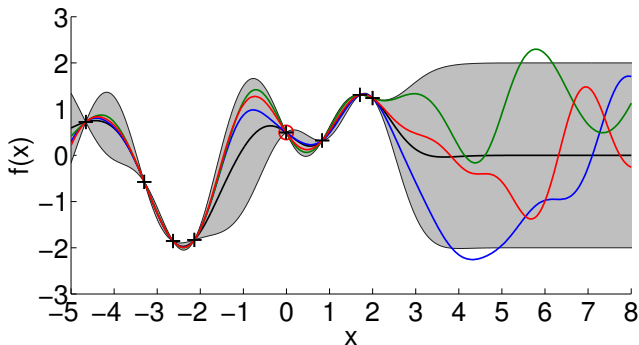
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration: Inference with Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

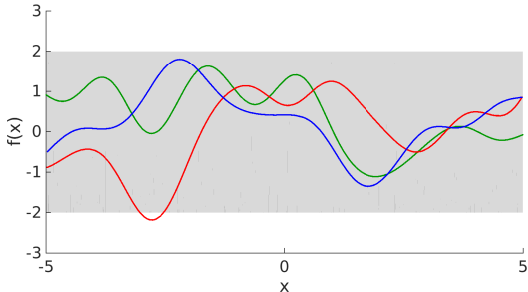
Covariance Function

- ▶ A Gaussian process is fully specified by a **mean function** m and a **kernel/covariance function** k
- ▶ The covariance function (kernel) is symmetric and positive semi-definite
- ▶ Covariance function encodes **high-level structural assumptions** about the latent function f (e.g., smoothness, differentiability, periodicity)

Gaussian Covariance Function

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

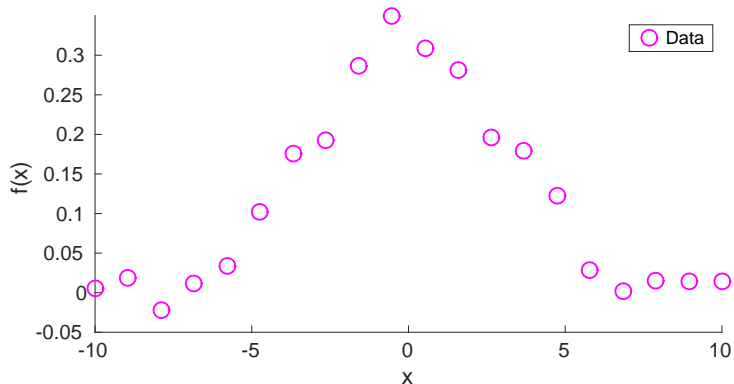
- ▶ σ_f : **Amplitude** of the latent function
- ▶ ℓ : **Length scale**. How far do we have to move in input space before the function value changes significantly
- ▶▶ **Smoothness parameter**



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

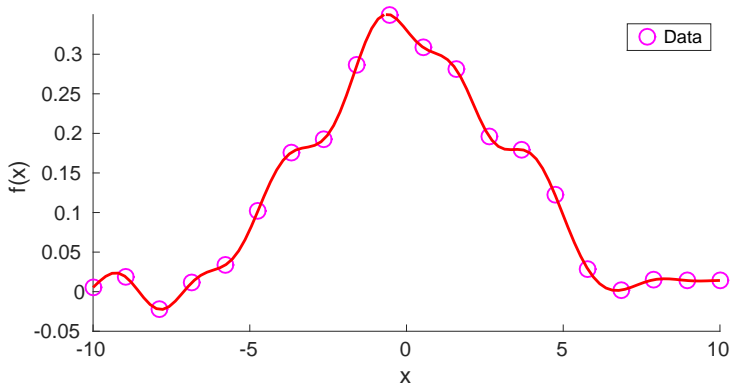
Length-Scales

Length scales determine how wiggly the function is and how much information we can transfer to other function values



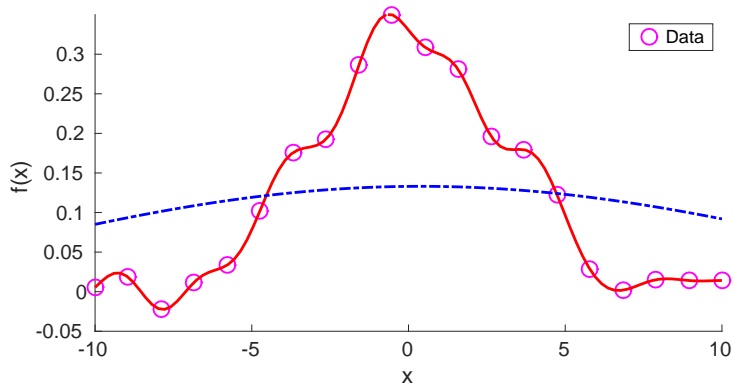
Length-Scales

Length scales determine how wiggly the function is and how much information we can transfer to other function values



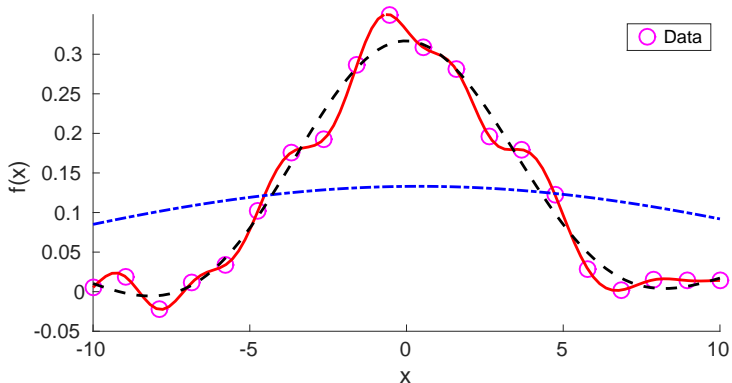
Length-Scales

Length scales determine how wiggly the function is and how much information we can transfer to other function values



Length-Scales

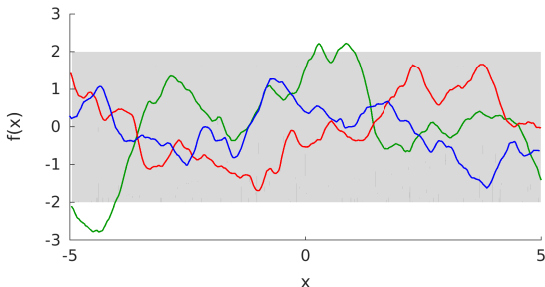
Length scales determine how wiggly the function is and how much information we can transfer to other function values



Matérn Covariance Function

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

- ▶ 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?

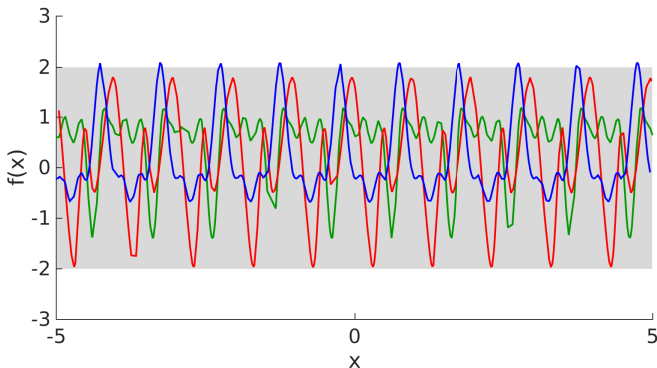


- ▶ Assumption on latent function: 1-times differentiable

Periodic Covariance Function

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

κ : Periodicity parameter



Meta-Parameters of a GP

The GP possesses a set of hyper-parameters:

- ▶ Parameters of the mean function
- ▶ Hyper-parameters of the covariance function (e.g., length-scales and signal variance)
- ▶ Likelihood parameters (e.g., noise variance σ_n^2)

Meta-Parameters of a GP

The GP possesses a set of hyper-parameters:

- ▶ Parameters of the mean function
 - ▶ Hyper-parameters of the covariance function (e.g., length-scales and signal variance)
 - ▶ Likelihood parameters (e.g., noise variance σ_n^2)
- ▶▶ Train a GP to find a good set of hyper-parameters

Meta-Parameters of a GP

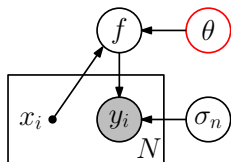
The GP possesses a set of hyper-parameters:

- ▶ Parameters of the mean function
 - ▶ Hyper-parameters of the covariance function (e.g., length-scales and signal variance)
 - ▶ Likelihood parameters (e.g., noise variance σ_n^2)
- ▶▶ Train a GP to find a good set of hyper-parameters
- ▶▶ Model selection to find good mean and covariance functions (can also be automated Automatic Statistician (Lloyd et al., 2014))

Gaussian Process Training: Hyper-Parameters

GP Training

Find good GP hyper-parameters θ (kernel and mean function parameters)



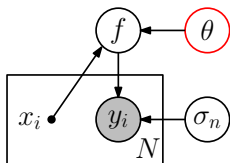
Gaussian Process Training: Hyper-Parameters

GP Training

Find good GP hyper-parameters θ (kernel and mean function parameters)

- ▶ Place a prior $p(\theta)$ on hyper-parameters
- ▶ Posterior over hyper-parameters:

$$p(\theta|\mathbf{X}, \mathbf{y}) = \frac{p(\theta) p(\mathbf{y}|\mathbf{X}, \theta)}{p(\mathbf{y}|\mathbf{X})}, \quad p(\mathbf{y}|\mathbf{X}, \theta) = \int p(\mathbf{y}|f(\mathbf{X}))p(f|\mathbf{X}, \theta)df$$



Gaussian Process Training: Hyper-Parameters

GP Training

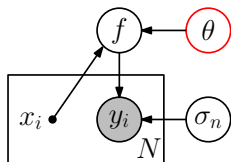
Find good GP hyper-parameters θ (kernel and mean function parameters)

- ▶ Place a prior $p(\theta)$ on hyper-parameters
- ▶ Posterior over hyper-parameters:

$$p(\theta|\mathbf{X}, \mathbf{y}) = \frac{p(\theta) p(\mathbf{y}|\mathbf{X}, \theta)}{p(\mathbf{y}|\mathbf{X})}, \quad p(\mathbf{y}|\mathbf{X}, \theta) = \int p(\mathbf{y}|f(\mathbf{X}))p(f|\mathbf{X}, \theta)df$$

- ▶ Choose hyper-parameters θ^* , such that

$$\theta^* \in \arg \max_{\theta} \log p(\theta) + \log p(\mathbf{y}|\mathbf{X}, \theta)$$



Gaussian Process Training: Hyper-Parameters

GP Training

Find good GP hyper-parameters θ (kernel and mean function parameters)

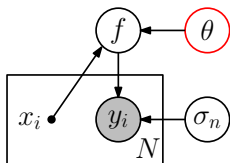
- ▶ Place a prior $p(\theta)$ on hyper-parameters
- ▶ Posterior over hyper-parameters:

$$p(\theta|\mathbf{X}, \mathbf{y}) = \frac{p(\theta) p(\mathbf{y}|\mathbf{X}, \theta)}{p(\mathbf{y}|\mathbf{X})}, \quad p(\mathbf{y}|\mathbf{X}, \theta) = \int p(\mathbf{y}|f(\mathbf{X}))p(f|\mathbf{X}, \theta)df$$

- ▶ Choose hyper-parameters θ^* , such that

$$\theta^* \in \arg \max_{\theta} \log p(\theta) + \log p(\mathbf{y}|\mathbf{X}, \theta)$$

- ▶▶ Maximize **marginal likelihood** if $p(\theta) = \mathcal{U}$ (uniform prior)



Training via Marginal Likelihood Maximization

GP Training

Maximize the evidence/marginal likelihood (probability of the data given the hyper-parameters, where the unwieldy f has been integrated out) ► Also called **Maximum Likelihood-Type-II**

Marginal likelihood:

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

Training via Marginal Likelihood Maximization

GP Training

Maximize the evidence/marginal likelihood (probability of the data given the hyper-parameters, where the unwieldy f has been integrated out) ► Also called **Maximum Likelihood-Type-II**

Marginal likelihood:

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

Learning the GP hyper-parameters:

$$\boldsymbol{\theta}^* \in \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$$

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \mathbf{y}^\top \mathbf{K}_{\boldsymbol{\theta}}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_{\boldsymbol{\theta}}| + \text{const}, \quad \mathbf{K}_{\boldsymbol{\theta}} := \mathbf{K} + \sigma_n^2 \mathbf{I}$$

Training via Marginal Likelihood Maximization

Log-marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}_\theta^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_\theta| + \text{const}, \quad \mathbf{K}_\theta := \mathbf{K} + \sigma_n^2 \mathbf{I}$$

Training via Marginal Likelihood Maximization

Log-marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}_\theta^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_\theta| + \text{const}, \quad \mathbf{K}_\theta := \mathbf{K} + \sigma_n^2 \mathbf{I}$$

- ▶ Automatic trade-off between data fit and model complexity

Training via Marginal Likelihood Maximization

Log-marginal likelihood:

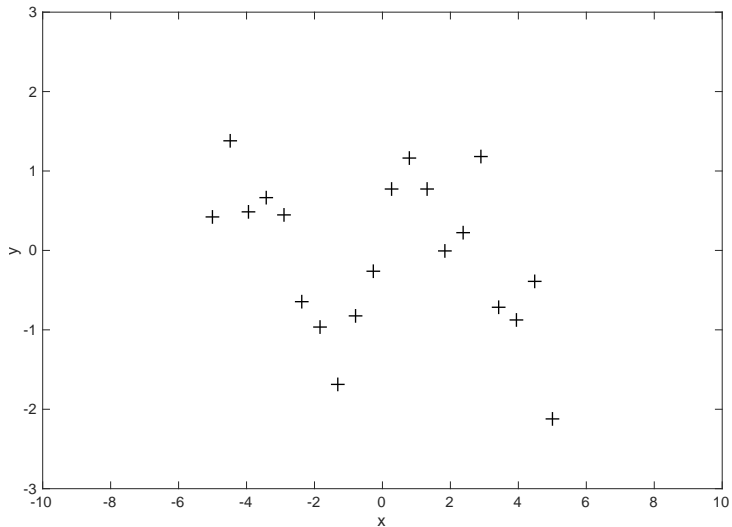
$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}_\theta^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_\theta| + \text{const}, \quad \mathbf{K}_\theta := \mathbf{K} + \sigma_n^2 \mathbf{I}$$

- ▶ Automatic trade-off between data fit and model complexity
- ▶ Gradient-based optimization of hyper-parameters $\boldsymbol{\theta}$:

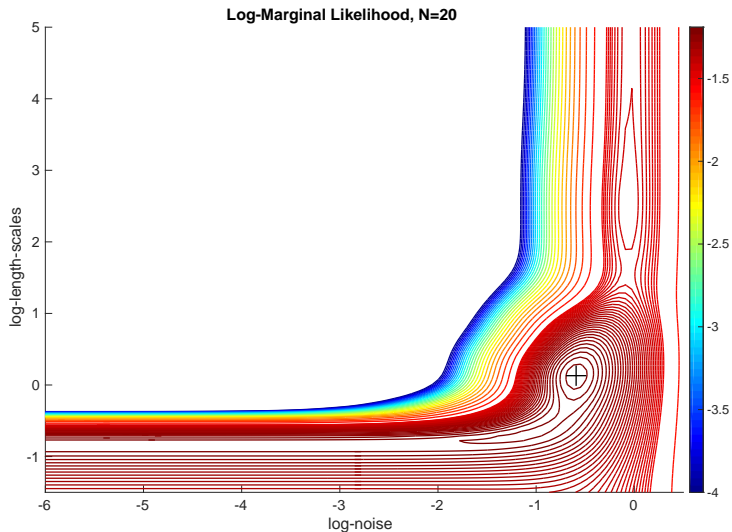
$$\begin{aligned} \frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_i} &= \frac{1}{2} \mathbf{y}^\top \mathbf{K}_\theta^{-1} \frac{\partial \mathbf{K}_\theta}{\partial \theta_i} \mathbf{K}_\theta^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left(\mathbf{K}_\theta^{-1} \frac{\partial \mathbf{K}_\theta}{\partial \theta_i} \right) \\ &= \frac{1}{2} \text{tr} \left((\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - \mathbf{K}_\theta^{-1}) \frac{\partial \mathbf{K}_\theta}{\partial \theta_i} \right), \end{aligned}$$

$$\boldsymbol{\alpha} := \mathbf{K}_\theta^{-1} \mathbf{y}$$

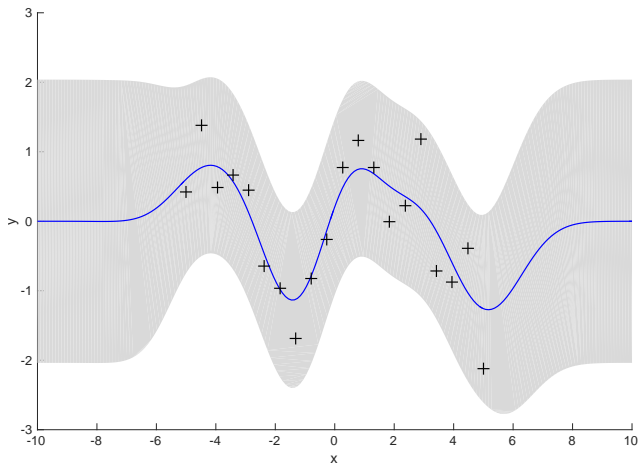
Example: Training Data



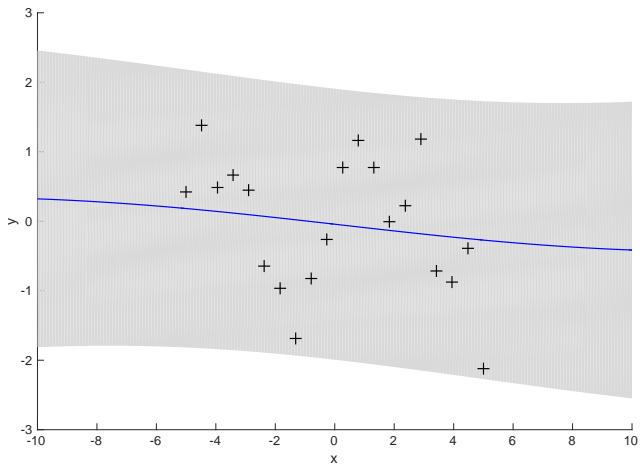
Example: Marginal Likelihood Contour



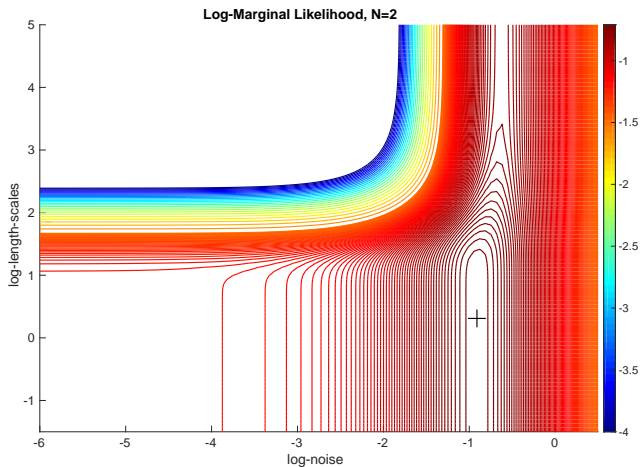
Example: Exploring the Modes (1)



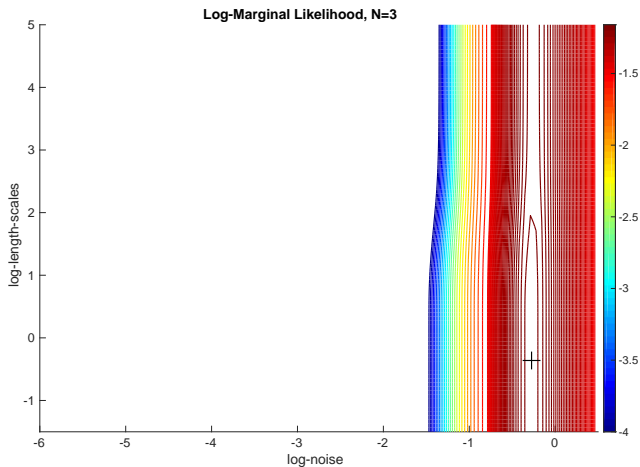
Example: Exploring the Modes (2)



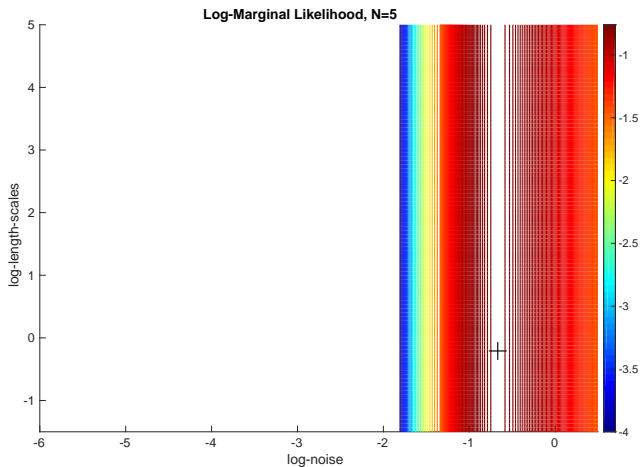
Marginal Likelihood (1)



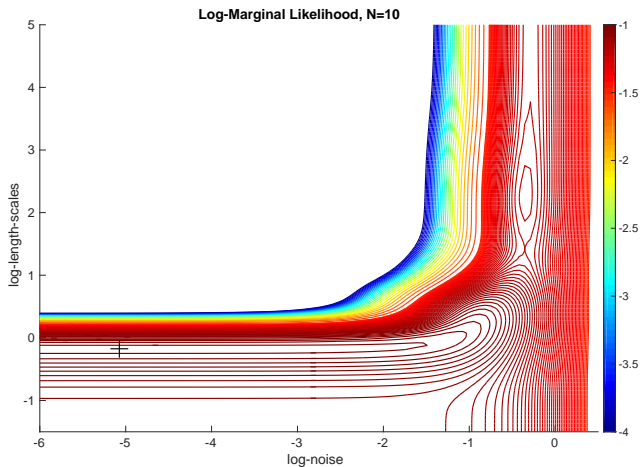
Marginal Likelihood (2)



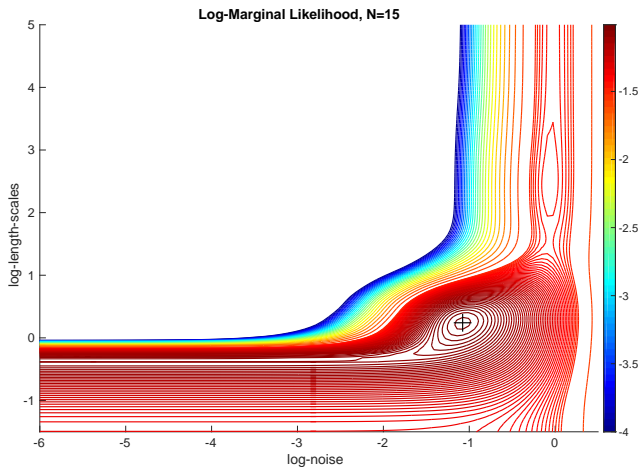
Marginal Likelihood (3)



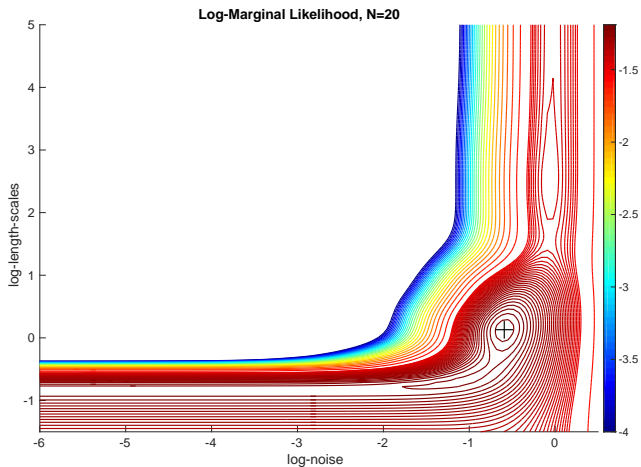
Marginal Likelihood (4)



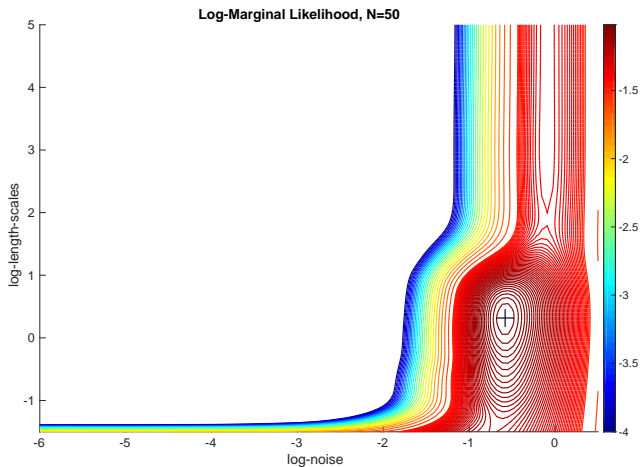
Marginal Likelihood (5)



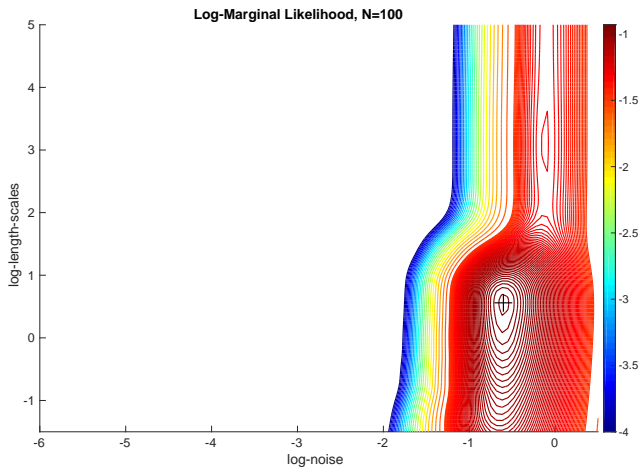
Marginal Likelihood (6)



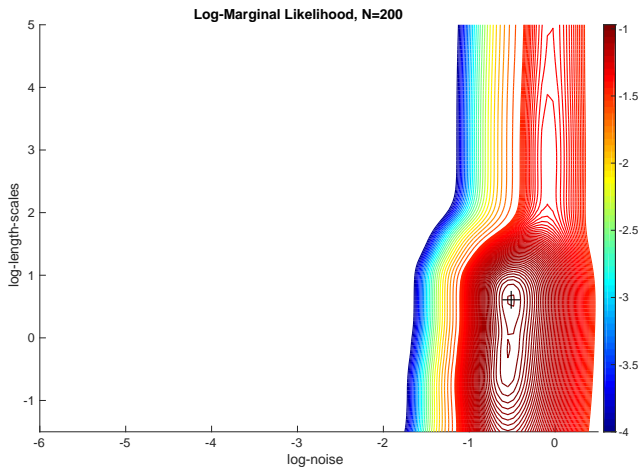
Marginal Likelihood (7)



Marginal Likelihood (8)



Marginal Likelihood (9)



Marginal Likelihood and Parameter Learning

- ▶ The marginal likelihood is non-convex

Marginal Likelihood and Parameter Learning

- ▶ The marginal likelihood is non-convex
- ▶ In particular in the very-small-data regime, a GP can end up in three different modes when optimizing the hyper-parameters:

Marginal Likelihood and Parameter Learning

- ▶ The marginal likelihood is non-convex
- ▶ In particular in the very-small-data regime, a GP can end up in three different modes when optimizing the hyper-parameters:
 - ▶ Overfitting (unlikely, but possible)

Marginal Likelihood and Parameter Learning

- The marginal likelihood is non-convex
- In particular in the very-small-data regime, a GP can end up in three different modes when optimizing the hyper-parameters:
 - Overfitting (unlikely, but possible)
 - Underfitting (everything is considered noise)

Marginal Likelihood and Parameter Learning

- The marginal likelihood is non-convex
- In particular in the very-small-data regime, a GP can end up in three different modes when optimizing the hyper-parameters:
 - Overfitting (unlikely, but possible)
 - Underfitting (everything is considered noise)
 - Good fit

Marginal Likelihood and Parameter Learning

- ▶ The marginal likelihood is non-convex
- ▶ In particular in the very-small-data regime, a GP can end up in three different modes when optimizing the hyper-parameters:
 - ▶ Overfitting (unlikely, but possible)
 - ▶ Underfitting (everything is considered noise)
 - ▶ Good fit
- ▶ Re-start hyper-parameter optimization from random initialization to mitigate the problem

Marginal Likelihood and Parameter Learning

- ▶ The marginal likelihood is non-convex
- ▶ In particular in the very-small-data regime, a GP can end up in three different modes when optimizing the hyper-parameters:
 - ▶ Overfitting (unlikely, but possible)
 - ▶ Underfitting (everything is considered noise)
 - ▶ Good fit
- ▶ Re-start hyper-parameter optimization from random initialization to mitigate the problem
- ▶ With increasing data set size the GP typically ends up in the “good-fit” mode. Overfitting (indicator: small length-scales and small noise variance) is very unlikely.

Marginal Likelihood and Parameter Learning

- ▶ The marginal likelihood is non-convex
- ▶ In particular in the very-small-data regime, a GP can end up in three different modes when optimizing the hyper-parameters:
 - ▶ Overfitting (unlikely, but possible)
 - ▶ Underfitting (everything is considered noise)
 - ▶ Good fit
- ▶ Re-start hyper-parameter optimization from random initialization to mitigate the problem
- ▶ With increasing data set size the GP typically ends up in the “good-fit” mode. Overfitting (indicator: small length-scales and small noise variance) is very unlikely.
- ▶ Ideally, we would integrate the hyper-parameters out
Why can we do not do this easily?

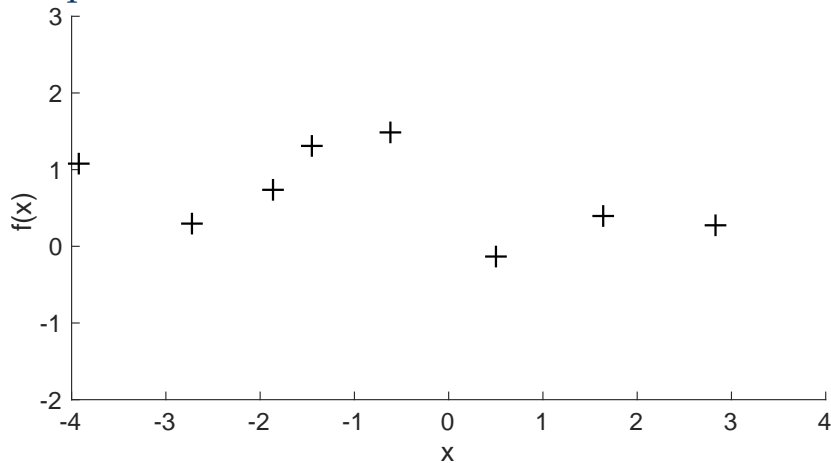
Model Selection—Mean Function and Kernel

- ▶ Assume we have a finite set of models M_i , each one specifying a mean function m_i and a kernel k_i . How do we find the best one?

Model Selection—Mean Function and Kernel

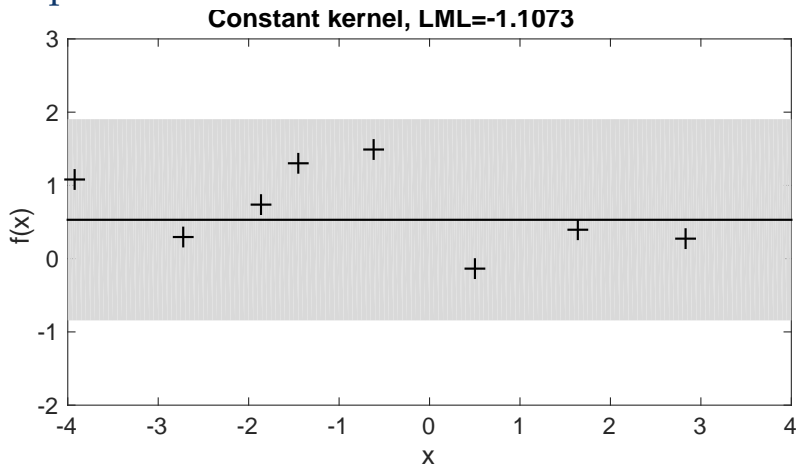
- ▶ Assume we have a finite set of models M_i , each one specifying a mean function m_i and a kernel k_i . How do we find the best one?
- ▶ Some options:
 - ▶ BIC, AIC (see CO-496)
 - ▶ Compare marginal likelihood values (assuming a uniform prior on the set of models)

Example



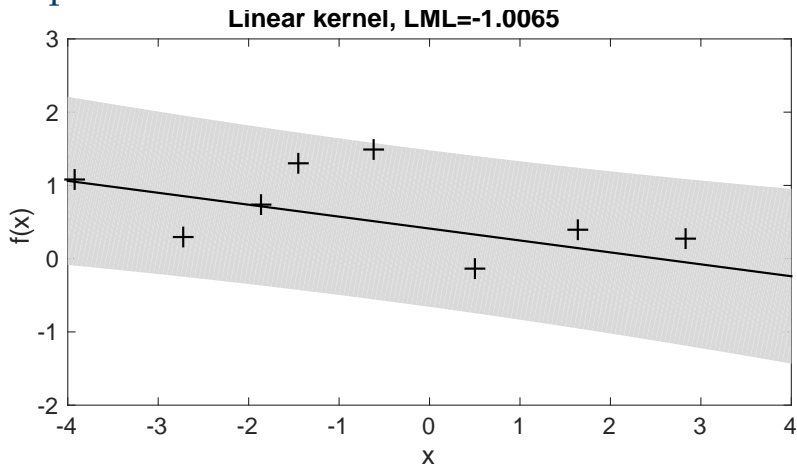
- ▶ Four different kernels (mean function fixed to $m \equiv 0$)
- ▶ MAP hyper-parameters for each kernel
- ▶ Log-marginal likelihood values for each (optimized) model

Example



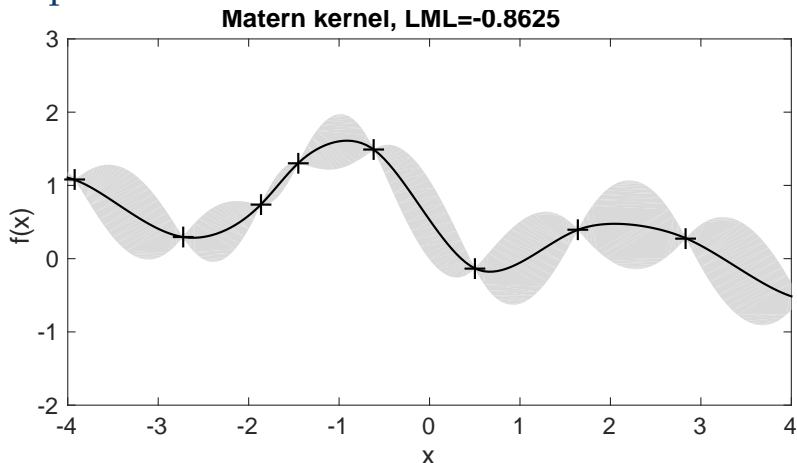
- ▶ Four different kernels (mean function fixed to $m \equiv 0$)
- ▶ MAP hyper-parameters for each kernel
- ▶ Log-marginal likelihood values for each (optimized) model

Example



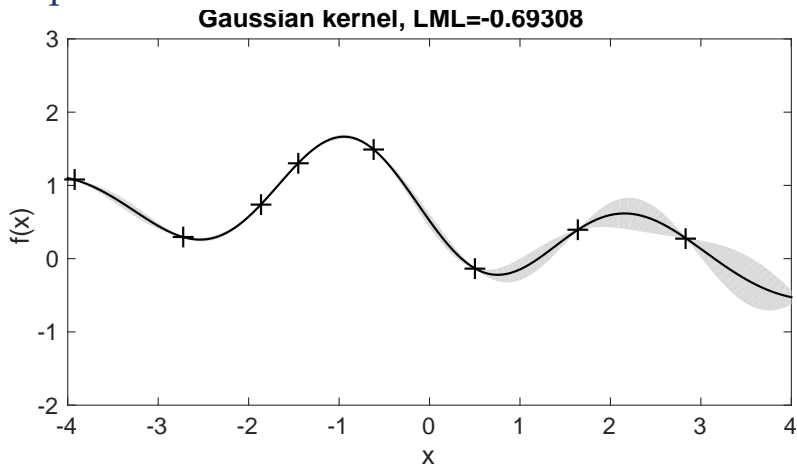
- ▶ Four different kernels (mean function fixed to $m \equiv 0$)
- ▶ MAP hyper-parameters for each kernel
- ▶ Log-marginal likelihood values for each (optimized) model

Example



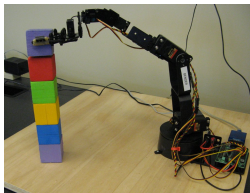
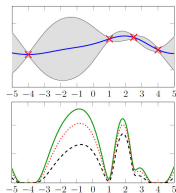
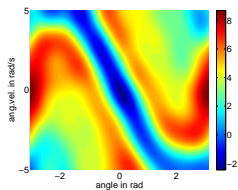
- ▶ Four different kernels (mean function fixed to $m \equiv 0$)
- ▶ MAP hyper-parameters for each kernel
- ▶ Log-marginal likelihood values for each (optimized) model

Example



- ▶ Four different kernels (mean function fixed to $m \equiv 0$)
- ▶ MAP hyper-parameters for each kernel
- ▶ Log-marginal likelihood values for each (optimized) model

Application Areas



- ▶ Reinforcement learning and robotics
 - ▶▶ Model value functions and/or dynamics with GPs
- ▶ Bayesian optimization (Experimental Design)
 - ▶▶ Model unknown utility functions with GPs
- ▶ Geostatistics
 - ▶▶ Spatial modeling (e.g., landscapes, resources)
- ▶ Sensor networks
- ▶ Time-series modeling and forecasting

Limitations of Gaussian Processes

Computational and memory complexity

Training set size: N

- ▶ Training scales in $\mathcal{O}(N^3)$
- ▶ Prediction (variances) scales in $\mathcal{O}(N^2)$
- ▶ Memory requirement: $\mathcal{O}(ND + N^2)$

▶▶ **Practical limit** $N \approx 10,000$

Tips and Tricks for Practitioners

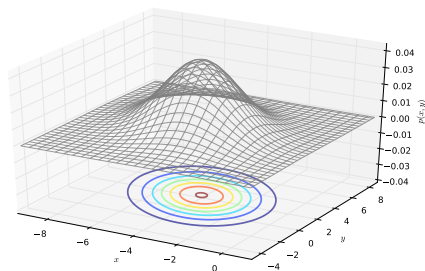
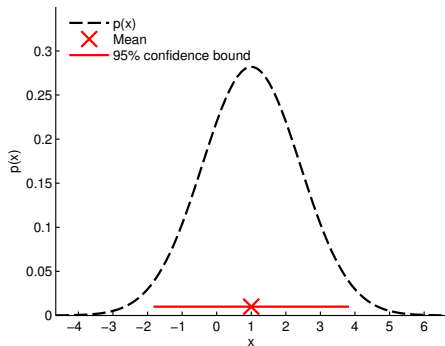
- ▶ To set initial hyper-parameters, use **domain knowledge** if possible.
- ▶ **Standardize** input data and set **initial length-scales** ℓ to ≈ 0.5 .
- ▶ Standardize targets y and set **initial signal variance** to $\sigma_f \approx 1$.
- ▶ Often useful: Set initial noise level relatively high (e.g., $\sigma_n \approx 0.5 \times \sigma_f$ amplitude, even if you think your data have low noise. The optimization surface for your other parameters will be easier to move in.
- ▶ When optimizing hyper-parameters, try **random restarts** or other tricks to avoid local optima are advised.
- ▶ Mitigate the problem of **numerical instability** (Cholesky decomposition of $\mathbf{K} + \sigma_n^2 \mathbf{I}$) by **penalizing high signal-to-noise ratios** σ_f/σ_n

Appendix

The Gaussian Distribution

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

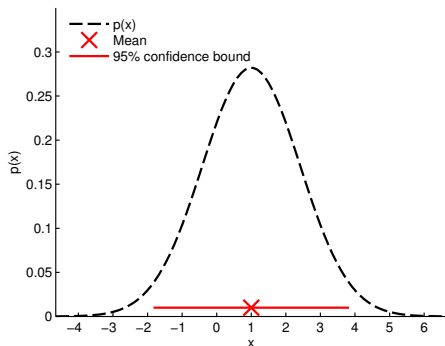
- ▶ Mean vector $\boldsymbol{\mu}$ ▶ Average of the data
- ▶ Covariance matrix $\boldsymbol{\Sigma}$ ▶ Spread of the data



The Gaussian Distribution

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

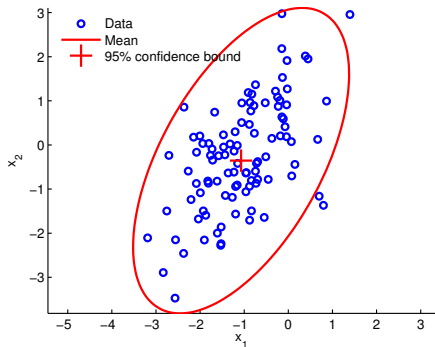
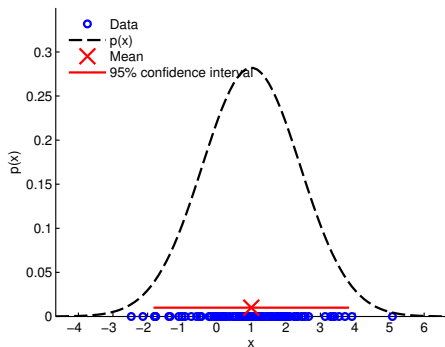
- ▶ Mean vector $\boldsymbol{\mu}$ ▶ Average of the data
- ▶ Covariance matrix $\boldsymbol{\Sigma}$ ▶ Spread of the data



The Gaussian Distribution

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

- ▶ Mean vector $\boldsymbol{\mu}$ ▶ Average of the data
- ▶ Covariance matrix $\boldsymbol{\Sigma}$ ▶ Spread of the data



Sampling from a Multivariate Gaussian

Objective

Generate a random sample $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ from a D -dimensional joint Gaussian with covariance matrix $\boldsymbol{\Sigma}$ and mean vector $\boldsymbol{\mu}$.

However, we only have access to a random number generator that can sample \mathbf{x} from $\mathcal{N}(\mathbf{0}, \mathbf{I})$...

Sampling from a Multivariate Gaussian

Objective

Generate a random sample $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ from a D -dimensional joint Gaussian with covariance matrix $\boldsymbol{\Sigma}$ and mean vector $\boldsymbol{\mu}$.

However, we only have access to a random number generator that can sample \mathbf{x} from $\mathcal{N}(\mathbf{0}, \mathbf{I})$...

Exploit that affine transformations $\mathbf{y} = \mathbf{Ax} + \mathbf{b}$ of a Gaussian random variable \mathbf{x} remain Gaussian

- ▶ Mean: $\mathbb{E}_x[\mathbf{Ax} + \mathbf{b}] = \mathbf{A}\mathbb{E}_x[\mathbf{x}] + \mathbf{b}$
- ▶ Covariance: $\mathbb{V}_x[\mathbf{Ax} + \mathbf{b}] = \mathbf{A}\mathbb{V}_x[\mathbf{x}]\mathbf{A}^\top$

Sampling from a Multivariate Gaussian

Objective

Generate a random sample $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ from a D -dimensional joint Gaussian with covariance matrix $\boldsymbol{\Sigma}$ and mean vector $\boldsymbol{\mu}$.

However, we only have access to a random number generator that can sample \mathbf{x} from $\mathcal{N}(\mathbf{0}, \mathbf{I})$...

Exploit that affine transformations $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ of a Gaussian random variable \mathbf{x} remain Gaussian

- ▶ Mean: $\mathbb{E}_x[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{E}_x[\mathbf{x}] + \mathbf{b}$
 - ▶ Covariance: $\mathbb{V}_x[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{V}_x[\mathbf{x}]\mathbf{A}^\top$
1. Find conditions for \mathbf{A}, \mathbf{b} to match the mean of \mathbf{y}
 2. Find conditions for \mathbf{A}, \mathbf{b} to match the covariance of \mathbf{y}

Sampling from a Multivariate Gaussian (2)

Objective

Generate a random sample $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ from a D -dimensional joint Gaussian with covariance matrix $\boldsymbol{\Sigma}$ and mean vector $\boldsymbol{\mu}$.

$\mathbf{x} = \text{randn}(D, 1);$ Sample $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 $\mathbf{y} = \text{chol}(\boldsymbol{\Sigma})' * \mathbf{x} + \boldsymbol{\mu};$ Scale \mathbf{x} and add offset

Here $\text{chol}(\boldsymbol{\Sigma})$ is the Cholesky factor \mathbf{L} , such that $\mathbf{L}^\top \mathbf{L} = \boldsymbol{\Sigma}$

Sampling from a Multivariate Gaussian (2)

Objective

Generate a random sample $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ from a D -dimensional joint Gaussian with covariance matrix $\boldsymbol{\Sigma}$ and mean vector $\boldsymbol{\mu}$.

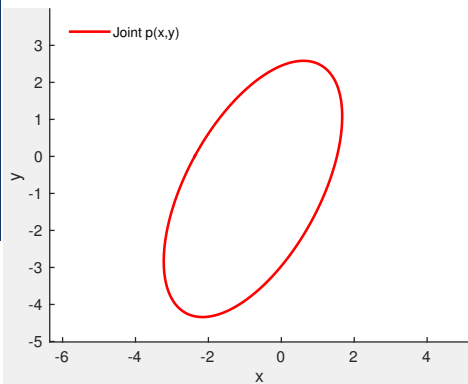
$\mathbf{x} = \text{randn}(D, 1);$ Sample $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 $\mathbf{y} = \text{chol}(\boldsymbol{\Sigma})' * \mathbf{x} + \boldsymbol{\mu};$ Scale \mathbf{x} and add offset

Here $\text{chol}(\boldsymbol{\Sigma})$ is the Cholesky factor \mathbf{L} , such that $\mathbf{L}^\top \mathbf{L} = \boldsymbol{\Sigma}$
Therefore, the mean and covariance of \mathbf{y} are

$$\mathbb{E}[\mathbf{y}] = \bar{\mathbf{y}} = \mathbb{E}[\mathbf{L}^\top \mathbf{x} + \boldsymbol{\mu}] = \mathbf{L}^\top \mathbb{E}[\mathbf{x}] + \boldsymbol{\mu} = \boldsymbol{\mu}$$

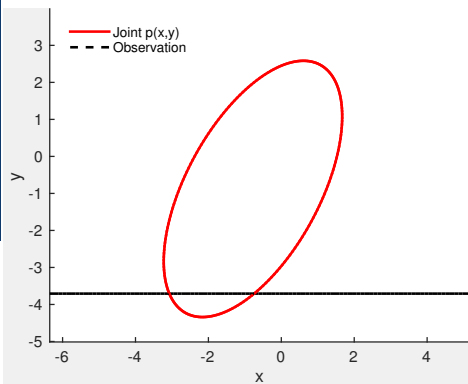
$$\text{Cov}[\mathbf{y}] = \mathbb{E}[(\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^\top] = \mathbb{E}[\mathbf{L}^\top \mathbf{x} \mathbf{x}^\top \mathbf{L}] = \mathbf{L}^\top \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \mathbf{L} = \mathbf{L}^\top \mathbf{L} = \boldsymbol{\Sigma}$$

Conditional



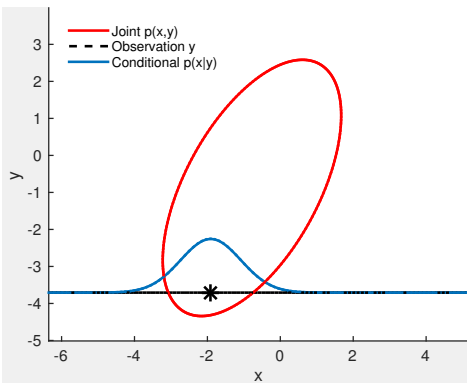
$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

Conditional



$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

Conditional



$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x|\mathbf{y}}, \boldsymbol{\Sigma}_{x|\mathbf{y}})$$

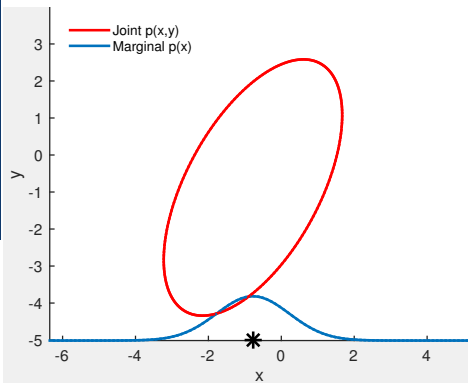
$$\boldsymbol{\mu}_{x|\mathbf{y}} = \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} (\mathbf{y} - \boldsymbol{\mu}_y)$$

$$\boldsymbol{\Sigma}_{x|\mathbf{y}} = \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx}$$

Conditional $p(\mathbf{x}|\mathbf{y})$ is also Gaussian

► Computationally convenient

Marginal

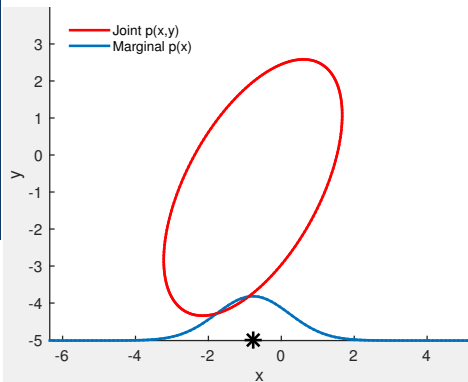


$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

Marginal distribution:

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \mathcal{N}(\mu_x, \Sigma_{xx}) \end{aligned}$$

Marginal



$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

Marginal distribution:

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \mathcal{N}(\mu_x, \Sigma_{xx}) \end{aligned}$$

- ▶ The marginal of a joint Gaussian distribution is Gaussian
- ▶ Intuitively: Ignore (integrate out) everything you are not interested in

The Gaussian Distribution in the Limit

Consider the **joint Gaussian distribution** $p(\mathbf{x}, \tilde{\mathbf{x}})$, where $\mathbf{x} \in \mathbb{R}^D$ and $\tilde{\mathbf{x}} \in \mathbb{R}^k, k \rightarrow \infty$ are random variables.

The Gaussian Distribution in the Limit

Consider the **joint Gaussian distribution** $p(\mathbf{x}, \tilde{\mathbf{x}})$, where $\mathbf{x} \in \mathbb{R}^D$ and $\tilde{\mathbf{x}} \in \mathbb{R}^k, k \rightarrow \infty$ are random variables.

Then

$$p(\mathbf{x}, \tilde{\mathbf{x}}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_{\tilde{\mathbf{x}}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{x\tilde{\mathbf{x}}} \\ \boldsymbol{\Sigma}_{\tilde{\mathbf{x}}x} & \boldsymbol{\Sigma}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} \end{bmatrix} \right)$$

where $\boldsymbol{\Sigma}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} \in \mathbb{R}^{k \times k}$ and $\boldsymbol{\Sigma}_{x\tilde{\mathbf{x}}} \in \mathbb{R}^{D \times k}, k \rightarrow \infty$.

The Gaussian Distribution in the Limit

Consider the **joint Gaussian distribution** $p(\mathbf{x}, \tilde{\mathbf{x}})$, where $\mathbf{x} \in \mathbb{R}^D$ and $\tilde{\mathbf{x}} \in \mathbb{R}^k, k \rightarrow \infty$ are random variables.

Then

$$p(\mathbf{x}, \tilde{\mathbf{x}}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_{\mathbf{x}} \\ \boldsymbol{\mu}_{\tilde{\mathbf{x}}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\mathbf{x}\mathbf{x}} & \boldsymbol{\Sigma}_{\mathbf{x}\tilde{\mathbf{x}}} \\ \boldsymbol{\Sigma}_{\tilde{\mathbf{x}}\mathbf{x}} & \boldsymbol{\Sigma}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} \end{bmatrix} \right)$$

where $\boldsymbol{\Sigma}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} \in \mathbb{R}^{k \times k}$ and $\boldsymbol{\Sigma}_{\mathbf{x}\tilde{\mathbf{x}}} \in \mathbb{R}^{D \times k}, k \rightarrow \infty$.

However, the **marginal remains finite**

$$p(\mathbf{x}) = \int p(\mathbf{x}, \tilde{\mathbf{x}}) d\tilde{\mathbf{x}} = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}\mathbf{x}})$$

where we integrate out an infinite number of random variables \tilde{x}_i .

Marginal and Conditional in the Limit

- ▶ In practice, we consider **finite training and test data** $\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}$

Marginal and Conditional in the Limit

- ▶ In practice, we consider **finite training and test data** $\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}$
- ▶ Then, $\mathbf{x} = \{\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}, \mathbf{x}_{\text{other}}\}$
($\mathbf{x}_{\text{other}}$ plays the role of $\tilde{\mathbf{x}}$ from previous slide)

Marginal and Conditional in the Limit

- ▶ In practice, we consider **finite training and test data** $\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}$
- ▶ Then, $\mathbf{x} = \{\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}, \mathbf{x}_{\text{other}}\}$
($\mathbf{x}_{\text{other}}$ plays the role of $\tilde{\mathbf{x}}$ from previous slide)

$$p(\mathbf{x}) = \mathcal{N} \left(\begin{bmatrix} \mu_{\text{train}} \\ \mu_{\text{test}} \\ \mu_{\text{other}} \end{bmatrix}, \begin{bmatrix} \Sigma_{\text{train}} & \Sigma_{\text{train,test}} & \Sigma_{\text{train,other}} \\ \Sigma_{\text{test,train}} & \Sigma_{\text{test}} & \Sigma_{\text{test,other}} \\ \Sigma_{\text{other,train}} & \Sigma_{\text{other,test}} & \Sigma_{\text{other}} \end{bmatrix} \right)$$

Marginal and Conditional in the Limit

- ▶ In practice, we consider **finite training and test data** $\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}$
- ▶ Then, $\mathbf{x} = \{\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}, \mathbf{x}_{\text{other}}\}$
($\mathbf{x}_{\text{other}}$ plays the role of $\tilde{\mathbf{x}}$ from previous slide)

$$p(\mathbf{x}) = \mathcal{N} \left(\begin{bmatrix} \mu_{\text{train}} \\ \mu_{\text{test}} \\ \mu_{\text{other}} \end{bmatrix}, \begin{bmatrix} \Sigma_{\text{train}} & \Sigma_{\text{train,test}} & \Sigma_{\text{train,other}} \\ \Sigma_{\text{test,train}} & \Sigma_{\text{test}} & \Sigma_{\text{test,other}} \\ \Sigma_{\text{other,train}} & \Sigma_{\text{other,test}} & \Sigma_{\text{other}} \end{bmatrix} \right)$$

$$p(\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}) = \int p(\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}, \mathbf{x}_{\text{other}}) d\mathbf{x}_{\text{other}}$$

Marginal and Conditional in the Limit

- ▶ In practice, we consider **finite training and test data** $\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}$
- ▶ Then, $\mathbf{x} = \{\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}, \mathbf{x}_{\text{other}}\}$
($\mathbf{x}_{\text{other}}$ plays the role of $\tilde{\mathbf{x}}$ from previous slide)

$$p(\mathbf{x}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_{\text{train}} \\ \boldsymbol{\mu}_{\text{test}} \\ \boldsymbol{\mu}_{\text{other}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\text{train}} & \boldsymbol{\Sigma}_{\text{train, test}} & \boldsymbol{\Sigma}_{\text{train, other}} \\ \boldsymbol{\Sigma}_{\text{test, train}} & \boldsymbol{\Sigma}_{\text{test}} & \boldsymbol{\Sigma}_{\text{test, other}} \\ \boldsymbol{\Sigma}_{\text{other, train}} & \boldsymbol{\Sigma}_{\text{other, test}} & \boldsymbol{\Sigma}_{\text{other}} \end{bmatrix} \right)$$

$$p(\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}) = \int p(\mathbf{x}_{\text{train}}, \mathbf{x}_{\text{test}}, \mathbf{x}_{\text{other}}) d\mathbf{x}_{\text{other}}$$

$$p(\mathbf{x}_{\text{test}} | \mathbf{x}_{\text{train}}) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\boldsymbol{\mu}_* = \boldsymbol{\mu}_{\text{test}} + \boldsymbol{\Sigma}_{\text{test, train}} \boldsymbol{\Sigma}_{\text{train}}^{-1} (\mathbf{x}_{\text{train}} - \boldsymbol{\mu}_{\text{train}})$$

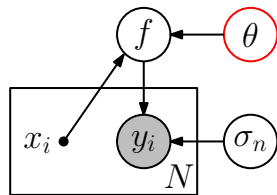
$$\boldsymbol{\Sigma}_* = \boldsymbol{\Sigma}_{\text{test}} - \boldsymbol{\Sigma}_{\text{test, train}} \boldsymbol{\Sigma}_{\text{train}}^{-1} \boldsymbol{\Sigma}_{\text{train, test}}$$

Gaussian Process Training: Hierarchical Inference

- Level-1 inference (posterior on f):

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

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



Gaussian Process Training: Hierarchical Inference

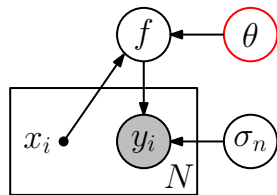
- ▶ Level-1 inference (posterior on f):

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

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

- ▶ Level-2 inference (posterior on $\boldsymbol{\theta}$)

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



GP as the Limit of an Infinite RBF Network

Consider the universal function approximator

$$f(x) = \sum_{i \in \mathbb{Z}} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \gamma_n \exp \left(-\frac{(x - (i + \frac{n}{N}))^2}{\lambda^2} \right), \quad x \in \mathbb{R}, \quad \lambda \in \mathbb{R}^+$$

with $\gamma_n \sim \mathcal{N}(0, 1)$ (random weights)

► Gaussian-shaped basis functions (with variance $\lambda^2/2$) everywhere on the real axis

GP as the Limit of an Infinite RBF Network

Consider the universal function approximator

$$f(x) = \sum_{i \in \mathbb{Z}} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \gamma_n \exp \left(-\frac{(x - (i + \frac{n}{N}))^2}{\lambda^2} \right), \quad x \in \mathbb{R}, \quad \lambda \in \mathbb{R}^+$$

with $\gamma_n \sim \mathcal{N}(0, 1)$ (random weights)

► Gaussian-shaped basis functions (with variance $\lambda^2/2$) everywhere on the real axis

$$f(x) = \sum_{i \in \mathbb{Z}} \int_i^{i+1} \gamma(s) \exp \left(-\frac{(x-s)^2}{\lambda^2} \right) ds = \int_{-\infty}^{\infty} \gamma(s) \exp \left(-\frac{(x-s)^2}{\lambda^2} \right) ds$$

GP as the Limit of an Infinite RBF Network

Consider the universal function approximator

$$f(x) = \sum_{i \in \mathbb{Z}} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \gamma_n \exp \left(-\frac{(x - (i + \frac{n}{N}))^2}{\lambda^2} \right), \quad x \in \mathbb{R}, \quad \lambda \in \mathbb{R}^+$$

with $\gamma_n \sim \mathcal{N}(0, 1)$ (random weights)

► Gaussian-shaped basis functions (with variance $\lambda^2/2$) everywhere on the real axis

$$f(x) = \sum_{i \in \mathbb{Z}} \int_i^{i+1} \gamma(s) \exp \left(-\frac{(x-s)^2}{\lambda^2} \right) ds = \int_{-\infty}^{\infty} \gamma(s) \exp \left(-\frac{(x-s)^2}{\lambda^2} \right) ds$$

► Mean: $\mathbb{E}[f(x)] = 0$

► Covariance: $\text{Cov}[f(x), f(x')] = \theta_1^2 \exp \left(-\frac{(x-x')^2}{2\lambda^2} \right)$ for suitable θ_1^2

► GP with mean 0 and Gaussian covariance function

References I

- [1] N. A. C. Cressie. *Statistics for Spatial Data*. Wiley-Interscience, 1993.
- [2] M. P. Deisenroth and S. Mohamed. Expectation Propagation in Gaussian Process Dynamical Systems. In *Advances in Neural Information Processing Systems*, pages 2618–2626, 2012.
- [3] M. P. Deisenroth and J. W. Ng. Distributed Gaussian Processes. In *Proceedings of the International Conference on Machine Learning*, 2015.
- [4] M. P. Deisenroth, C. E. Rasmussen, and J. Peters. Gaussian Process Dynamic Programming. *Neurocomputing*, 72(7–9):1508–1524, March 2009.
- [5] 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.
- [6] 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 C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, pages 3156–3164. Curran Associates, Inc., 2013.
- [7] J. Kocijan, R. Murray-Smith, C. E. Rasmussen, and A. Girard. Gaussian Process Model Based Predictive Control. In *Proceedings of the 2004 American Control Conference (ACC 2004)*, pages 2214–2219, Boston, MA, USA, June–July 2004.
- [8] 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, February 2008.
- [9] 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.
- [10] 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.
- [11] J. Quiñero-Candela and C. E. Rasmussen. A Unifying View of Sparse Approximate Gaussian Process Regression. *Journal of Machine Learning Research*, 6(2):1939–1960, 2005.
- [12] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA, 2006.
- [13] S. Roberts, M. A. Osborne, M. Ebdon, S. Reece, N. Gibson, and S. Aigrain. Gaussian Processes for Time Series Modelling. *Philosophical Transactions of the Royal Society (Part A)*, 371(1984), February 2013.