

Gaussian Processes for Big Data Problems

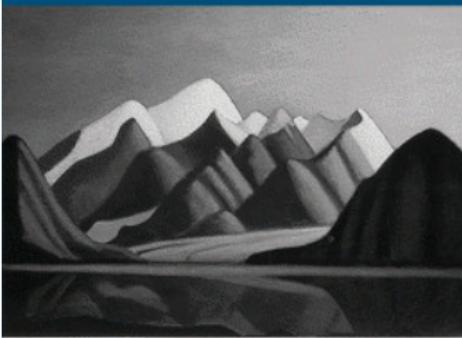
Marc Deisenroth

Department of Computing
Imperial College London

<http://wp.doc.ic.ac.uk/sml/marc-deisenroth>

Machine Learning Summer School, Chalmers University
14 April 2015

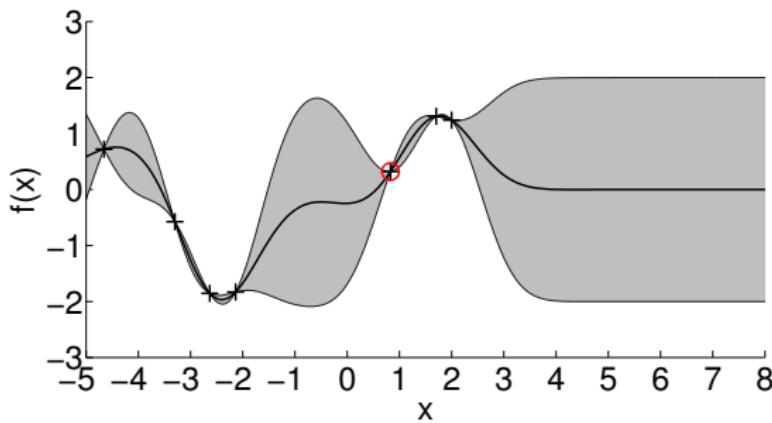
Gaussian Processes for Machine Learning



Carl Edward Rasmussen and Christopher K. I. Williams

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

Problem Setting

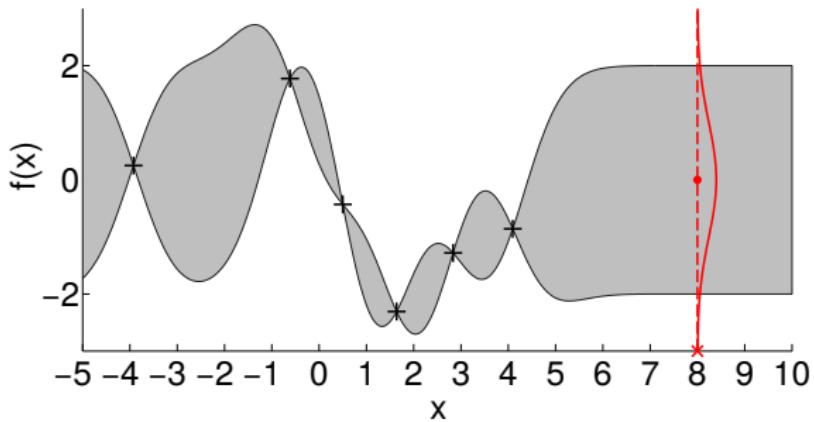


Objective

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

► Probabilistic regression problem

Problem Setting

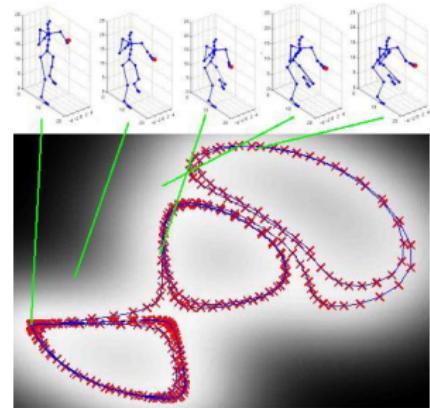
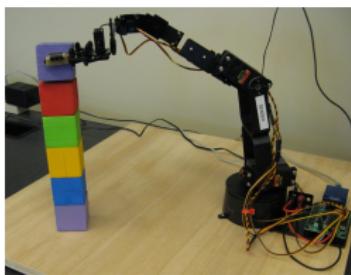
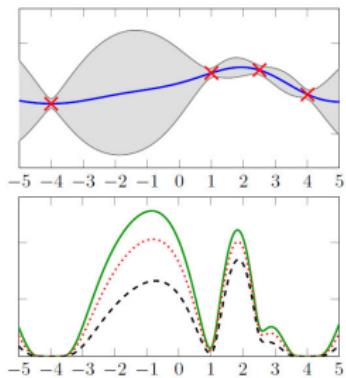


Objective

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

(Some) Relevant Application Areas



- Global black-box optimization and experimental design
- Autonomous learning in robotics
- Probabilistic dimensionality reduction and data visualization

Table of Contents

Introduction

 Gaussian Distribution

Gaussian Processes

 Definition and Derivation

 Inference

 Covariance Functions

 Model Selection

 Limitations

Scaling Gaussian Processes to Large Data Sets

 Sparse Gaussian Processes

 Distributed Gaussian Processes

Table of Contents

Introduction

Gaussian Distribution

Gaussian Processes

Definition and Derivation

Inference

Covariance Functions

Model Selection

Limitations

Scaling Gaussian Processes to Large Data Sets

Sparse Gaussian Processes

Distributed Gaussian Processes

Key Concepts in Probability Theory

Two fundamental rules:

$$p(x) = \int p(x,y)dy \quad \text{Sum rule/Marginalization property}$$

$$p(x,y) = p(y|x)p(x) \quad \text{Product rule}$$

Key Concepts in Probability Theory

Two fundamental rules:

$$p(x) = \int p(x,y)dy \quad \text{Sum rule/Marginalization property}$$

$$p(x,y) = p(y|x)p(x) \quad \text{Product rule}$$

Bayes' Theorem (Probabilistic Inverse)

$$p(x|y) = \frac{p(y|x) p(x)}{p(y)}, \quad x : \text{hypothesis}, \quad y : \text{measurement}$$

Key Concepts in Probability Theory

Two fundamental rules:

$$p(x) = \int p(x,y)dy \quad \text{Sum rule/Marginalization property}$$

$$p(x,y) = p(y|x)p(x) \quad \text{Product rule}$$

Bayes' Theorem (Probabilistic Inverse)

$$p(x|y) = \frac{p(y|x) p(x)}{p(y)}, \quad x : \text{hypothesis}, \quad y : \text{measurement}$$

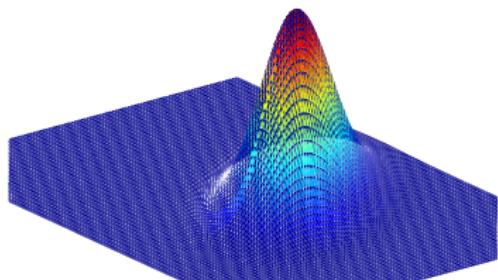
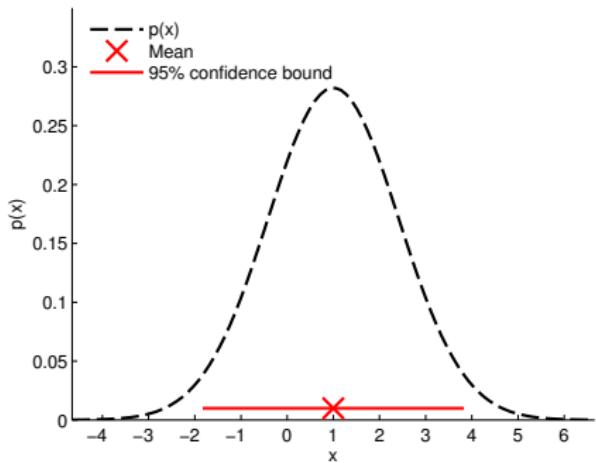
- ▶ Posterior belief
- ▶ Prior belief
- ▶ Likelihood (measurement model)
- ▶ Marginal likelihood (normalization constant)



The Gaussian Distribution

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

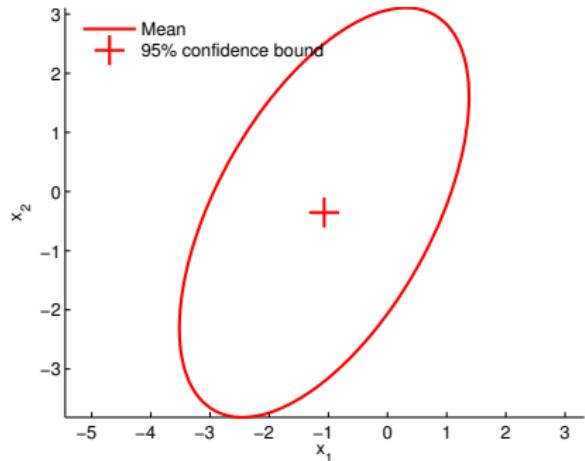
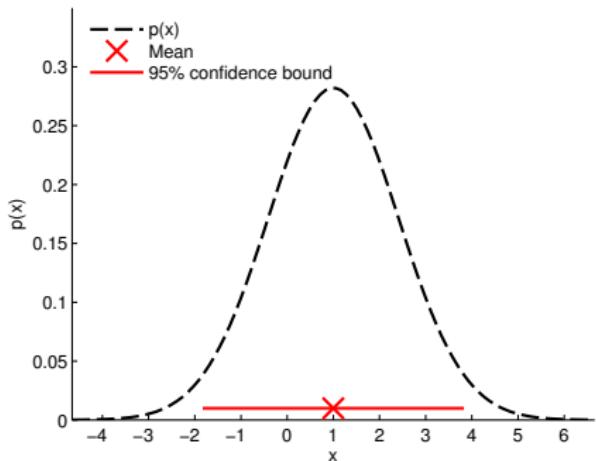
- ▶ Mean vector μ ► Average of the data
- ▶ Covariance matrix Σ ► Spread of the data



The Gaussian Distribution

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

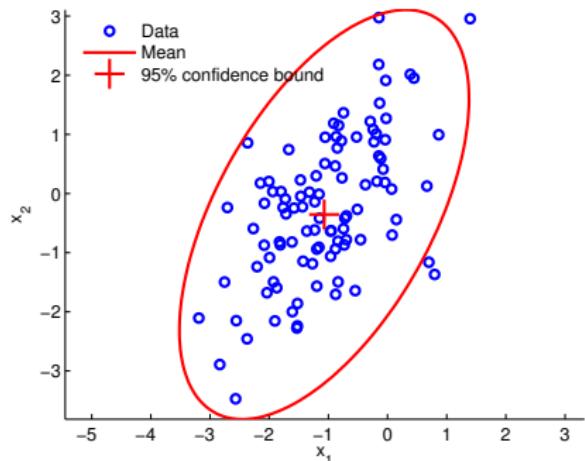
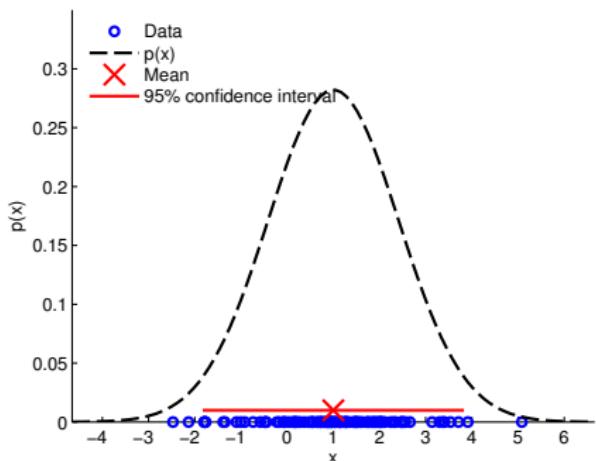
- ▶ Mean vector μ ➡ Average of the data
- ▶ Covariance matrix Σ ➡ Spread of the data



The Gaussian Distribution

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

- ▶ Mean vector μ ➡ Average of the data
- ▶ Covariance matrix Σ ➡ Spread of the data



Sampling from a Multivariate Gaussian

Objective

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

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

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 $\mathbf{y} = A\mathbf{x} + \mathbf{b}$ of Gaussians remain Gaussian

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

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 Gaussians remain Gaussian

- Mean: $\mathbb{E}_{\mathbf{x}}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{E}_{\mathbf{x}}[\mathbf{x}] + \mathbf{b}$
 - Covariance: $\mathbb{V}_{\mathbf{x}}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{V}_{\mathbf{x}}[\mathbf{x}]\mathbf{A}^{\top}$
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 $y \sim \mathcal{N}(\mu, \Sigma)$ from a D -dimensional joint Gaussian with covariance matrix Σ and mean vector μ .

```
x = randn(D, 1);           Sample x ~ N(0, I)
y = chol(Sigma)' * x + mu; Scale x
```

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

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 $\boldsymbol{\mu}$.

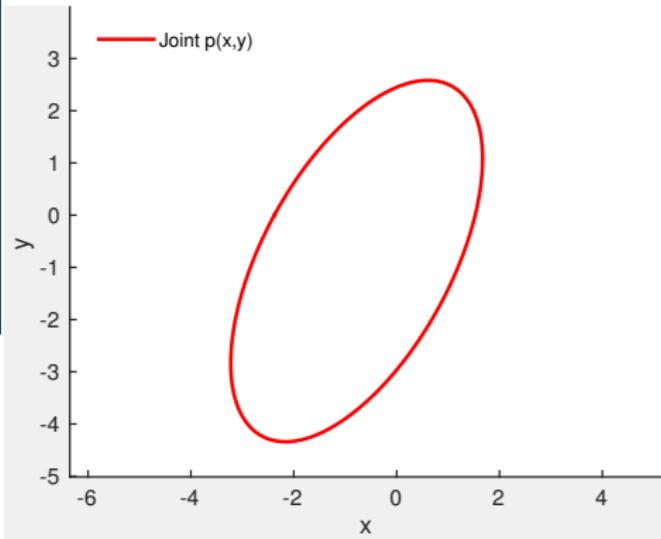
```
x = randn(D, 1);           Sample  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
y = chol(Sigma)' * x + mu; Scale  $\mathbf{x}$ 
```

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

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

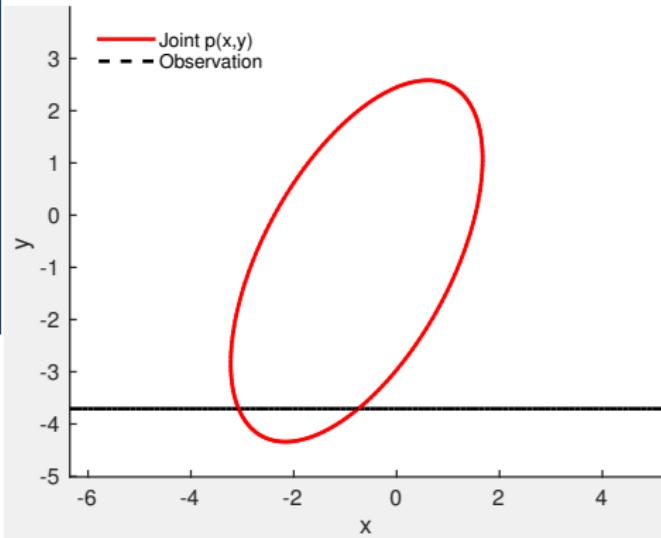
$$\text{Cov}[\mathbf{y}] = \mathbb{E}[(\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^\top] = \mathbb{E}[\mathbf{L}^\top \mathbf{x} \mathbf{x}^\top \mathbf{L}] = \mathbf{L}^\top \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \mathbf{L} = \mathbf{L}^\top \boldsymbol{\Sigma} \mathbf{L} = \boldsymbol{\Sigma}$$

Conditional



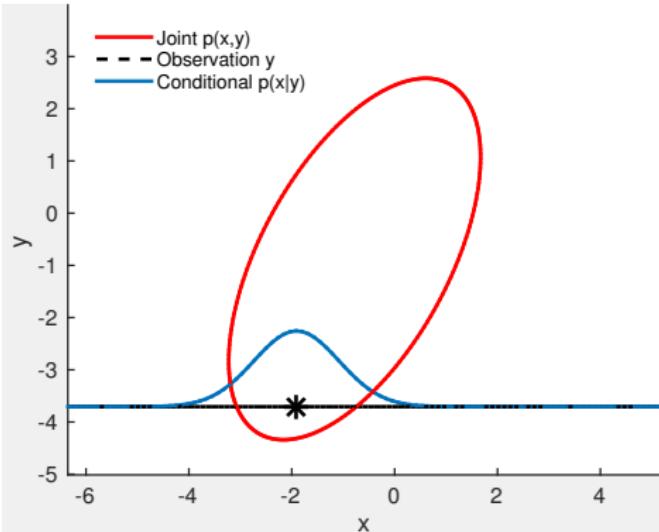
$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

Conditional



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

Conditional



Conditional $p(x|y)$ is also Gaussian
► Computationally convenient

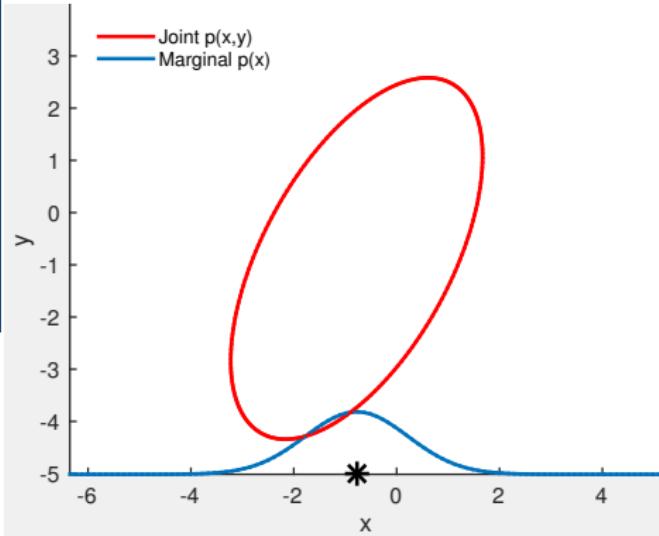
$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

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

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

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

Marginal



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

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}) \end{aligned}$$

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

Table of Contents

Introduction

 Gaussian Distribution

Gaussian Processes

 Definition and Derivation

 Inference

 Covariance Functions

 Model Selection

 Limitations

Scaling Gaussian Processes to Large Data Sets

 Sparse Gaussian Processes

 Distributed Gaussian Processes

Gaussian Process

Definition

A **Gaussian process** (GP) is a collection of random variables x_1, x_2, \dots , any finite number of which is Gaussian distributed.

Gaussian Process

Definition

A **Gaussian process** (GP) is a collection of random variables x_1, x_2, \dots , any finite number of which is Gaussian distributed.

Unexpected (?) example: Linear dynamical system

$$x_{t+1} = Ax_t + w, \quad w \sim \mathcal{N}(0, Q)$$

x_1, x_2, \dots is a collection of random variables

Gaussian Process

Definition

A **Gaussian process** (GP) is a collection of random variables x_1, x_2, \dots , any finite number of which is Gaussian distributed.

Unexpected (?) example: Linear dynamical system

$$x_{t+1} = Ax_t + w, \quad w \sim \mathcal{N}(0, Q)$$

x_1, x_2, \dots is a collection of random variables,
any finite number of which is Gaussian distributed

Gaussian Process

Definition

A **Gaussian process** (GP) is a collection of random variables x_1, x_2, \dots , any finite number of which is Gaussian distributed.

Unexpected (?) example: Linear dynamical system

$$x_{t+1} = Ax_t + w, \quad w \sim \mathcal{N}(0, Q)$$

x_1, x_2, \dots is a collection of random variables,
any finite number of which is Gaussian distributed
► This is a GP

Gaussian Process

Definition

A **Gaussian process** (GP) is a collection of random variables x_1, x_2, \dots , any finite number of which is Gaussian distributed.

Unexpected (?) example: Linear dynamical system

$$x_{t+1} = Ax_t + w, \quad w \sim \mathcal{N}(0, Q)$$

x_1, x_2, \dots is a collection of random variables,
any finite number of which is Gaussian distributed
► This is a GP

But we will use the Gaussian process for a different purpose, not for time series

The Gaussian in the Limit

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

The Gaussian in the Limit

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

Then

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

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

The Gaussian in the Limit

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

Then

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

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

However, the marginal remains finite

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

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

Marginal and Conditional in the Limit

- In practice, we consider finite training and test data $x_{\text{train}}, x_{\text{test}}$

Marginal and Conditional in the Limit

- In practice, we consider finite training and test data $x_{\text{train}}, x_{\text{test}}$
- Then, $x = \{x_{\text{train}}, x_{\text{test}}, x_{\text{other}}\}$
(x_{other} plays the role of \tilde{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)

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

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)

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

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

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)

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

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

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

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

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

Back to Regression: Distribution over Functions

- We are not really interested in a distribution on x , but rather in a distribution over function values $f(x)$

Back to Regression: Distribution over Functions

- We are not really interested in a distribution on x , but rather in a distribution over function values $f(x)$
- Let's replace x with function values $f = f(x)$
 - ▶ (Treat a function as a long vector of function values)

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

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

Back to Regression: Distribution over Functions

- We are not really interested in a distribution on x , but rather in a distribution over function values $f(x)$
- Let's replace x with function values $f = f(x)$
 - ▶ (Treat a function as a long vector of function values)

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

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

- Again, the marginal remains finite

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

Marginal and Conditional over Functions

Define $f_* := f_{\text{test}}$, $f := f_{\text{train}}$.

Marginal and Conditional over Functions

Define $f_* := f_{\text{test}}$, $f := f_{\text{train}}$.

- Marginal

$$p(f_*, f) = \mathcal{N} \left(\begin{bmatrix} \mu_f \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{f*} \\ \Sigma_{*f} & \Sigma_{**} \end{bmatrix} \right)$$

Marginal and Conditional over Functions

Define $f_* := f_{\text{test}}$, $f := f_{\text{train}}$.

- Marginal

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

- Conditional (predictive distribution)

$$p(f_* | f) = \mathcal{N}(\mathbf{m}, \mathbf{S})$$

$$\mathbf{m} = \boldsymbol{\mu}_* + \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} (f - \boldsymbol{\mu})$$

$$\mathbf{S} = \boldsymbol{\Sigma}_{**} - \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} \boldsymbol{\Sigma}_{f*}$$

Marginal and Conditional over Functions

Define $f_* := f_{\text{test}}$, $f := f_{\text{train}}$.

- Marginal

$$p(f_*, f) = \mathcal{N} \left(\begin{bmatrix} \mu_f \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{f*} \\ \Sigma_{*f} & \Sigma_{**} \end{bmatrix} \right)$$

- Conditional (predictive distribution)

$$p(f_* | f) = \mathcal{N}(\mathbf{m}, \mathbf{S})$$

$$\mathbf{m} = \mu_* + \Sigma_{*f} \Sigma_{ff}^{-1} (f - \mu)$$

$$\mathbf{S} = \Sigma_{**} - \Sigma_{*f} \Sigma_{ff}^{-1} \Sigma_{f*}$$

- We need to compute (cross-)covariances between **unknown** function values

Marginal and Conditional over Functions

Define $f_* := f_{\text{test}}$, $f := f_{\text{train}}$.

- Marginal

$$p(f_*, f) = \mathcal{N} \left(\begin{bmatrix} \mu_f \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{f*} \\ \Sigma_{*f} & \Sigma_{**} \end{bmatrix} \right)$$

- Conditional (predictive distribution)

$$p(f_* | f) = \mathcal{N}(\mathbf{m}, \mathbf{S})$$

$$\mathbf{m} = \mu_* + \Sigma_{*f} \Sigma_{ff}^{-1} (f - \mu)$$

$$\mathbf{S} = \Sigma_{**} - \Sigma_{*f} \Sigma_{ff}^{-1} \Sigma_{f*}$$

- ▶ We need to compute (cross-)covariances between **unknown** function values
- ▶ The **kernel trick** helps us out

Kernelization

$$p(f_*, f) = \mathcal{N} \left(\begin{bmatrix} \mu_f \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{f*} \\ \Sigma_{*f} & \Sigma_{**} \end{bmatrix} \right)$$

$$p(f_* | f) = \mathcal{N}(\boldsymbol{m}, \boldsymbol{S})$$

$$\boldsymbol{m} = \boldsymbol{\mu}_* + \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} (\boldsymbol{f} - \boldsymbol{\mu})$$

$$\boldsymbol{S} = \boldsymbol{\Sigma}_{**} - \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} \boldsymbol{\Sigma}_{f*}$$

Kernelization

$$p(f_* | f) = \mathcal{N} \left(\begin{bmatrix} \mu_f \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{f*} \\ \Sigma_{*f} & \Sigma_{**} \end{bmatrix} \right)$$
$$p(f_* | f) = \mathcal{N}(f_* | \mathbf{m}, S)$$
$$\mathbf{m} = \boldsymbol{\mu}_* + \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} (f - \boldsymbol{\mu})$$
$$S = \boldsymbol{\Sigma}_{**} - \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} \boldsymbol{\Sigma}_{f*}$$

- A **kernel function** k is symmetric and positive definite.

Kernelization

$$p(f_* | f) = \mathcal{N} \left(\begin{bmatrix} \mu_f \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{f*} \\ \Sigma_{*f} & \Sigma_{**} \end{bmatrix} \right)$$
$$p(f_* | f) = \mathcal{N}(f_* | \boldsymbol{m}, S)$$
$$\boldsymbol{m} = \boldsymbol{\mu}_* + \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} (f - \boldsymbol{\mu})$$
$$S = \boldsymbol{\Sigma}_{**} - \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} \boldsymbol{\Sigma}_{f*}$$

- A **kernel function** k is symmetric and positive definite.
- Kernel computes covariances between unknown function values $f(\mathbf{x}_i)$ and $f(\mathbf{x}_j)$ by just looking at the corresponding inputs $\mathbf{x}_i, \mathbf{x}_j$

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

Kernelization

$$p(f_*|f) = \mathcal{N}\left(\begin{bmatrix} \mu_f \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{f*} \\ \Sigma_{*f} & \Sigma_{**} \end{bmatrix}\right)$$
$$p(f_*|f) = \mathcal{N}(f_* | \mathbf{m}, S)$$
$$\mathbf{m} = \boldsymbol{\mu}_* + \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} (f - \boldsymbol{\mu})$$
$$S = \boldsymbol{\Sigma}_{**} - \boldsymbol{\Sigma}_{*f} \boldsymbol{\Sigma}_{ff}^{-1} \boldsymbol{\Sigma}_{f*}$$

- A **kernel function** k is symmetric and positive definite.
- Kernel computes covariances between unknown function values $f(\mathbf{x}_i)$ and $f(\mathbf{x}_j)$ by just looking at the corresponding inputs $\mathbf{x}_i, \mathbf{x}_j$

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

- This yields the predictive distribution

$$p(f_*|f, \mathbf{X}, \mathbf{X}_*) = \mathcal{N}(f_* | \mathbf{m}, S)$$

$$\mathbf{m} = \boldsymbol{\mu}_* + k(\mathbf{X}_*, \mathbf{X}) k(\mathbf{X}, \mathbf{X})^{-1} (f - \boldsymbol{\mu})$$

$$S = \boldsymbol{\Sigma}_{**} - k(\mathbf{X}_*, \mathbf{X}) k(\mathbf{X}, \mathbf{X})^{-1} k(\mathbf{X}, \mathbf{X}_*)$$

GP Regression as a Bayesian Inference Problem

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

GP Regression as a Bayesian Inference Problem

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

Training data: X, y . Bayes' theorem yields

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

GP Regression as a Bayesian Inference Problem

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

Training data: X, y . Bayes' theorem yields

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

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

GP Regression as a Bayesian Inference Problem

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

Training data: X, y . Bayes' theorem yields

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

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

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

GP Regression as a Bayesian Inference Problem

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

Training data: X, y . Bayes' theorem yields

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

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

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

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

GP Regression as a Bayesian Inference Problem

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

Training data: X, y . Bayes' theorem yields

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

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

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

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

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

GP Regression as a Bayesian Inference Problem

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

Training data: X, y . Bayes' theorem yields

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

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

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

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

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

$$m_{\text{post}}(x_i) = m(x_i) + k(X, x_i)^\top (K + \sigma_\varepsilon^2 I)^{-1} (y - m(x_i))$$

$$k_{\text{post}}(x_i, x_j) = k(x_i, x_j) - k(X, x_i)^\top (K + \sigma_\varepsilon^2 I)^{-1} k(X, x_j)$$

GP Regression as a Bayesian Inference Problem (2)

Posterior Gaussian process:

$$m_{\text{post}}(\mathbf{x}_i) = m(\mathbf{x}_i) + k(\mathbf{X}, \mathbf{x}_i)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} (\mathbf{y} - m(\mathbf{x}_i))$$

$$k_{\text{post}}(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) - k(\mathbf{X}, \mathbf{x}_i)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_j)$$

GP Regression as a Bayesian Inference Problem (2)

Posterior Gaussian process:

$$m_{\text{post}}(\mathbf{x}_i) = m(\mathbf{x}_i) + k(\mathbf{X}, \mathbf{x}_i)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} (\mathbf{y} - m(\mathbf{x}_i))$$

$$k_{\text{post}}(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) - k(\mathbf{X}, \mathbf{x}_i)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_j)$$

Predictive distribution $p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ at test inputs \mathbf{x}_* :

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

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

$$\mathbb{V}[f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*] = k_{\text{post}}(\mathbf{x}_*, \mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{X}, \mathbf{x}_*)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_*)$$

GP Regression as a Bayesian Inference Problem (2)

Posterior Gaussian process:

$$m_{\text{post}}(\mathbf{x}_i) = m(\mathbf{x}_i) + k(\mathbf{X}, \mathbf{x}_i)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} (\mathbf{y} - m(\mathbf{x}_i))$$

$$k_{\text{post}}(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) - k(\mathbf{X}, \mathbf{x}_i)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_j)$$

Predictive distribution $p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ at test inputs \mathbf{x}_* :

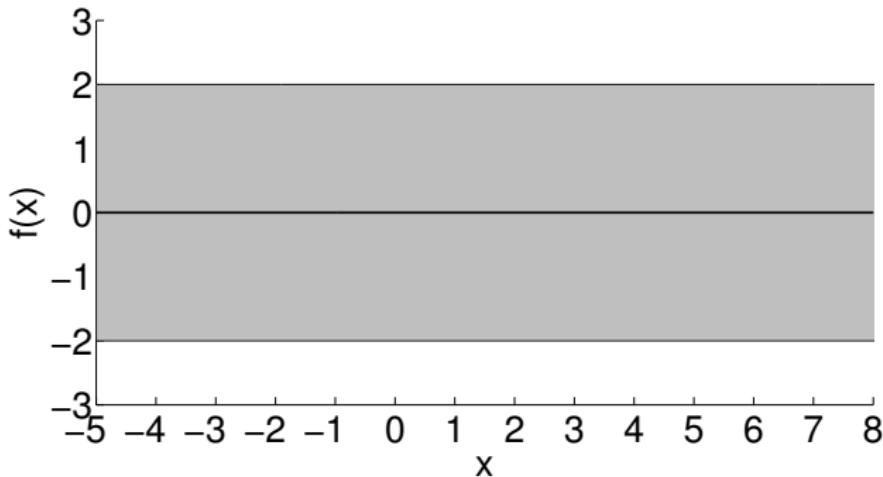
$$p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \mathcal{N}(\mathbb{E}[f_*], \mathbb{V}[f_*])$$

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

$$\mathbb{V}[f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*] = k_{\text{post}}(\mathbf{x}_*, \mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{X}, \mathbf{x}_*)^{\top} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_*)$$

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

Illustration



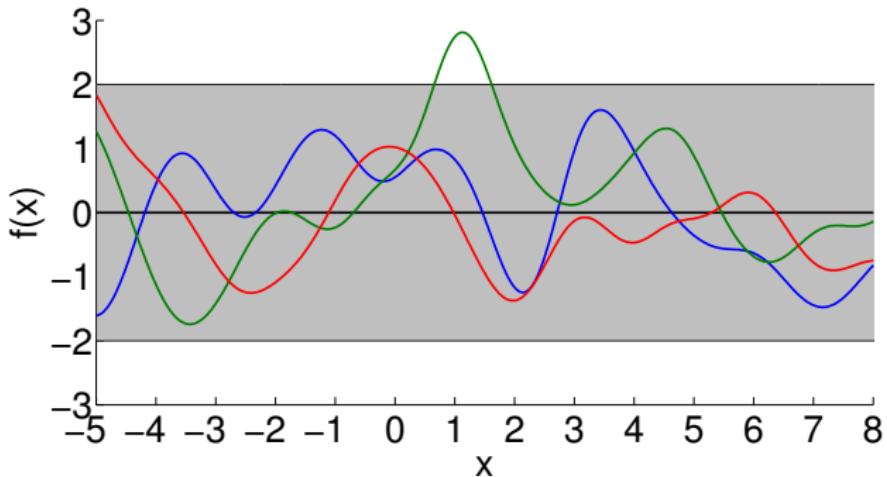
Prior belief about the function

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



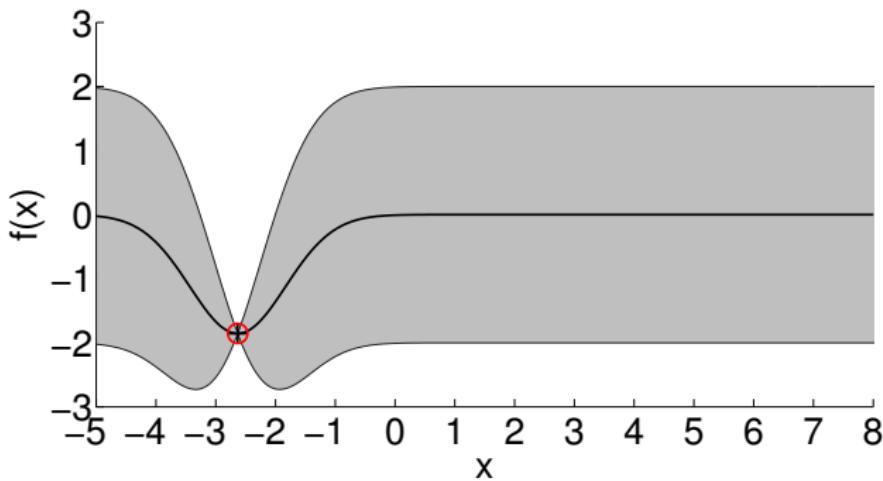
Prior belief about the function

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



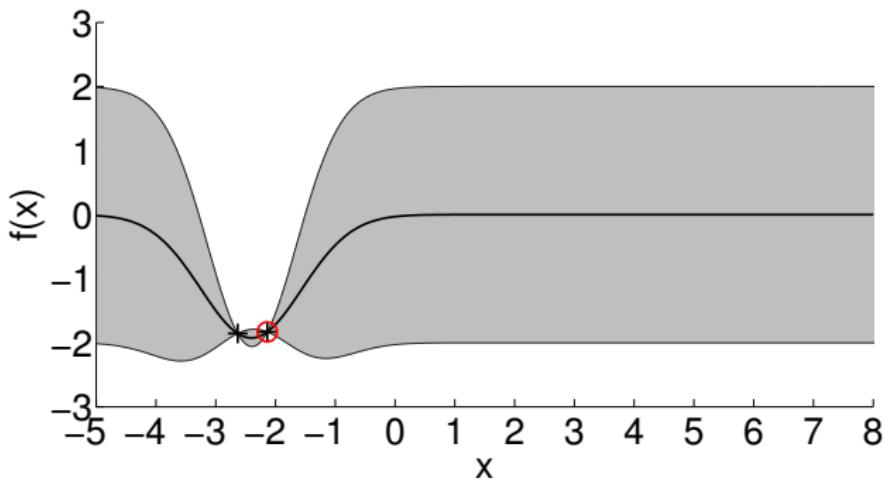
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



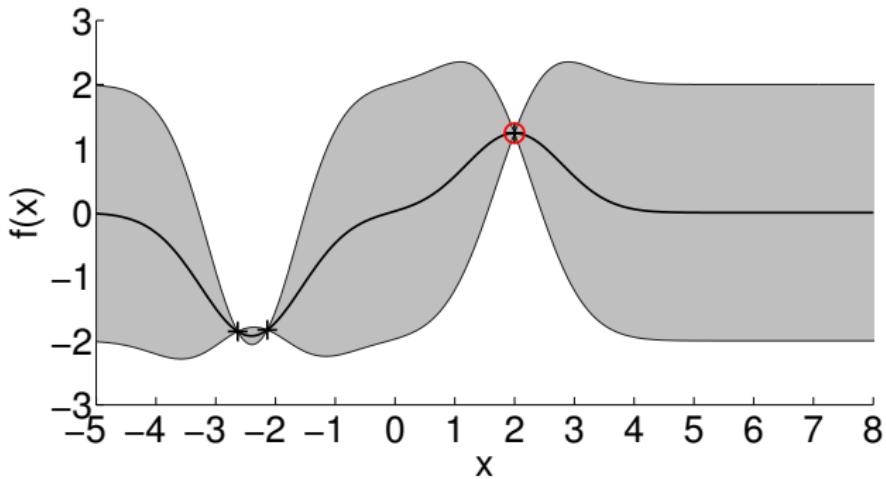
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



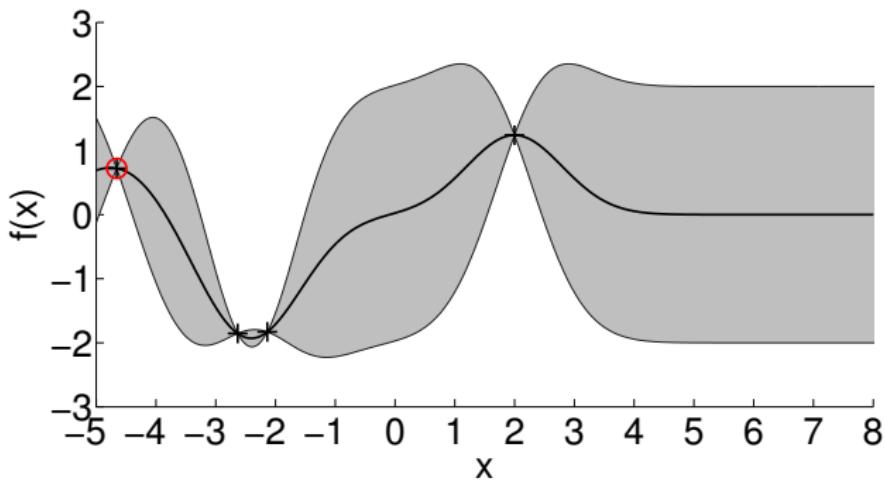
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



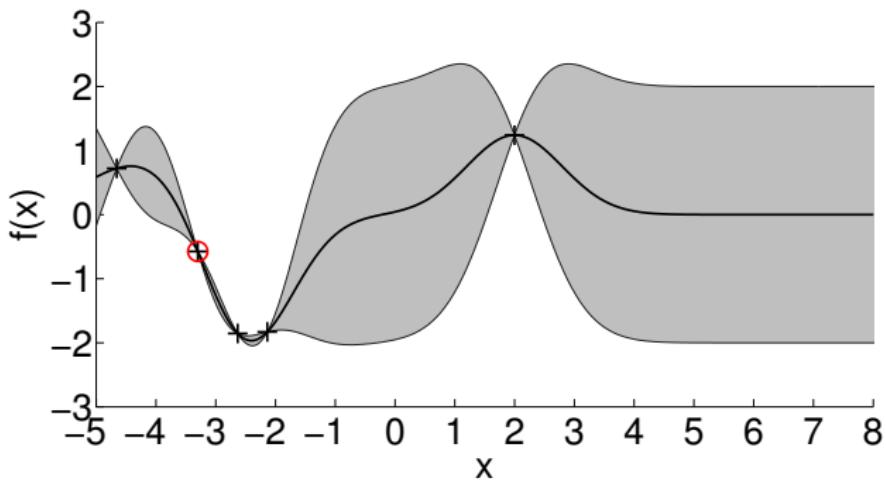
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



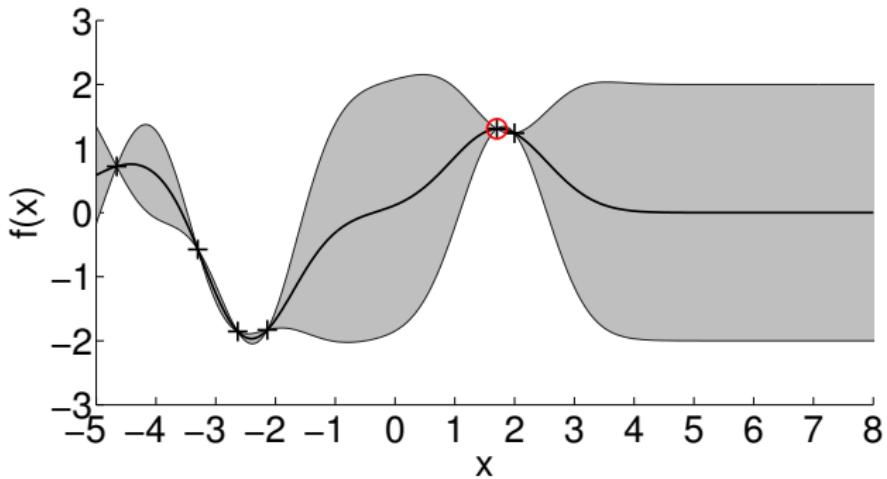
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



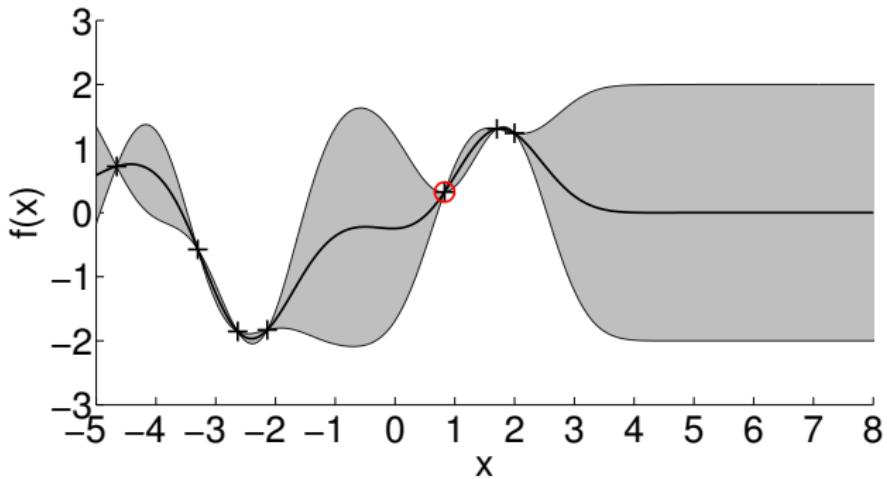
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



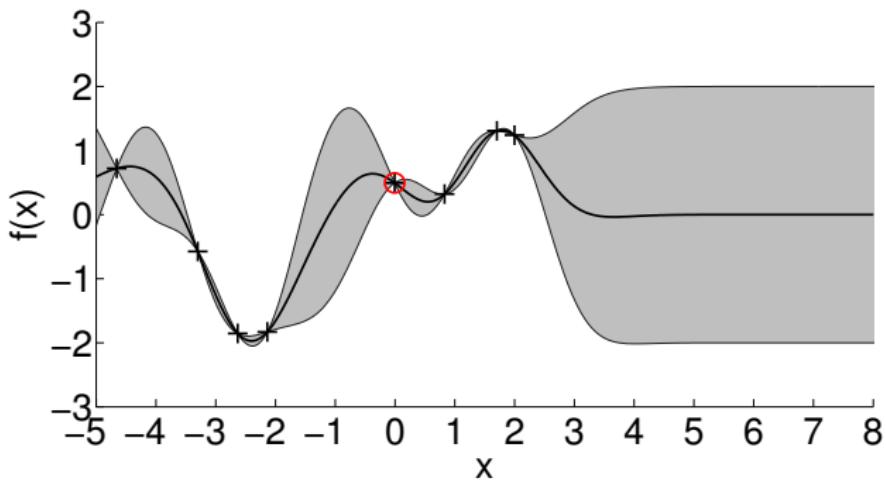
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



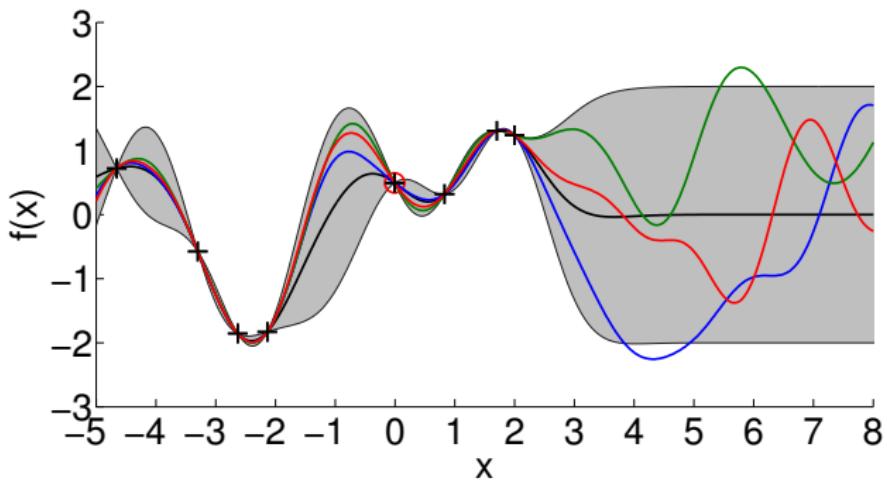
Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

Illustration



Posterior belief about the function

Predictive (marginal) mean and variance:

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

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

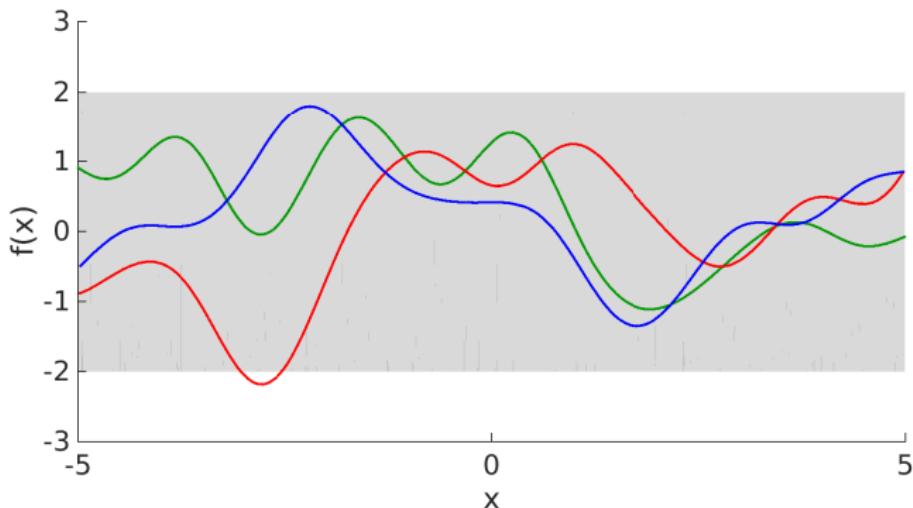
Covariance Function

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

Gaussian Covariance Function

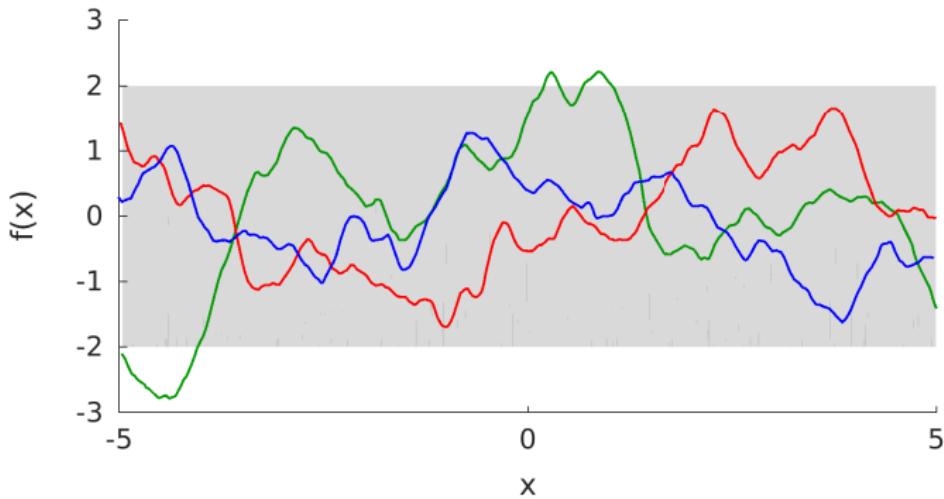
$$k_{Gauss}(\mathbf{x}_i, \mathbf{x}_j) = \theta_1^2 \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) / \theta_2^2\right)$$

- ▶ θ_1 : Amplitude of the latent function
- ▶ θ_2 : Length scale. How far do we have to move in input space before the function value changes significantly
- ▶ Smoothness parameter



Matérn Covariance Function

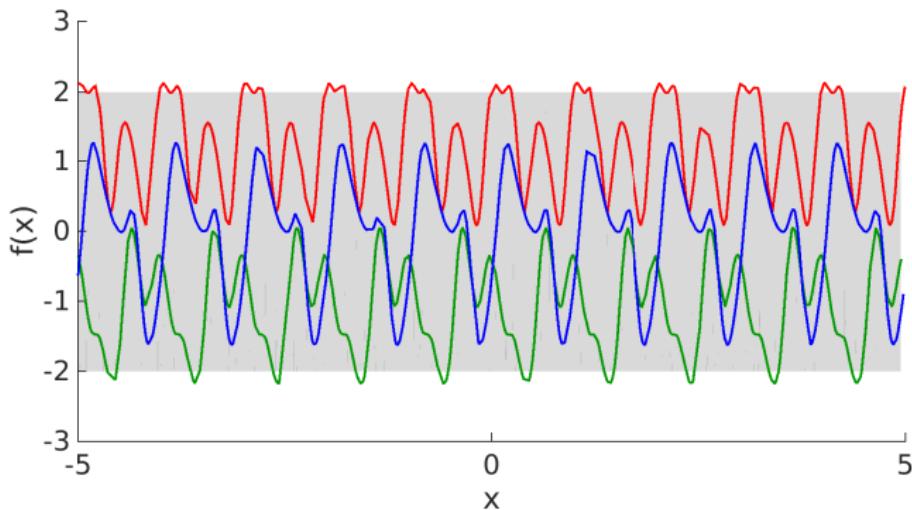
$$k_{Mat,3/2}(x_i, x_j) = \theta_1^2 \left(1 + \frac{\sqrt{3}(x_i - x_j)}{\theta_2} \right) \exp \left(-\frac{\sqrt{3}(x_i - x_j)}{\theta_2} \right)$$



Periodic Covariance Function

$$k_{per}(x_i, x_j) = \theta_1^2 \exp\left(-\frac{2 \sin^2\left(\frac{\theta_3(x_i - x_j)}{2\pi}\right)}{\theta_2^2}\right)$$
$$= k_{Gauss}(\mathbf{u}(x_i), \mathbf{u}(x_j)), \quad \mathbf{u}(x) = \begin{bmatrix} \cos(x) \\ \sin(x) \end{bmatrix}$$

θ_3 : Periodicity parameter



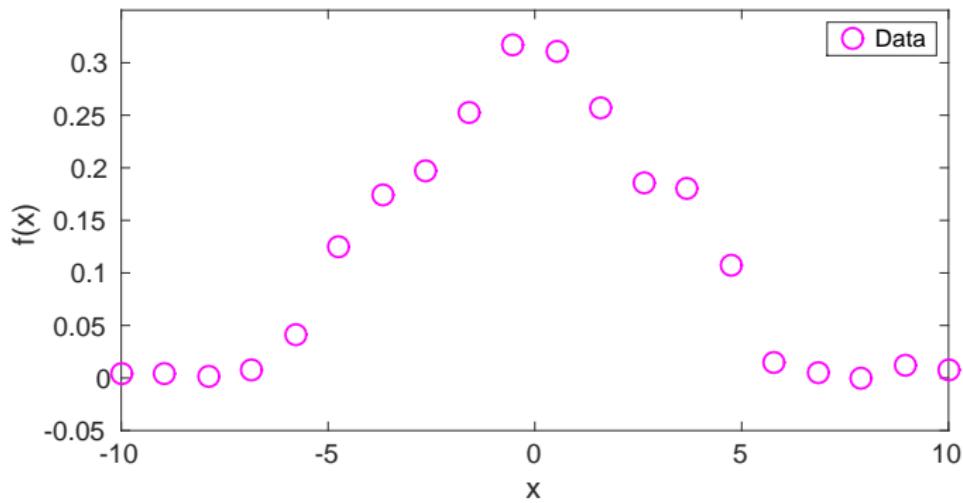
Model Selection

A GP is fully specified by a mean function m and a covariance function (kernel) k . Both functions possess parameters $\{\theta_m, \theta_k\} =: \theta$

- How do we find good parameters θ ?
 - How do we choose m and k ?
- Model selection

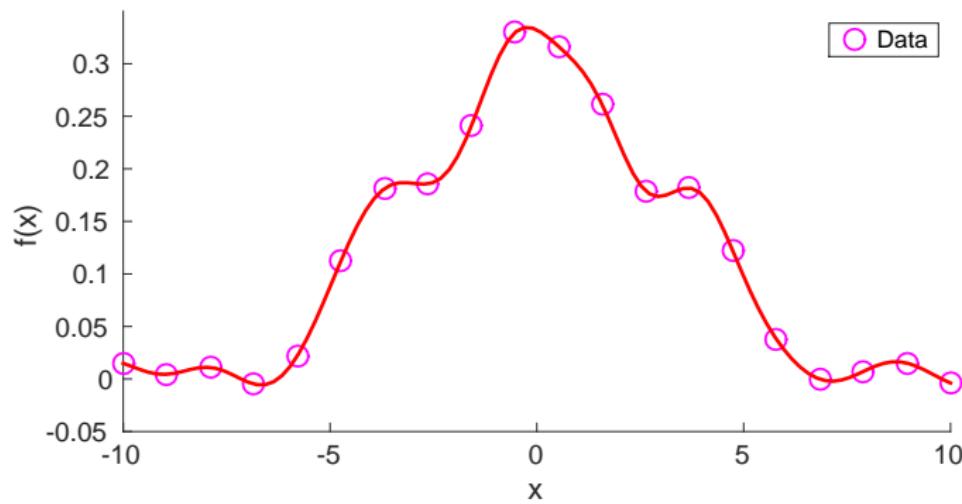
Model Selection I: Length-Scales

Length scales determine how wiggly the function is



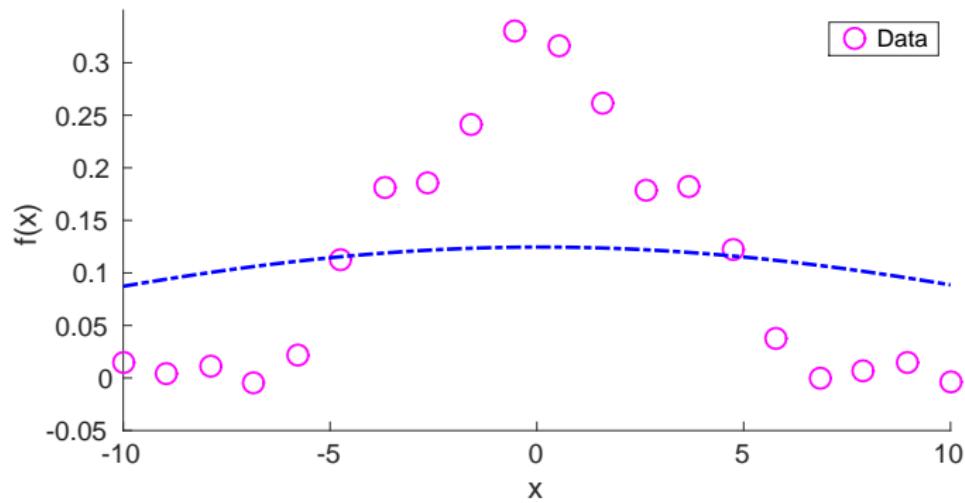
Model Selection I: Length-Scales

Length scales determine how wiggly the function is



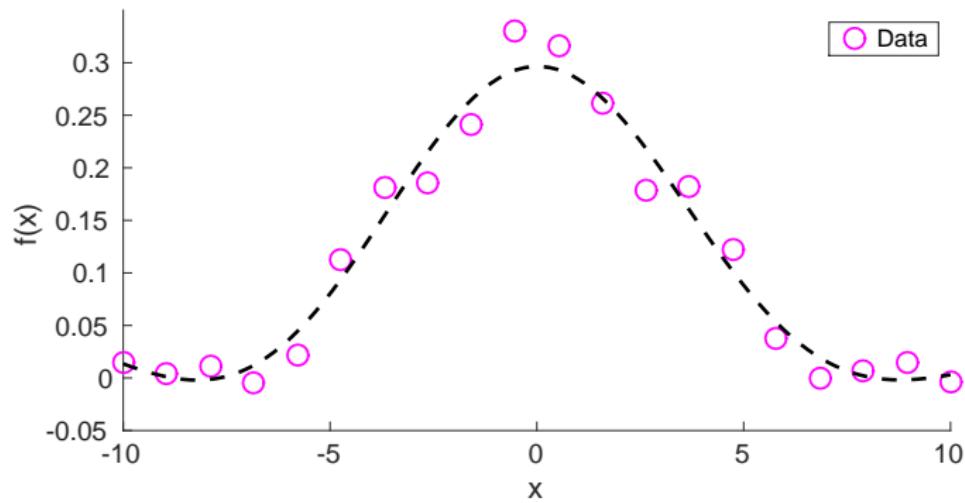
Model Selection I: Length-Scales

Length scales determine how wiggly the function is



Model Selection I: Length-Scales

Length scales determine how wiggly the function is



Model Selection: Hyper-Parameters

GP Training

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

Model Selection: Hyper-Parameters

GP Training

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

- Assign a prior $p(\theta)$ over hyper-parameters
- Posterior over hyper-parameters:

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

Model Selection: Hyper-Parameters

GP Training

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

- Assign a prior $p(\theta)$ over hyper-parameters
- Posterior over hyper-parameters:

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

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

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

Model Selection: Hyper-Parameters

GP Training

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

- Assign a prior $p(\theta)$ over hyper-parameters
- Posterior over hyper-parameters:

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

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

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

- 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)

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

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

Training via Marginal Likelihood Maximization

GP Training

Maximize the evidence/marginal likelihood (probability of the data given the hyper-parameters)

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

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

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

Training via Marginal Likelihood Maximization

GP Training

Maximize the evidence/marginal likelihood (probability of the data given the hyper-parameters)

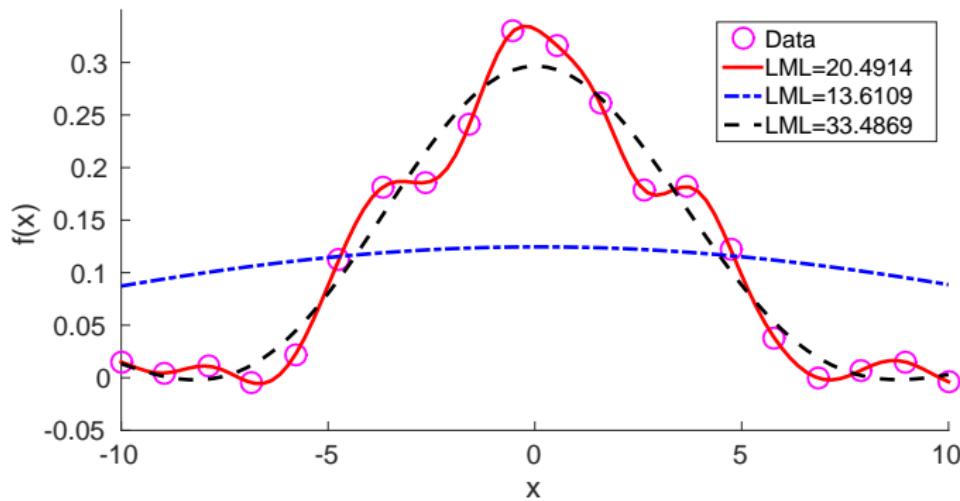
$$\theta^* \in \arg \max_{\theta} \log p(\mathbf{y}|\mathbf{X}, \theta)$$

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

- Automatic trade-off between **data fit** and **model complexity**
- Gradient-based optimization possible:

$$\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \theta)}{\partial \theta} = \frac{1}{2}\mathbf{y}^\top \mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \theta} \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \theta} \right)$$

Example



- ▶ Compromise between fitting the data and simplicity of the model
- ▶ Not fitting the data is explained by a larger noise variance (treated as an additional hyper-parameter and learned jointly)s

Model Selection II—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 II—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?
- ▶ Assign a prior $p(M_i)$ for each model
- ▶ Posterior model probability:

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

Model Selection II—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?
- ▶ Assign a prior $p(M_i)$ for each model
- ▶ Posterior model probability:

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

- ▶ Choose MAP model M_* , such that

$$M_* \in \arg \max_M \log p(M) + \log p(\mathbf{y}|\mathbf{X}, M)$$

Model Selection II—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?
- ▶ Assign a prior $p(M_i)$ for each model
- ▶ Posterior model probability:

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

- ▶ Choose MAP model M_* , such that

$$M_* \in \arg \max_M \log p(M) + \log p(\mathbf{y}|\mathbf{X}, M)$$

► Compare marginal likelihoods if $p(M_i) = 1/|M|$

Model Selection II—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?
- ▶ Assign a prior $p(M_i)$ for each model
- ▶ Posterior model probability:

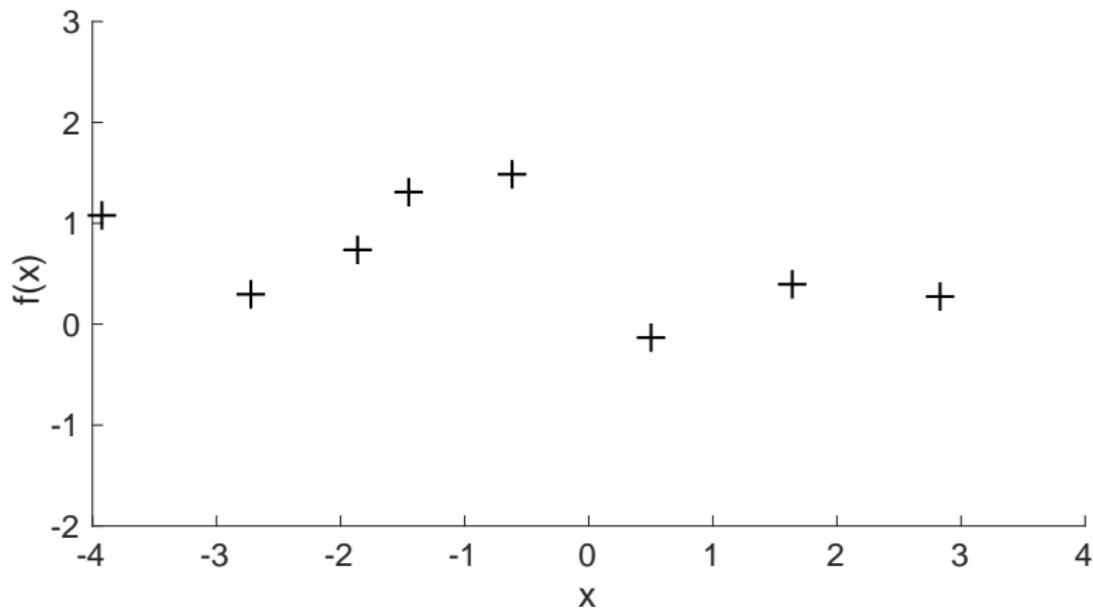
$$p(M_i|\mathbf{X}, \mathbf{y}) = \frac{p(M_i)p(\mathbf{y}|\mathbf{X}, M_i)}{p(\mathbf{y}|\mathbf{X})}$$

- ▶ Choose MAP model M_* , such that

$$M_* \in \arg \max_M \log p(M) + \log p(\mathbf{y}|\mathbf{X}, M)$$

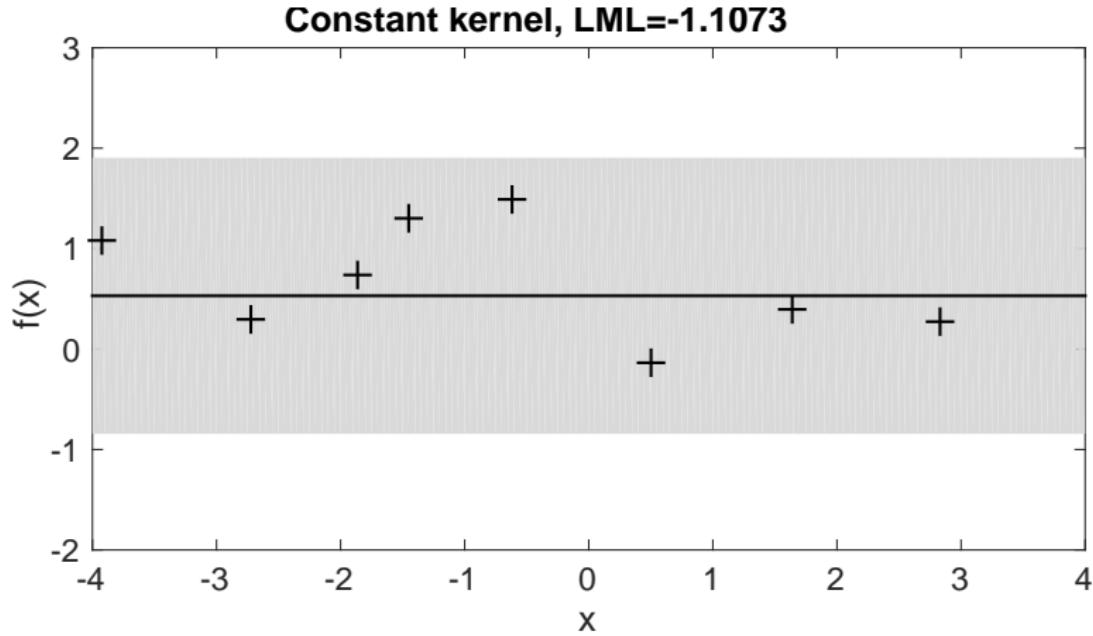
- ▶ Compare marginal likelihoods if $p(M_i) = 1/|M|$
- ▶ Marginal likelihood is central for model selection

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

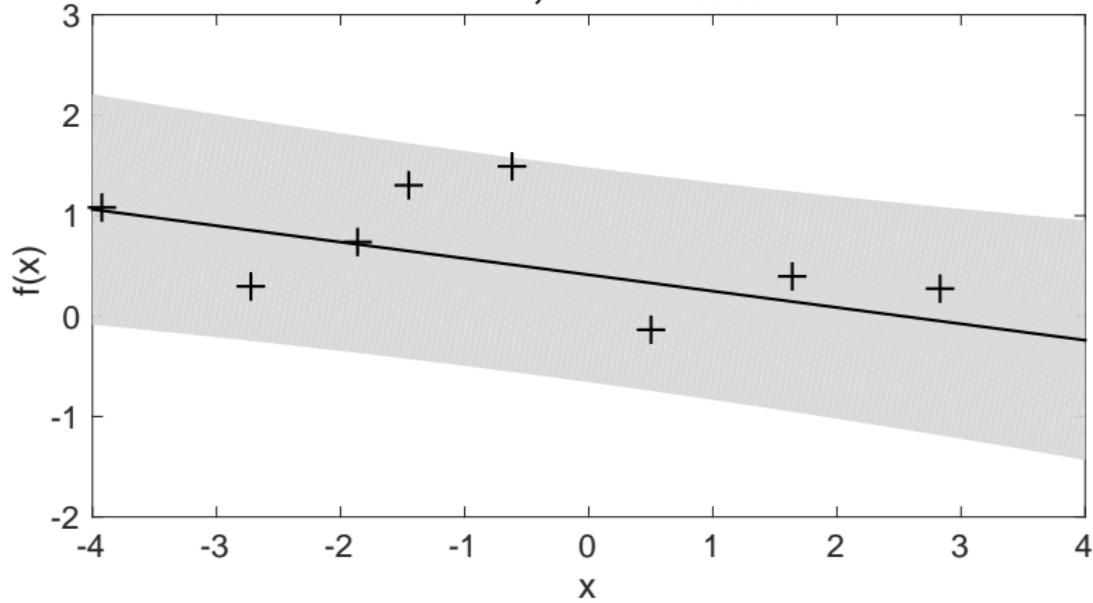
Example



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

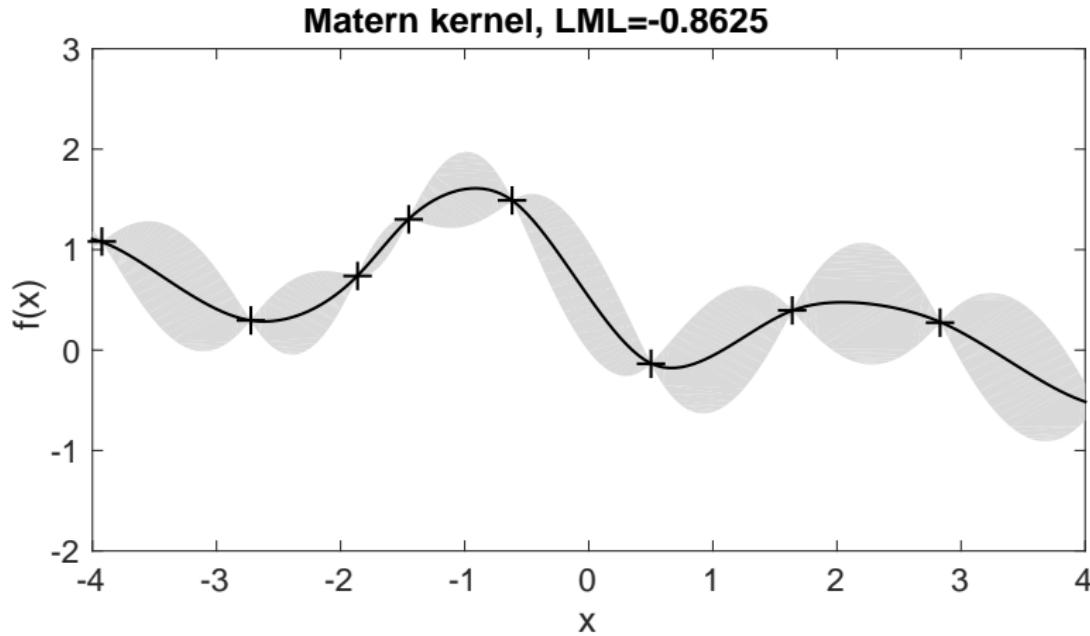
Example

Linear kernel, LML=-1.0065



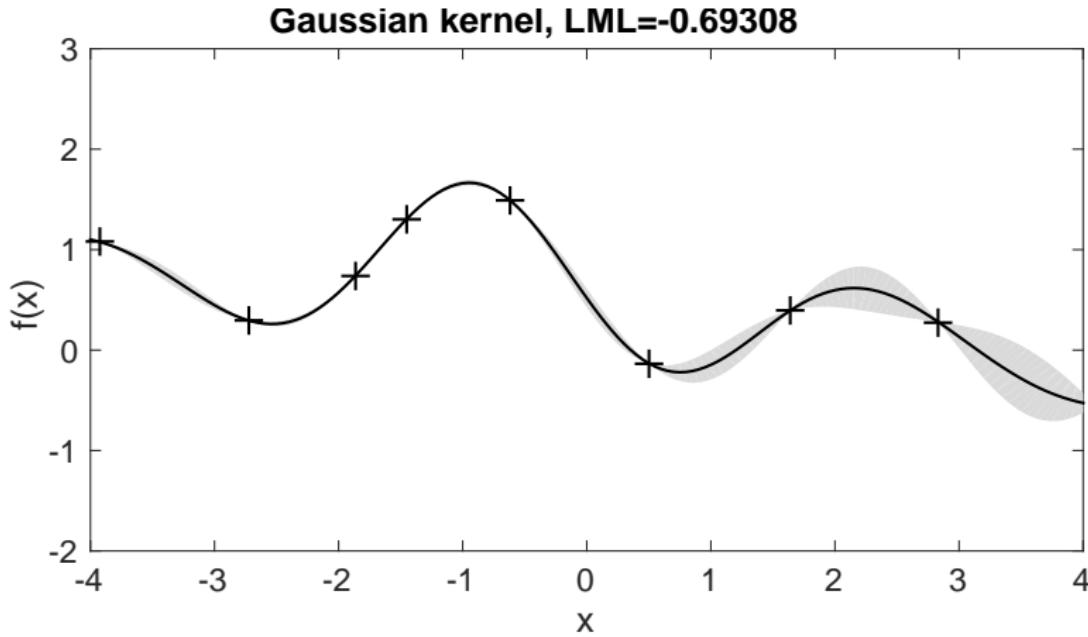
- Four different kernels (mean function fixed to $m \equiv 0$)
- MAP hyper-parameters for each kernel
- Log-marginal likelihood values for each (optimized) model

Example



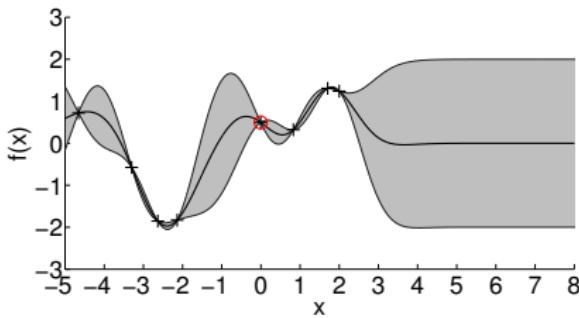
- Four different kernels (mean function fixed to $m \equiv 0$)
- MAP hyper-parameters for each kernel
- Log-marginal likelihood values for each (optimized) model

Example



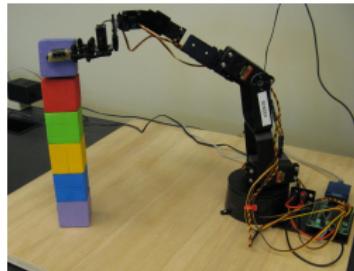
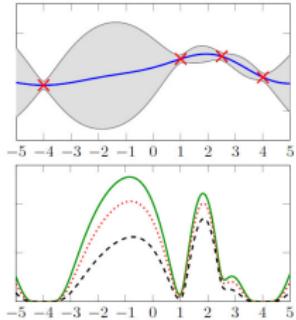
- Four different kernels (mean function fixed to $m \equiv 0$)
- MAP hyper-parameters for each kernel
- Log-marginal likelihood values for each (optimized) model

Quick Summary



- Gaussian processes are extremely **flexible** models
- **Consistent confidence bounds!**
- Based on simple manipulations of plain Gaussian distributions
- Few hyper-parameters ➡ **Generally easy to train**
- First choice for black-box regression

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

Limitations of Gaussian Processes

Computational and memory complexity

- Training scales in $\mathcal{O}(N^3)$
 - Prediction (variances) scales in $\mathcal{O}(N^2)$
 - Memory requirement: $\mathcal{O}(ND + N^2)$
- Practical limit $N \approx 10,000$

Table of Contents

Introduction

 Gaussian Distribution

Gaussian Processes

 Definition and Derivation

 Inference

 Covariance Functions

 Model Selection

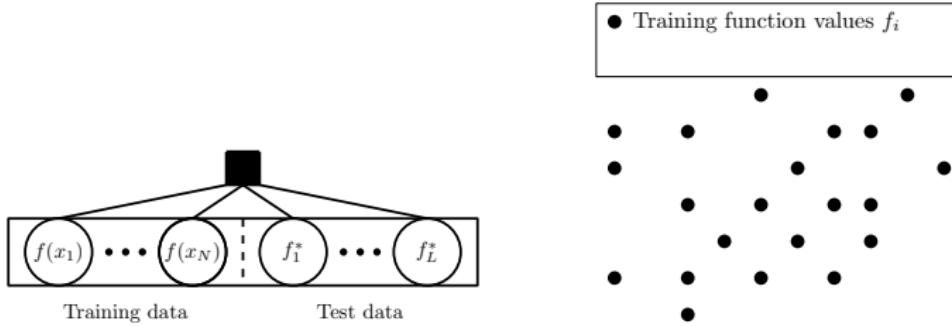
 Limitations

Scaling Gaussian Processes to Large Data Sets

 Sparse Gaussian Processes

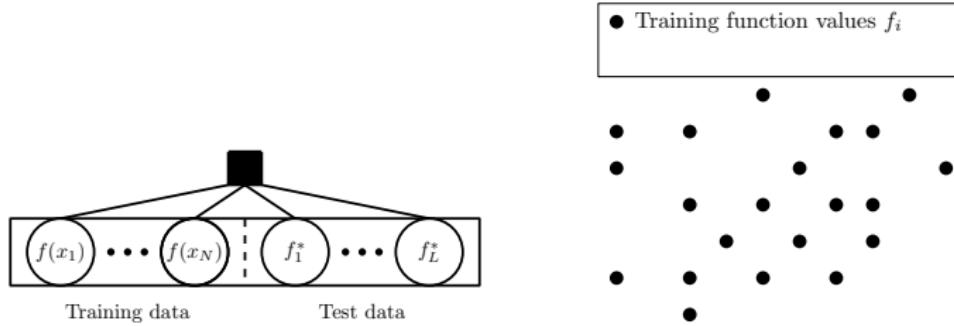
 Distributed Gaussian Processes

GP Factor Graph



- ▶ Probabilistic graphical model (factor graph) of a GP
- ▶ All function values are jointly Gaussian distributed (e.g., training and test function values)

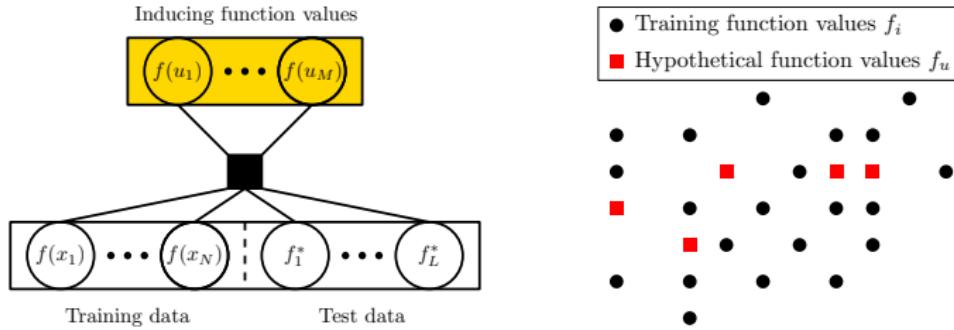
GP Factor Graph



- ▶ Probabilistic graphical model (factor graph) of a GP
- ▶ All function values are jointly Gaussian distributed (e.g., training and test function values)
- ▶ GP prior

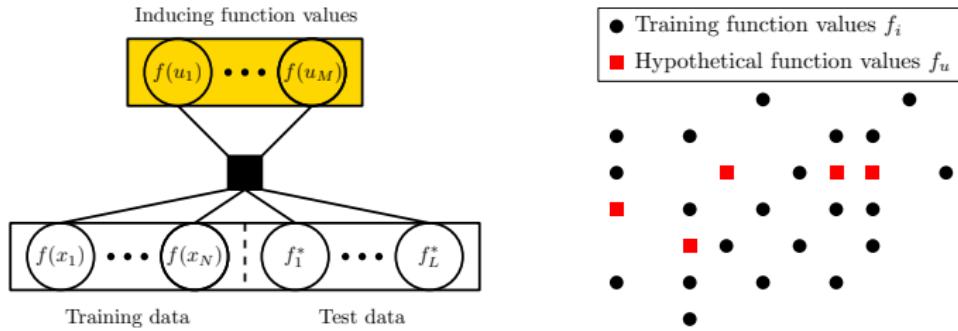
$$p(f, f_*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} K_{ff} & K_{f*} \\ K_{*f} & K_{**} \end{bmatrix} \right)$$

Inducing Variables



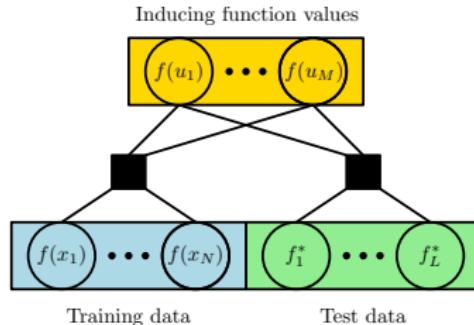
- ▶ Introduce **inducing function values** f_u
- ▶ “Hypothetical” function values

Inducing Variables



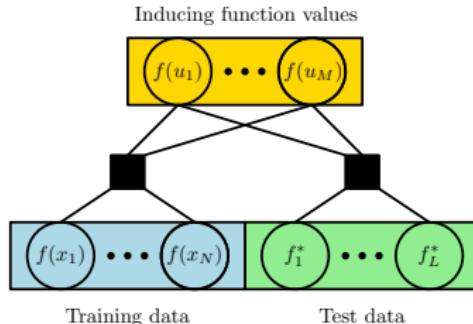
- ▶ Introduce **inducing function values** f_u
 - ▶ “Hypothetical” function values
- ▶ All function values are still jointly Gaussian distributed (e.g., training, test and inducing function values)

Central Approximation Scheme



- Approximation: Training and test set are **conditionally independent given the inducing function values**: $f \perp\!\!\!\perp f_* | f_u$

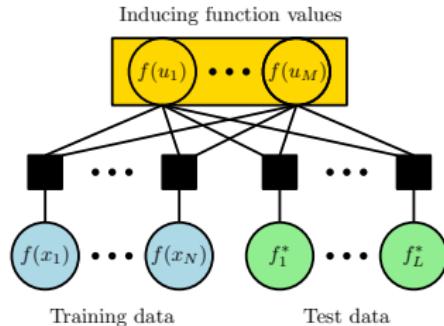
Central Approximation Scheme



- Approximation: Training and test set are **conditionally independent given the inducing function values**: $f \perp\!\!\!\perp f_* | f_u$
- Then, the effective GP prior is

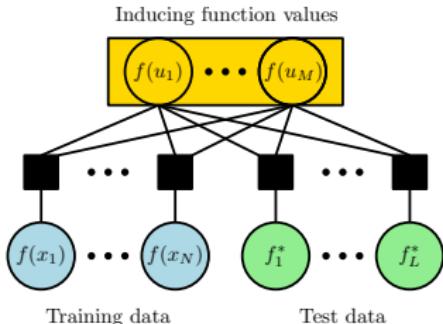
$$\begin{aligned} q(f, f_*) &= \int p(f|f_u)p(f_*|f_u)p(f_u)df_u \\ &= \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{Q}_{f*} \\ \mathbf{Q}_{*f} & \mathbf{K}_{**} \end{bmatrix} \right), \quad \mathbf{Q}_{ab} := \mathbf{K}_{af_u} \mathbf{K}_{f_u f_u}^{-1} \mathbf{K}_{f_u b} \end{aligned}$$

FI(T)C Sparse Approximation



- Assume that training (and test sets) are fully independent given the inducing variables (Snelson & Ghahramani, 2006)

FI(T)C Sparse Approximation

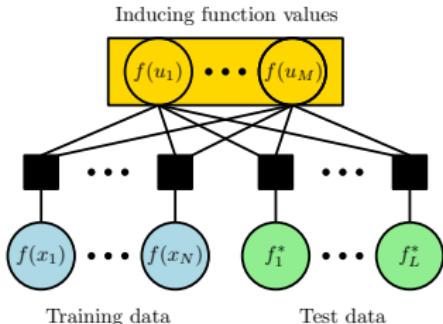


- Assume that training (and test sets) are fully independent given the inducing variables (Snelson & Ghahramani, 2006)
- Effective GP prior with this approximation

$$q(f, f_*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} Q_{ff} - \text{diag}(Q_{ff} - K_{ff}) & Q_{f*} \\ Q_{*f} & K_{**} \end{bmatrix} \right)$$

- $Q_{**} - \text{diag}(Q_{**} - K_{**})$ can be used instead of K_{**} ➤ FIC

FI(T)C Sparse Approximation



- Assume that training (and test sets) are fully independent given the inducing variables (Snelson & Ghahramani, 2006)
- Effective GP prior with this approximation

$$q(f, f_*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} Q_{ff} - \text{diag}(Q_{ff} - K_{ff}) & Q_{f*} \\ Q_{*f} & K_{**} \end{bmatrix} \right)$$

- $Q_{**} - \text{diag}(Q_{**} - K_{**})$ can be used instead of K_{**} ➤ FIC
- Training: $\mathcal{O}(NM^2)$, Prediction: $\mathcal{O}(M^2)$

Inducing Inputs

- FI(T)C sparse approximation relies on inducing function values $f(\mathbf{u}_i)$, where \mathbf{u}_i are the corresponding inputs

Inducing Inputs

- FI(T)C sparse approximation relies on inducing function values $f(\mathbf{u}_i)$, where \mathbf{u}_i are the corresponding inputs
- These **inputs are unknown** a priori ➤ Find “optimal” ones

Inducing Inputs

- FI(T)C sparse approximation relies on inducing function values $f(\mathbf{u}_i)$, where \mathbf{u}_i are the corresponding inputs
- These **inputs are unknown** a priori ➤ Find “optimal” ones
- Find them by **maximizing the FI(T)C marginal likelihood with respect to the inducing inputs**
(and the standard hyper-parameters):

$$\mathbf{u}_{1:M}^* \in \arg \max_{\mathbf{u}_{1:M}} q_{\text{FITC}}(\mathbf{y} | \mathbf{X}, \mathbf{u}_{1:M}, \boldsymbol{\theta})$$

Inducing Inputs

- FI(T)C sparse approximation relies on inducing function values $f(\mathbf{u}_i)$, where \mathbf{u}_i are the corresponding inputs
- These **inputs are unknown** a priori ➤ Find “optimal” ones
- Find them by **maximizing the FI(T)C marginal likelihood with respect to the inducing inputs**
(and the standard hyper-parameters):

$$\mathbf{u}_{1:M}^* \in \arg \max_{\mathbf{u}_{1:M}} q_{\text{FITC}}(\mathbf{y} | \mathbf{X}, \mathbf{u}_{1:M}, \theta)$$

- Intuitively: The **marginal likelihood** is not only **parameterized** by the hyper-parameters θ , but also **by the inducing inputs $\mathbf{u}_{1:M}$** .

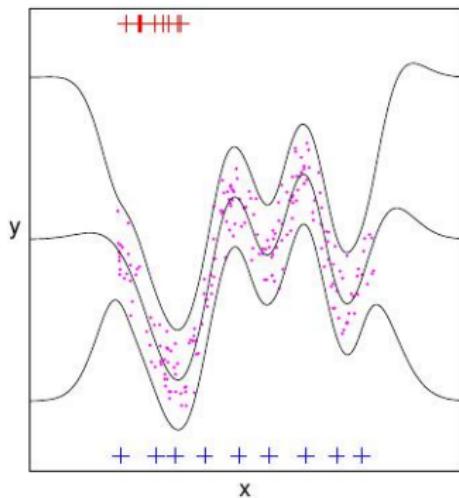
Inducing Inputs

- FI(T)C sparse approximation relies on inducing function values $f(\mathbf{u}_i)$, where \mathbf{u}_i are the corresponding inputs
- These **inputs are unknown** a priori ➤ Find “optimal” ones
- Find them by **maximizing the FI(T)C marginal likelihood with respect to the inducing inputs**
(and the standard hyper-parameters):

$$\mathbf{u}_{1:M}^* \in \arg \max_{\mathbf{u}_{1:M}} q_{\text{FITC}}(\mathbf{y} | \mathbf{X}, \mathbf{u}_{1:M}, \theta)$$

- Intuitively: The **marginal likelihood** is not only **parameterized** by the hyper-parameters θ , but also **by the inducing inputs $\mathbf{u}_{1:M}$** .
- End up with a **high-dimensional non-convex optimization** problem with **MD additional parameters**

FITC Example



- ▶ Pink: Original data
- ▶ Red crosses: Initialization of inducing inputs
- ▶ Blue crosses: Location of inducing inputs after optimization

Figure from Ed Snelson

- ▶ Efficient compression of the original data set

Summary Sparse Gaussian Processes

- Sparse approximations typically approximate a GP with N data points by a model with $M \ll N$ data points

Summary Sparse Gaussian Processes

- Sparse approximations typically approximate a GP with N data points by a model with $M \ll N$ data points
- Selection of these M data points can be tricky and may involve non-trivial computations (e.g., optimizing inducing inputs)

Summary Sparse Gaussian Processes

- Sparse approximations typically approximate a GP with N data points by a model with $M \ll N$ data points
- Selection of these M data points can be tricky and may involve non-trivial computations (e.g., optimizing inducing inputs)
- Simple (random) subset selection is fast and generally robust (Chalupka et al., 2013)

Summary Sparse Gaussian Processes

- Sparse approximations typically approximate a GP with N data points by a model with $M \ll N$ data points
- Selection of these M data points can be tricky and may involve non-trivial computations (e.g., optimizing inducing inputs)
- Simple (random) subset selection is fast and generally robust (Chalupka et al., 2013)
- Computational complexity: $\mathcal{O}(M^3)$ or $\mathcal{O}(NM^2)$ for training

Summary Sparse Gaussian Processes

- Sparse approximations typically approximate a GP with N data points by a model with $M \ll N$ data points
- Selection of these M data points can be tricky and may involve non-trivial computations (e.g., optimizing inducing inputs)
- Simple (random) subset selection is fast and generally robust (Chalupka et al., 2013)
- Computational complexity: $\mathcal{O}(M^3)$ or $\mathcal{O}(NM^2)$ for training
- Practical limit $M \leq 10^4$. Often: $M \in \mathcal{O}(10^2)$ in the case of inducing variables

Summary Sparse Gaussian Processes

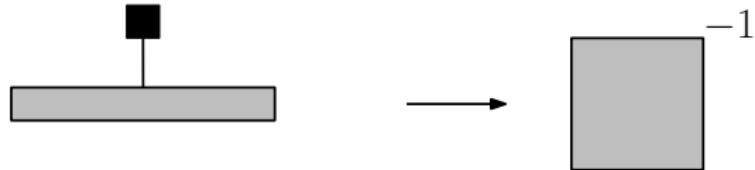
- Sparse approximations typically approximate a GP with N data points by a model with $M \ll N$ data points
- Selection of these M data points can be tricky and may involve non-trivial computations (e.g., optimizing inducing inputs)
- Simple (random) subset selection is fast and generally robust (Chalupka et al., 2013)
- Computational complexity: $\mathcal{O}(M^3)$ or $\mathcal{O}(NM^2)$ for training
- Practical limit $M \leq 10^4$. Often: $M \in \mathcal{O}(10^2)$ in the case of inducing variables
- If we set $M = N/100$, i.e., each inducing function value summarizes 100 real function values, our practical limit is $N \in \mathcal{O}(10^6)$

Summary Sparse Gaussian Processes

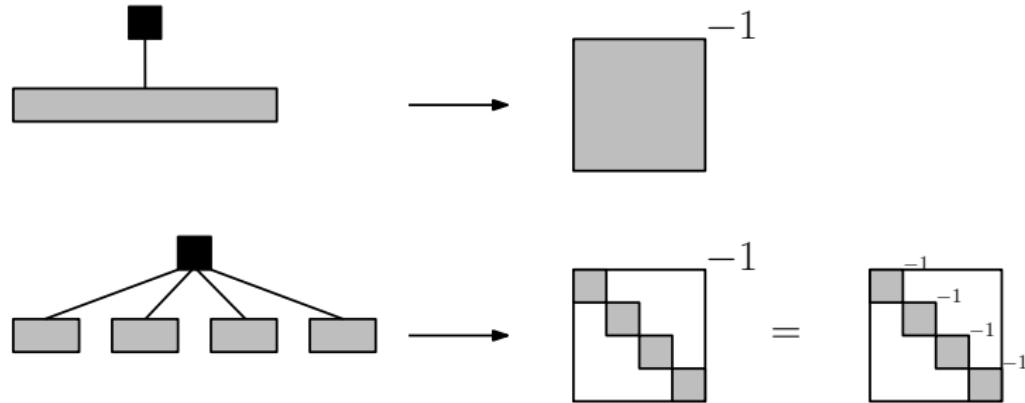
- Sparse approximations typically approximate a GP with N data points by a model with $M \ll N$ data points
- Selection of these M data points can be tricky and may involve non-trivial computations (e.g., optimizing inducing inputs)
- Simple (random) subset selection is fast and generally robust (Chalupka et al., 2013)
- Computational complexity: $\mathcal{O}(M^3)$ or $\mathcal{O}(NM^2)$ for training
- Practical limit $M \leq 10^4$. Often: $M \in \mathcal{O}(10^2)$ in the case of inducing variables
- If we set $M = N/100$, i.e., each inducing function value summarizes 100 real function values, our practical limit is $N \in \mathcal{O}(10^6)$

► Let's try something different

An Orthogonal Approximation: Distributed GPs

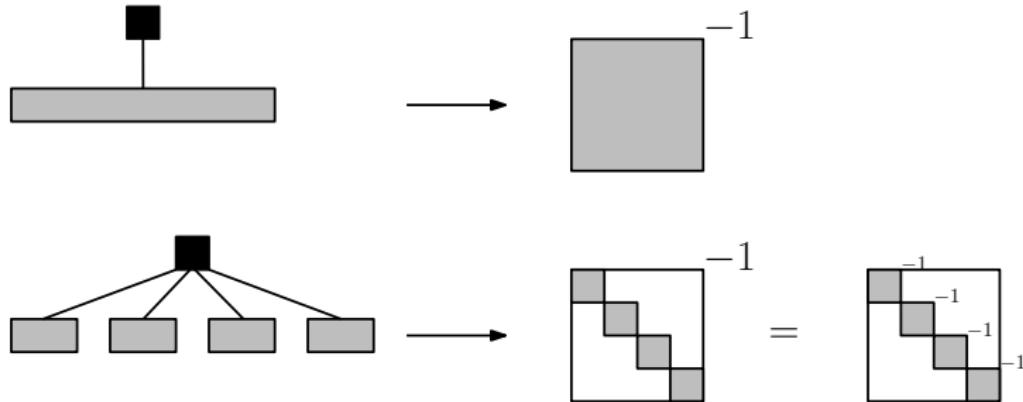


An Orthogonal Approximation: Distributed GPs



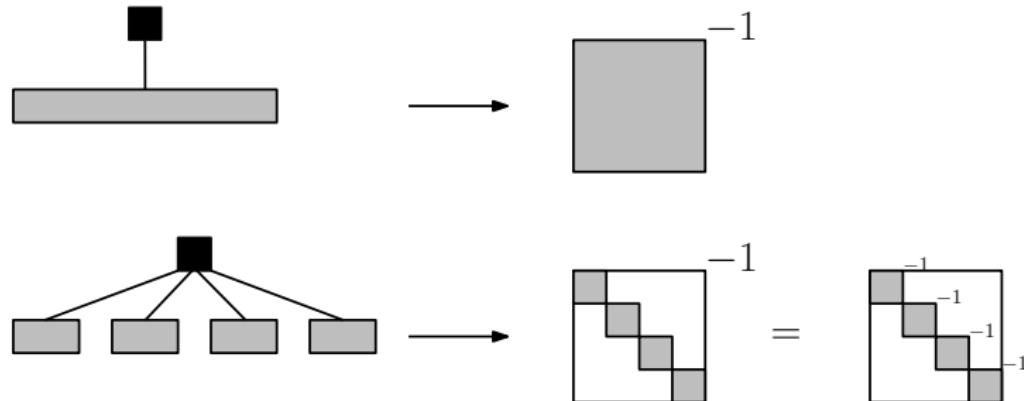
- Randomly split the full data set into M chunks

An Orthogonal Approximation: Distributed GPs



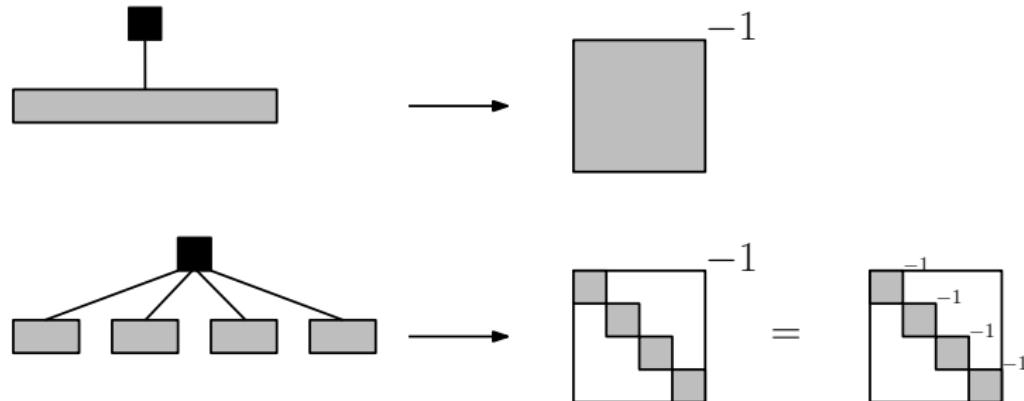
- Randomly split the full data set into M chunks
- Place M **independent GP models** (experts) on these small chunks

An Orthogonal Approximation: Distributed GPs



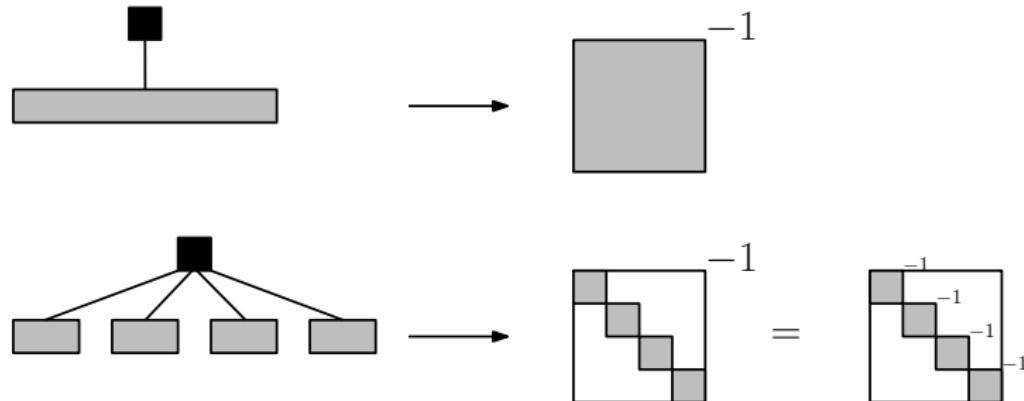
- Randomly split the full data set into M chunks
- Place M **independent GP models** (experts) on these small chunks
- Independent computations can be distributed

An Orthogonal Approximation: Distributed GPs



- Randomly split the full data set into M chunks
- Place M **independent GP models** (experts) on these small chunks
- Independent computations can be distributed
- Block-diagonal approximation of kernel matrix K

An Orthogonal Approximation: Distributed GPs



- Randomly split the full data set into M chunks
- Place M **independent GP models** (experts) on these small chunks
- Independent computations can be distributed
- Block-diagonal approximation of kernel matrix \mathbf{K}
- Combine independent computations to an overall result

Training the Distributed GP

- Split data set of size N into M chunks of size P
- Independence of experts ➤ Factorization of marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) \approx \sum_{k=1}^M \log p_k(\mathbf{y}^{(k)}|\mathbf{X}^{(k)}, \boldsymbol{\theta})$$

Training the Distributed GP

- Split data set of size N into M chunks of size P
- Independence of experts ➤ Factorization of marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{X}, \theta) \approx \sum_{k=1}^M \log p_k(\mathbf{y}^{(k)}|\mathbf{X}^{(k)}, \theta)$$

- Distributed optimization and training straightforward

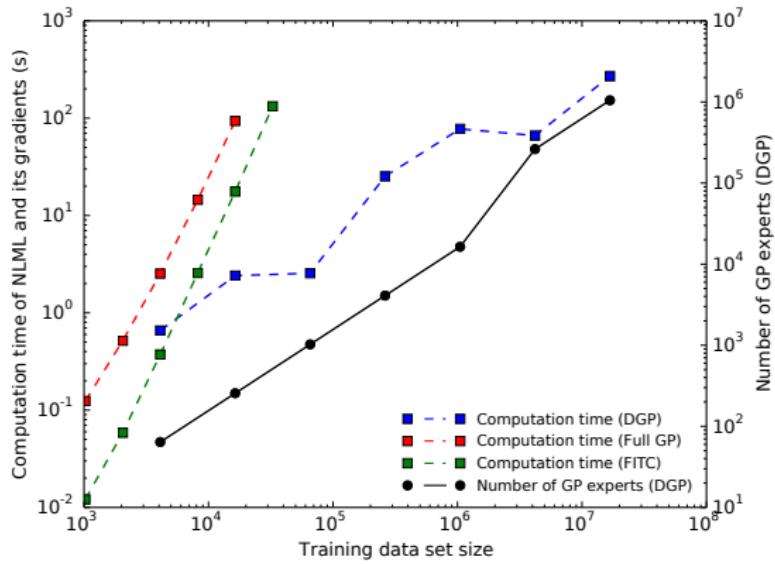
Training the Distributed GP

- ▶ Split data set of size N into M chunks of size P
- ▶ Independence of experts ➤ Factorization of marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{X}, \theta) \approx \sum_{k=1}^M \log p_k(\mathbf{y}^{(k)}|\mathbf{X}^{(k)}, \theta)$$

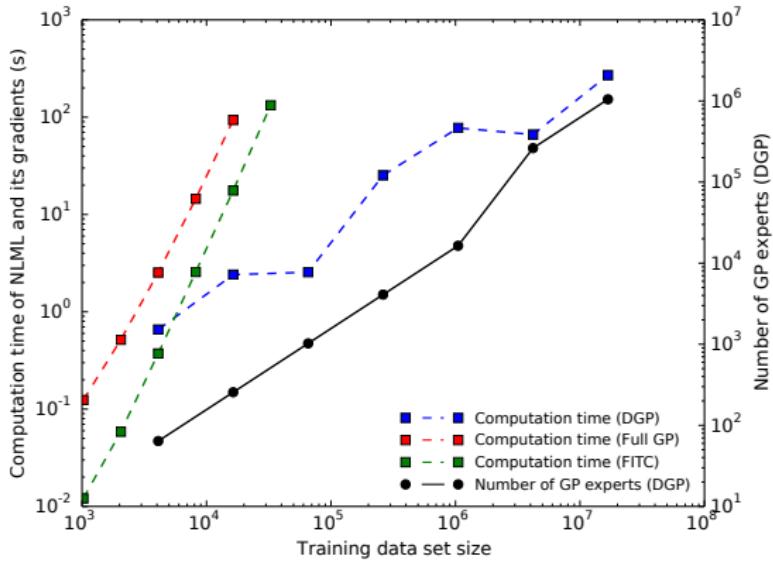
- ▶ Distributed optimization and training straightforward
- ▶ Computational complexity: $\mathcal{O}(MP^3)$ [instead of $\mathcal{O}(N^3)$]
But distributed over many machines
- ▶ Memory footprint: $\mathcal{O}(MP^2 + ND)$ [instead of $\mathcal{O}(N^2 + ND)$]

Empirical Training Time



- NLML is proportional to training time

Empirical Training Time



- NLML is proportional to training time
- Full GP (16K training points) \approx sparse GP (50K training points)
 \approx distributed GP (16M training points)

► Push practical limit by order(s) of magnitude

Practical Training Times

- Training* with $N = 10^6, D = 1$ on a laptop: ≈ 30 min
- Training* with $N = 5 \times 10^6, D = 8$ on a workstation: ≈ 4 hours

Practical Training Times

- Training* with $N = 10^6, D = 1$ on a laptop: ≈ 30 min
- Training* with $N = 5 \times 10^6, D = 8$ on a workstation: ≈ 4 hours

*: Maximize the marginal likelihood, stop when converged**

: Convergence often after 30–80 line searches*

Practical Training Times

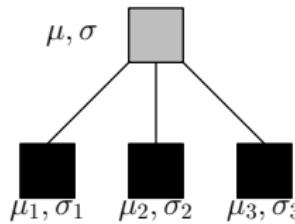
- Training* with $N = 10^6, D = 1$ on a laptop: ≈ 30 min
- Training* with $N = 5 \times 10^6, D = 8$ on a workstation: ≈ 4 hours

*: Maximize the marginal likelihood, stop when converged**

: Convergence often after 30–80 line searches*

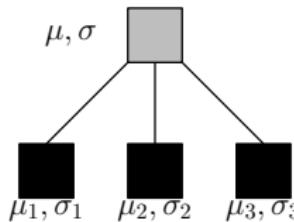
***: Line search $\approx 2\text{--}3$ evaluations of marginal likelihood and its gradient (usually $\mathcal{O}(N^3)$)

Predictions with the Distributed GP



- ▶ Prediction of each GP expert is Gaussian $\mathcal{N}(\mu_i, \sigma_i^2)$
- ▶ How to combine them to an overall prediction $\mathcal{N}(\mu, \sigma^2)$?

Predictions with the Distributed GP



- ▶ Prediction of each GP expert is Gaussian $\mathcal{N}(\mu_i, \sigma_i^2)$
- ▶ How to combine them to an overall prediction $\mathcal{N}(\mu, \sigma^2)$?

► Product-of-GP-experts

- ▶ PoE (product of experts) ► (Ng & Deisenroth, 2014)
- ▶ gPoE (generalized product of experts) ► (Cao & Fleet, 2014)
- ▶ BCM (Bayesian Committee Machine) ► (Tresp, 2000)
- ▶ rBCM (robust BCM) ► (Deisenroth & Ng, 2015)

Objectives

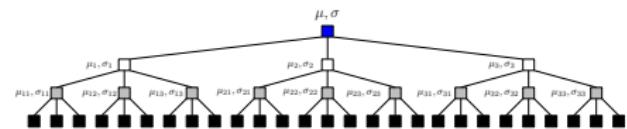
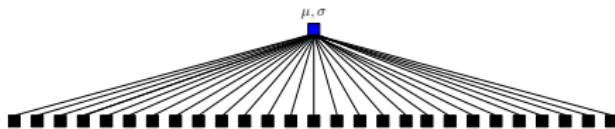


Figure: Two computational graphs

- Scale to large data sets ✓

Objectives

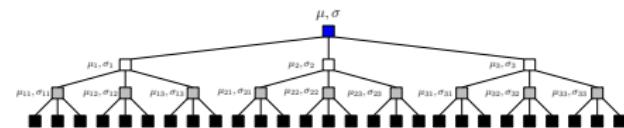
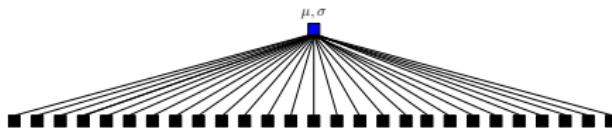


Figure: Two computational graphs

- Scale to large data sets ✓
- Good approximation of full GP (“ground truth”)

Objectives

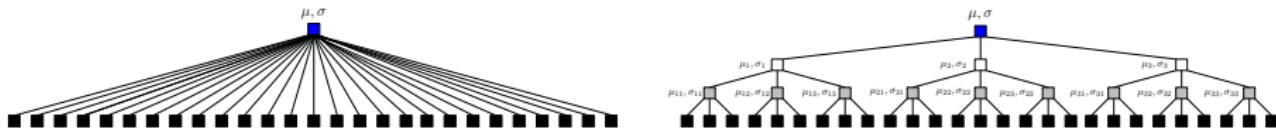


Figure: Two computational graphs

- Scale to large data sets ✓
- Good approximation of full GP (“ground truth”)
- Predictions independent of computational graph
 - ▶ Runs on heterogeneous computing infrastructures (laptop, cluster, ...)

Objectives

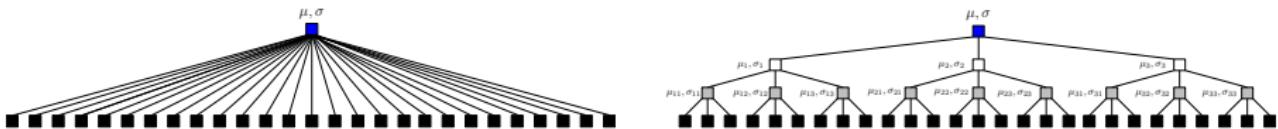
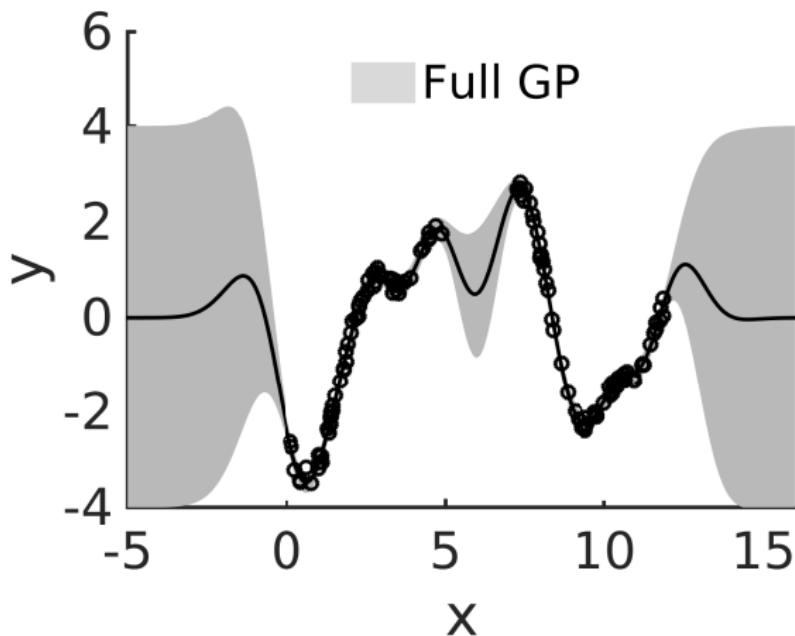


Figure: Two computational graphs

- Scale to large data sets ✓
- Good approximation of full GP (“ground truth”)
- Predictions independent of computational graph
 - ▶ Runs on heterogeneous computing infrastructures (laptop, cluster, ...)
- Reasonable predictive variances

Running Example



- ▶ Investigate various product-of-experts models
Same training procedure, but different mechanisms for predictions

Product of GP Experts

- Prediction model (independent predictors):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) ,$$

$$p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) = \mathcal{N}(f_* | \mu_k(\mathbf{x}_*), \sigma_k^2(\mathbf{x}_*))$$

Product of GP Experts

- ▶ Prediction model (independent predictors):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) ,$$

$$p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) = \mathcal{N}(f_* | \mu_k(\mathbf{x}_*), \sigma_k^2(\mathbf{x}_*))$$

- ▶ Predictive precision (inverse variance) and mean:

$$(\sigma_*^{\text{poe}})^{-2} = \sum_k \sigma_k^{-2}(\mathbf{x}_*)$$

$$\mu_*^{\text{poe}} = (\sigma_*^{\text{poe}})^2 \sum_k \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

Product of GP Experts

- ▶ Prediction model (independent predictors):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}),$$

$$p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) = \mathcal{N}(f_* | \mu_k(\mathbf{x}_*), \sigma_k^2(\mathbf{x}_*))$$

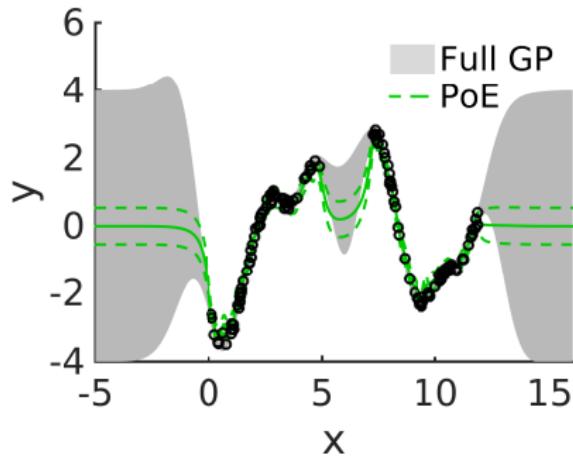
- ▶ Predictive precision (inverse variance) and mean:

$$(\sigma_*^{\text{poe}})^{-2} = \sum_k \sigma_k^{-2}(\mathbf{x}_*)$$

$$\mu_*^{\text{poe}} = (\sigma_*^{\text{poe}})^2 \sum_k \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

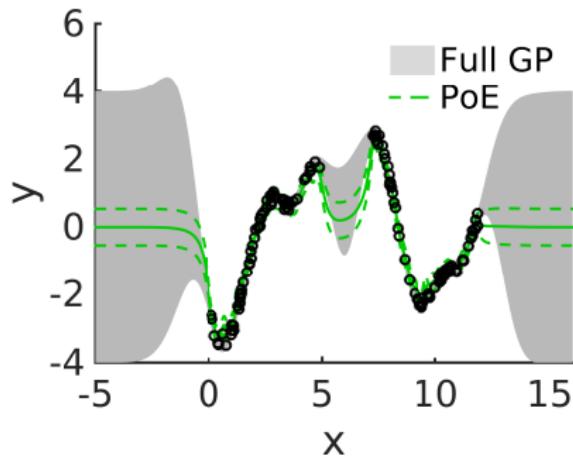
- ▶ Independent of the computational graph ✓

Product of GP Experts



- Unreasonable variances for $M > 1$:

Product of GP Experts



- Unreasonable variances for $M > 1$:

$$(\sigma_*^{\text{poe}})^{-2} = \sum_k \sigma_k^{-2}(x_*)$$

- The more experts the more certain the prediction, even if every expert itself is very uncertain \times ➡ Cannot fall back to the prior

Generalized Product of GP Experts

- Weight the responsibility of each expert in PoE with β_k

Generalized Product of GP Experts

- Weight the responsibility of each expert in PoE with β_k
- Prediction model (independent predictors):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \prod_{k=1}^M p_k^{\beta_k} (f_* | \mathbf{x}_*, \mathcal{D}^{(k)})$$

$$p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) = \mathcal{N}(f_* | \mu_k(\mathbf{x}_*), \sigma_k^2(\mathbf{x}_*))$$

Generalized Product of GP Experts

- Weight the responsibility of each expert in PoE with β_k
- Prediction model (independent predictors):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \prod_{k=1}^M p_k^{\beta_k}(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})$$

$$p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) = \mathcal{N}(f_* | \mu_k(\mathbf{x}_*), \sigma_k^2(\mathbf{x}_*))$$

- Predictive precision and mean:

$$(\sigma_*^{\text{gpoe}})^{-2} = \sum_k \beta_k \sigma_k^{-2}(\mathbf{x}_*)$$

$$\mu_*^{\text{gpoe}} = (\sigma_*^{\text{gpoe}})^2 \sum_k \beta_k \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

Generalized Product of GP Experts

- Weight the responsibility of each expert in PoE with β_k
- Prediction model (independent predictors):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \prod_{k=1}^M p_k^{\beta_k} (f_* | \mathbf{x}_*, \mathcal{D}^{(k)})$$

$$p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) = \mathcal{N}(f_* | \mu_k(\mathbf{x}_*), \sigma_k^2(\mathbf{x}_*))$$

- Predictive precision and mean:

$$(\sigma_*^{\text{gpoe}})^{-2} = \sum_k \beta_k \sigma_k^{-2}(\mathbf{x}_*)$$

$$\mu_*^{\text{gpoe}} = (\sigma_*^{\text{gpoe}})^2 \sum_k \beta_k \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

- With $\sum_k \beta_k = 1$, the model can fall back to the prior ✓
“Log-opinion pool” model (Heskes, 1998)

Generalized Product of GP Experts

- Weight the responsibility of each expert in PoE with β_k
- Prediction model (independent predictors):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \prod_{k=1}^M p_k^{\beta_k} (f_* | \mathbf{x}_*, \mathcal{D}^{(k)})$$

$$p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)}) = \mathcal{N}(f_* | \mu_k(\mathbf{x}_*), \sigma_k^2(\mathbf{x}_*))$$

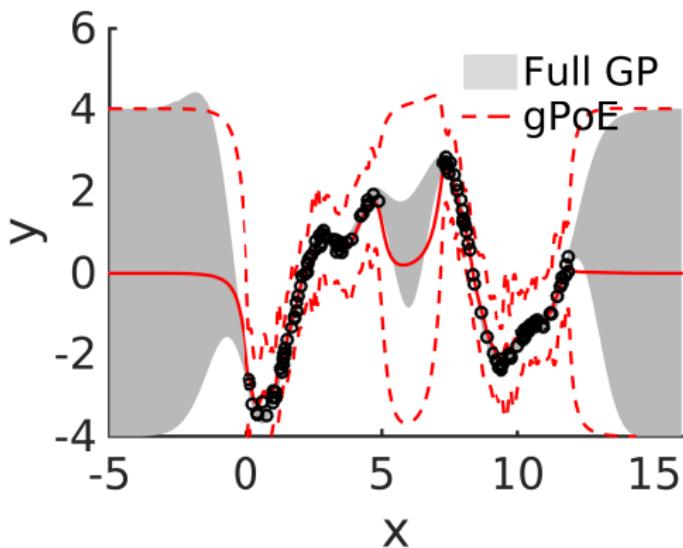
- Predictive precision and mean:

$$(\sigma_*^{\text{gpoe}})^{-2} = \sum_k \beta_k \sigma_k^{-2}(\mathbf{x}_*)$$

$$\mu_*^{\text{gpoe}} = (\sigma_*^{\text{gpoe}})^2 \sum_k \beta_k \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

- With $\sum_k \beta_k = 1$, the model can fall back to the prior ✓
“Log-opinion pool” model (Heskes, 1998)
- Independent of computational graph for $\beta_k = 1/M$ ✓

Generalized Product of GP Experts



- ▶ Same mean as PoE
- ▶ Model no longer overconfident and falls back to prior ✓
- ▶ **Very conservative variances ✗**

Bayesian Committee Machine

- Apply Bayes' theorem when combining predictions (and not only for computing predictions)

Bayesian Committee Machine

- Apply Bayes' theorem when combining predictions (and not only for computing predictions)
- Prediction model ($\mathcal{D}^{(j)} \perp\!\!\!\perp \mathcal{D}^{(k)} | f_*$):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \frac{\prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})}{p^{M-1}(f_*)}$$

Bayesian Committee Machine

- Apply Bayes' theorem when combining predictions (and not only for computing predictions)
- Prediction model $(\mathcal{D}^{(j)} \perp\!\!\!\perp \mathcal{D}^{(k)} | f_*)$:

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \frac{\prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})}{p^{M-1}(f_*)}$$

- Predictive precision and mean:

$$\begin{aligned} (\sigma_*^{\text{bcm}})^{-2} &= \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) - (M-1)\sigma_{**}^{-2} \\ \mu_*^{\text{bcm}} &= (\sigma_*^{\text{bcm}})^2 \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*) \end{aligned}$$

Bayesian Committee Machine

- Apply Bayes' theorem when combining predictions (and not only for computing predictions)
- Prediction model ($\mathcal{D}^{(j)} \perp\!\!\!\perp \mathcal{D}^{(k)} | f_*$):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \frac{\prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})}{p^{M-1}(f_*)}$$

- Predictive precision and mean:

$$\begin{aligned} (\sigma_*^{\text{bcm}})^{-2} &= \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) - (M-1)\sigma_{**}^{-2} \\ \mu_*^{\text{bcm}} &= (\sigma_*^{\text{bcm}})^2 \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*) \end{aligned}$$

- Product of GP experts, divided by $M-1$ times the prior

Bayesian Committee Machine

- Apply Bayes' theorem when combining predictions (and not only for computing predictions)
- Prediction model ($\mathcal{D}^{(j)} \perp\!\!\!\perp \mathcal{D}^{(k)} | f_*$):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \frac{\prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})}{p^{M-1}(f_*)}$$

- Predictive precision and mean:

$$(\sigma_*^{\text{bcm}})^{-2} = \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) - (M-1)\sigma_{**}^{-2}$$

$$\mu_*^{\text{bcm}} = (\sigma_*^{\text{bcm}})^2 \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

- Product of GP experts, divided by $M-1$ times the prior
- Guaranteed to fall back to the prior outside data regime ✓

Bayesian Committee Machine

- Apply Bayes' theorem when combining predictions (and not only for computing predictions)
- Prediction model ($\mathcal{D}^{(j)} \perp\!\!\!\perp \mathcal{D}^{(k)} | f_*$):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \frac{\prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})}{p^{M-1}(f_*)}$$

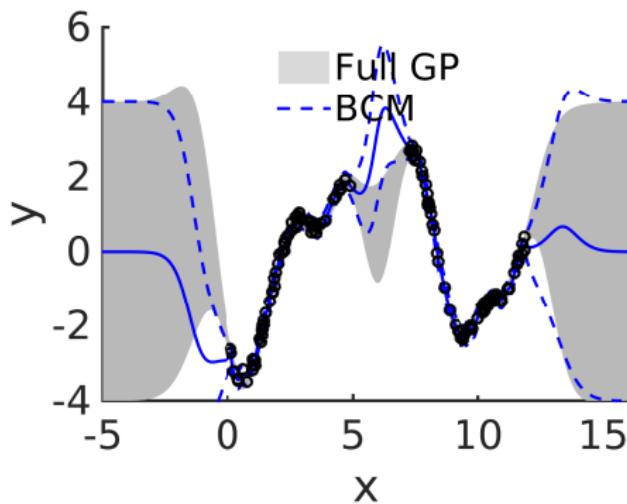
- Predictive precision and mean:

$$(\sigma_*^{\text{bcm}})^{-2} = \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) - (M-1)\sigma_{**}^{-2}$$

$$\mu_*^{\text{bcm}} = (\sigma_*^{\text{bcm}})^2 \sum_{k=1}^M \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

- Product of GP experts, divided by $M-1$ times the prior
- Guaranteed to fall back to the prior outside data regime ✓
- Independent of computational graph ✓

Bayesian Committee Machine



- Variance estimates are about right ✓
 - When leaving the data regime, the BCM can produce junk ✗
- Robustify

Robust Bayesian Committee Machine

- Merge gPoE (weighting of experts) with the BCM (Bayes' theorem when combining predictions)

Robust Bayesian Committee Machine

- Merge gPoE (weighting of experts) with the BCM (Bayes' theorem when combining predictions)
- Prediction model (conditional independence $\mathcal{D}^{(j)} \perp\!\!\!\perp \mathcal{D}^{(k)} | f_*$):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \frac{\prod_{k=1}^M p_k^{\beta_k}(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})}{p^{\sum_k \beta_k - 1}(f_*)}$$

Robust Bayesian Committee Machine

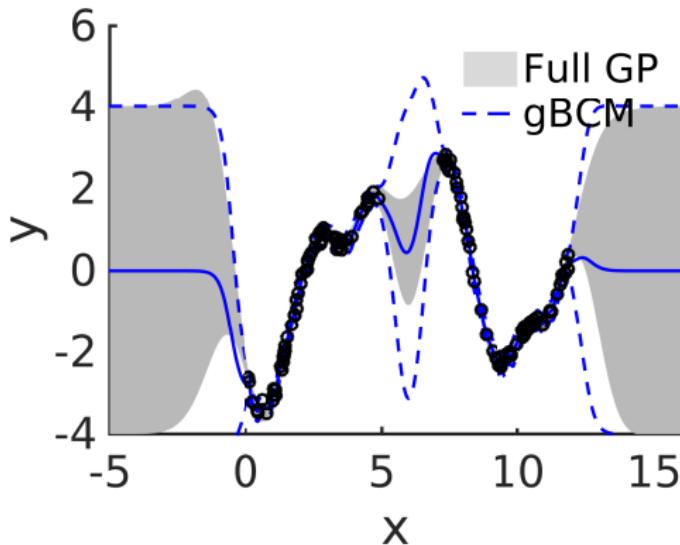
- Merge gPoE (weighting of experts) with the BCM (Bayes' theorem when combining predictions)
- Prediction model (conditional independence $\mathcal{D}^{(j)} \perp\!\!\!\perp \mathcal{D}^{(k)} | f_*$):

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = \frac{\prod_{k=1}^M p_k^{\beta_k}(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})}{p^{\sum_k \beta_k - 1}(f_*)}$$

- Predictive precision and mean:

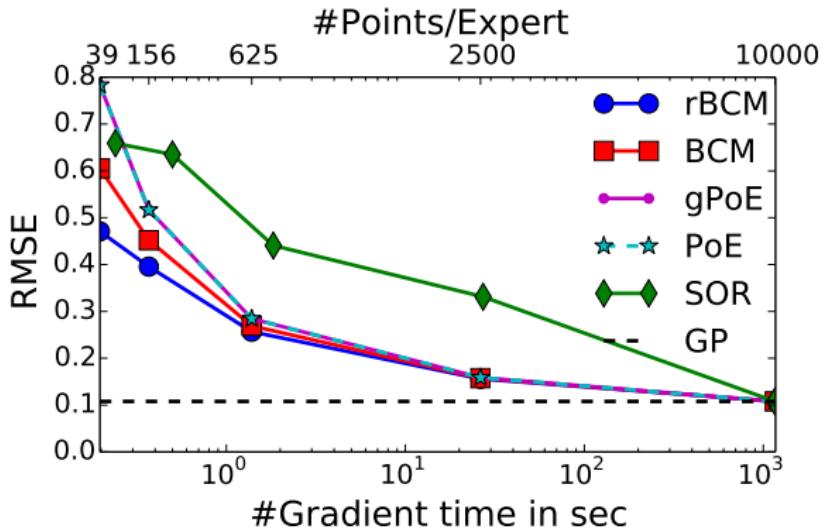
$$(\sigma_*^{\text{rbcm}})^{-2} = \sum_{k=1}^M \beta_k \sigma_k^{-2}(\mathbf{x}_*) + (1 - \sum_{k=1}^M \beta_k) \sigma_{**}^{-2},$$
$$\mu_*^{\text{rbcm}} = (\sigma_*^{\text{rbcm}})^2 \sum_k \beta_k \sigma_k^{-2}(\mathbf{x}_*) \mu_k(\mathbf{x}_*)$$

Robust Bayesian Committee Machine



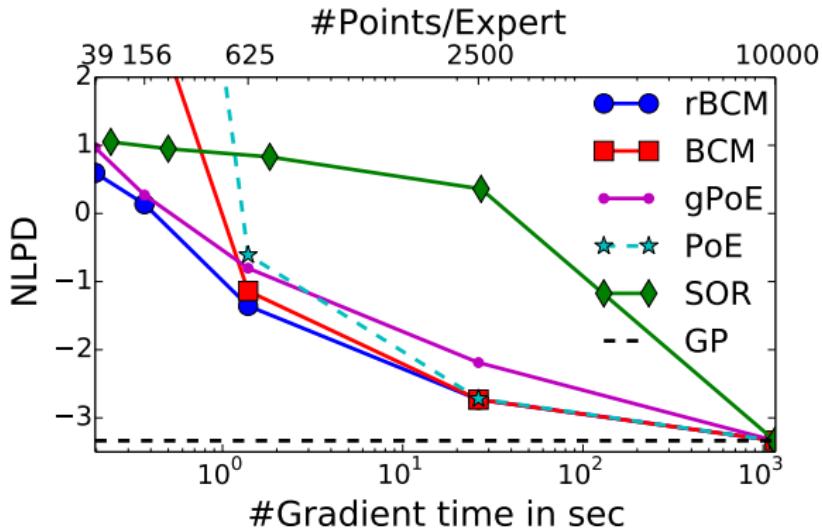
- Does not break down in case of weak experts ➤ Robustified ✓
- Robust version of BCM ➤ Reasonable predictions ✓
- Independent of computational graph (for all choices of β_k) ✓

Empirical Approximation Error



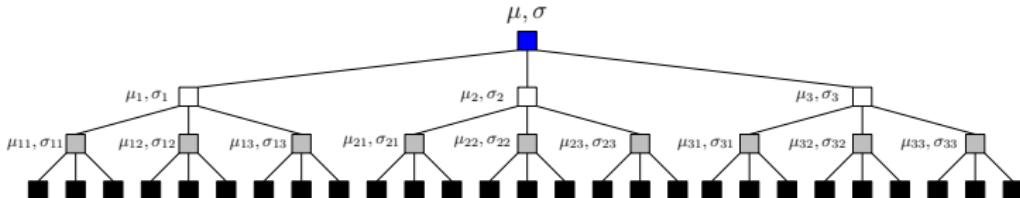
- Simulated robot arm data (10K training, 10K test)
- Hyper-parameters of ground-truth full GP
- RMSE as a function of the training time
- Sparse GP (SOR) performs worse than any distributed GP
- rBCM performs best with “weak” GP experts

Empirical Approximation Error (2)



- NLPD as a function of the training time ➤ Mean and variance
- BCM and PoE are not robust for weak experts
- gPoE suffers from too conservative variances
- rBCM consistently outperforms other methods

Summary: Distributed Gaussian Processes



- Scale Gaussian processes to large data (beyond 10^6)
- Model conceptually straightforward and easy to train
- Key: Distributed computation
- Currently tested with $N \in \mathcal{O}(10^7)$
- Scales to arbitrarily large data sets (with enough computing power)

m.deisenroth@imperial.ac.uk

Thank you for your attention

Expressiveness of the Standard Model

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

Expressiveness of the Standard Model

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

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

Expressiveness of the Standard Model

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

Two-level inference

- ▶ level-1 inference:

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

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

- ▶ level-2 inference

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

References I

- [1] E. Brochu, V. M. Cora, and N. de Freitas. A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning. Technical Report TR-2009-023, Department of Computer Science, University of British Columbia, 2009.
- [2] 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>, October 2014.
- [3] K. Chalupka, C. K. I. Williams, and I. Murray. A Framework for Evaluating Approximate Methods for Gaussian Process Regression. *Journal of Machine Learning Research*, 14:333–350, February 2013.
- [4] N. A. C. Cressie. *Statistics for Spatial Data*. Wiley-Interscience, 1993.
- [5] M. P. Deisenroth. *Efficient Reinforcement Learning using Gaussian Processes*, volume 9 of *Karlsruhe Series on Intelligent Sensor-Actuator-Systems*. KIT Scientific Publishing, November 2010. ISBN 978-3-86644-569-7.
- [6] M. P. Deisenroth, D. Fox, and C. E. Rasmussen. Gaussian Processes for Data-Efficient Learning in Robotics and Control. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):408–423, February 2015.
- [7] M. P. Deisenroth and J. Ng. Distributed Gaussian Processes. <http://arxiv.org/abs/1502.02843>, February 2015.
- [8] Y. Gal, M. van der Wilk, and C. E. Rasmussen. Distributed Variational Inference in Sparse Gaussian Process Regression and Latent Variable Models. In *Advances in Neural Information Processing Systems*. 2014.
- [9] J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian Processes for Big Data. In A. Nicholson and P. Smyth, editors, *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2013.
- [10] T. Heskes. Selecting Weighting Factors in Logarithmic Opinion Pools. In *Advances in Neural Information Processing Systems*, pages 266–272. Morgan Kaufman, 1998.
- [11] J. Kern. *Bayesian Process-Convolution Approaches to Specifying Spatial Dependence Structure*. PhD thesis, Institut of Statistics and Decision Sciences, Duke University, 2000.
- [12] 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.

References II

- [13] N. Lawrence. Probabilistic Non-linear Principal Component Analysis with Gaussian Process Latent Variable Models. *Journal of Machine Learning Research*, 6:1783–1816, November 2005.
- [14] 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.
- [15] J. Ng and M. P. Deisenroth. Hierarchical Mixture-of-Experts Model for Large-Scale Gaussian Process Regression. <http://arxiv.org/abs/1412.3078>, December 2014.
- [16] 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.
- [17] 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.
- [18] 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.
- [19] 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.
- [20] M. K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2009.
- [21] V. Tresp. A Bayesian Committee Machine. *Neural Computation*, 12(11):2719–2741, 2000.