**Imperial College
London**

# Gaussian Processes

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

**Marc Deisenroth**

Department of Computing
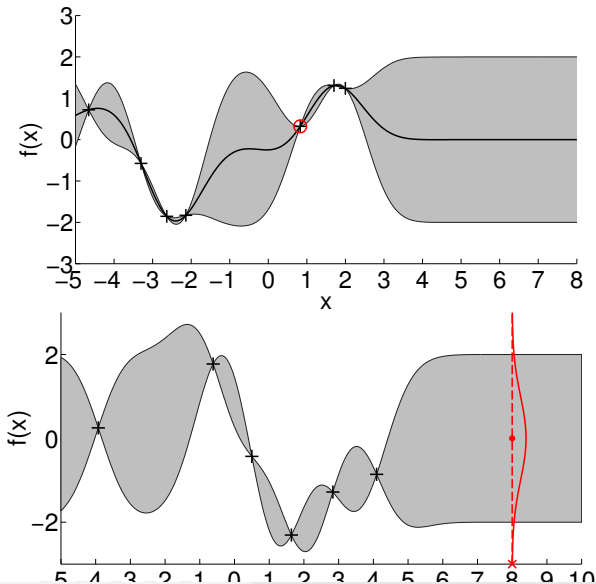Imperial College London

February 13, 2017

Gaussian Processes for Machine Learning

Carl Edward Rasmussen and Christopher K. I. Williams

http://www.gaussianprocess.org/

# Problem Setting

# Recap from CO-496: Bayesian Linear Regression

‣ Linear Regression Model:

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

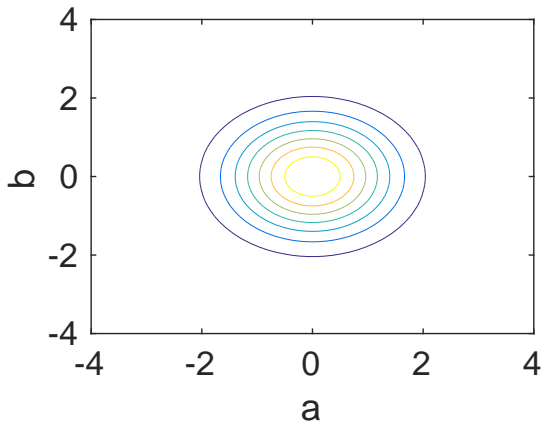‣ Integrating out the parameters when predicting leads to a distribution over functions:

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

# Sampling from the Prior over Functions

Consider a linear regression setting

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

# Sampling from the Prior over Functions

Consider a linear regression setting

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

# Sampling from the Prior over Functions
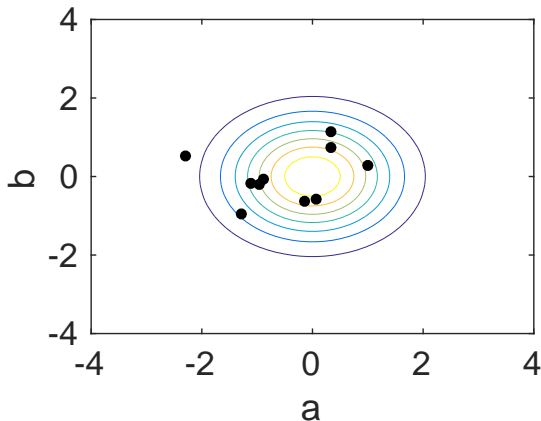
Consider a linear regression setting

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

# Sampling from the Prior over Functions
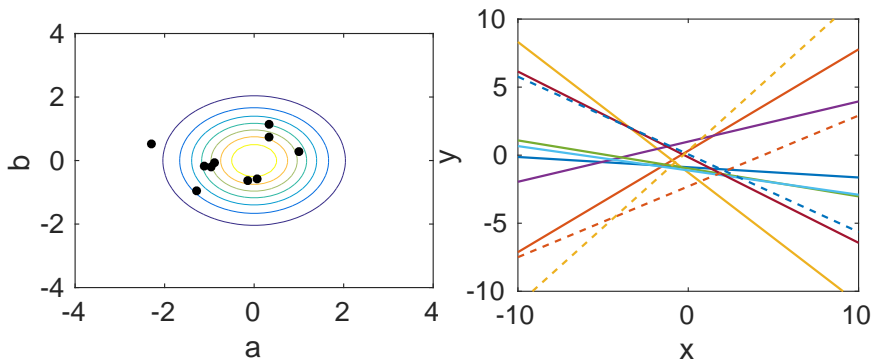
Consider a linear regression setting

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

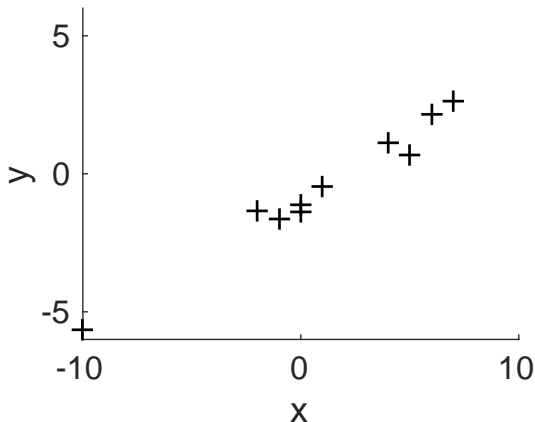# Sampling from the Posterior over Functions

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

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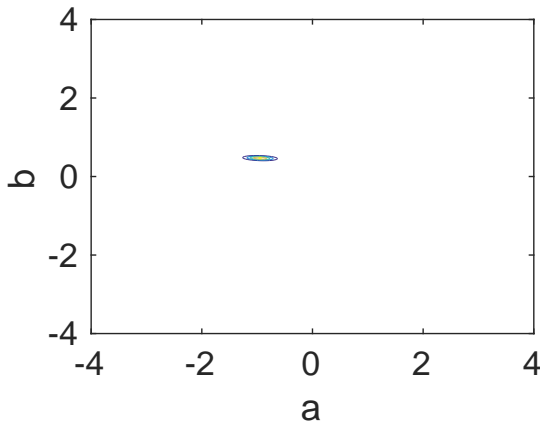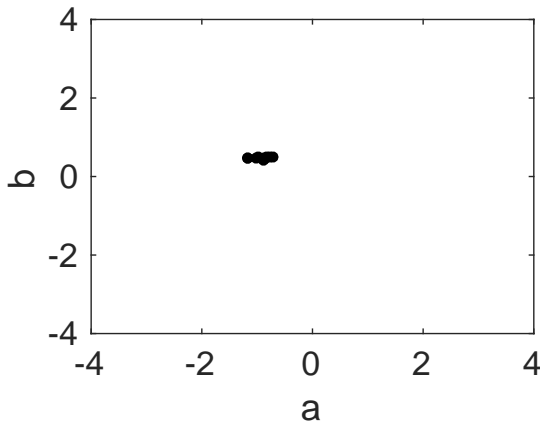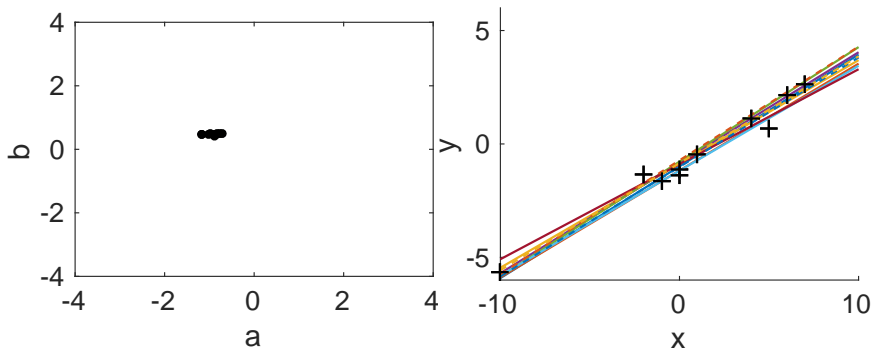$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

# 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

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

where

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

for given "centers" $\boldsymbol{\mu}_i$

# Illustration: Fitting a Radial Basis Function Network

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



‣ Place Gaussian-shaped basis functions $\phi_i$ at 25 input locations $\mu_i$, linearly spaced in the interval $[-5, 3]$

# Samples from the RBF Prior

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

# Samples from the RBF Posterior

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

# 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 $f = [f_1, f_2, f_3, ...]$
- 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 $f_1, f_2, \ldots$, any finite number of which is Gaussian distributed.

- A Gaussian distribution is specified by a mean vector $\mu$ and a covariance matrix $\Sigma$
- 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:

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

# GP Regression as a Bayesian Inference Problem

## Objective

For a set of observations $y_i = f(x_i) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$, find a (posterior) distribution over functions $p(f|X, y)$ 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): $p(y|f, X) = \mathcal{N}(f(X), \sigma_n^2 I)$

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

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

# 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 $f_i = f(\boldsymbol{x}_i)$

‣ Consider a finite number of $N$ function values $\boldsymbol{f}$ and all other (infinitely many) function values $\tilde{\boldsymbol{f}}$. Informally:

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

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

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

‣ Key property: The marginal remains finite

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

# Training and Test Marginal

‣ In practice, we always have finite training and test inputs $x_{\text{train}}, x_{\text{test}}$.

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

‣ Then, we obtain the finite marginal

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

# GP Regression as a Bayesian Inference Problem (ctd.)

Posterior over functions (with training data $X, y$):

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

Using the properties of Gaussians, we obtain

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

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

$$K = k(X, X)$$

Marginal likelihood:

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

# GP Predictions (1)

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

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

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

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

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

# GP Predictions (2)

Prior:

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

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

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

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

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

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

# Illustration: Inference with Gaussian Processes



Prior belief about the function

Predictive (marginal) mean and variance:

$$
\begin{aligned}
\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{x}_*, \varnothing] &= m(\boldsymbol{x}_*) = 0 \\
\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{x}_*, \varnothing] &= \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*)
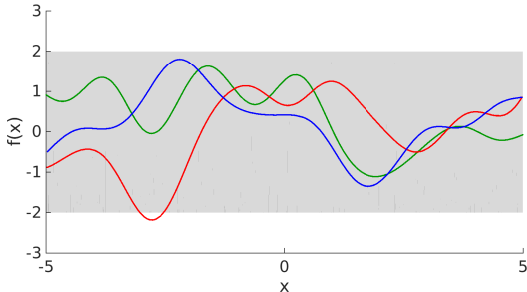\end{aligned}
$$

# 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

$$k_{Gauss}(\boldsymbol{x}_i, \boldsymbol{x}_j) = \sigma_f^2 \exp\left(-(\boldsymbol{x}_i - \boldsymbol{x}_j)^\top(\boldsymbol{x}_i - \boldsymbol{x}_j)/\ell^2\right)$$
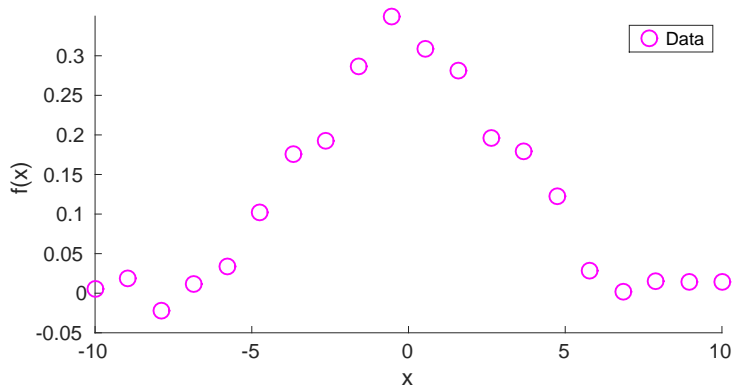
- $\sigma_f$: Amplitude of the latent function
- $\ell$: Length scale. How far do we have to move in input space before the function value changes significantly

  ▶▶ **Smoothness parameter**



- Assumption on latent function: Smooth ($\infty$ differentiable)
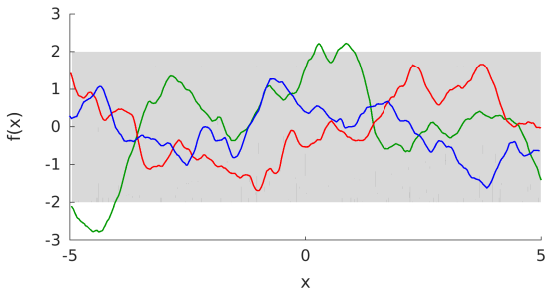
# Length-Scales

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

# Matérn Covariance Function

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

- $\sigma_f$: Amplitude of the latent function
- $\ell$: Length scale. How far do we have to move in input space before the function value changes significantly?
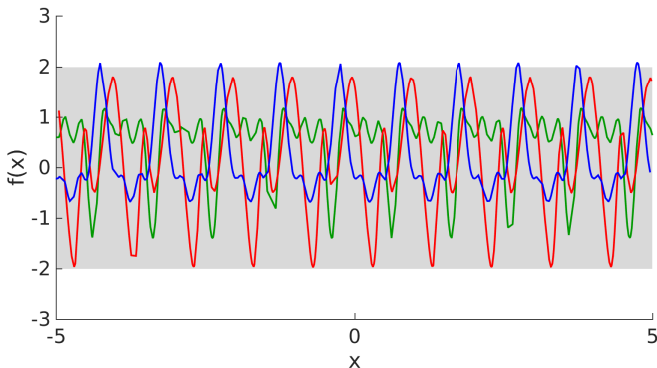


- Assumption on latent function: 1-times differentiable

# Periodic Covariance Function

$$k_{per}(x_i, x_j) = \sigma_f^2 \exp\left(-\frac{2\sin^2\left(\frac{\kappa(x_i - x_j)}{2\pi}\right)}{\ell^2}\right)$$

$$= k_{Gauss}(\boldsymbol{u}(x_i), \boldsymbol{u}(x_j)), \quad \boldsymbol{u}(x) = \begin{bmatrix} \cos(\kappa x) \\ \sin(\kappa x) \end{bmatrix}$$

$\kappa$: 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 $\sigma_n^2$)

▶▶ 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

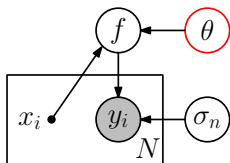Find good GP hyper-parameters $\boldsymbol{\theta}$ (kernel and mean function parameters)

- Place a prior $p(\boldsymbol{\theta})$ on hyper-parameters

- Posterior over hyper-parameters:

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

- Choose hyper-parameters $\boldsymbol{\theta}^*$, such that

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

▶▶ Maximize marginal likelihood if $p(\boldsymbol{\theta}) = \mathcal{U}$ (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:

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

Learning the GP hyper-parameters:

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

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

# Training via Marginal Likelihood Maximization
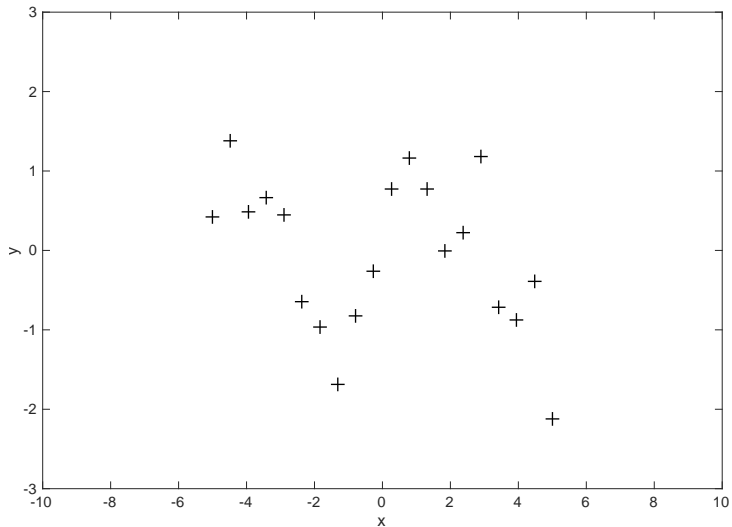
Log-marginal likelihood:

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

▸ Automatic trade-off between data fit and model complexity

▸ Gradient-based optimization of hyper-parameters $\boldsymbol{\theta}$:

$$\frac{\partial \log p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{\theta})}{\partial \theta_i} = \tfrac{1}{2}\boldsymbol{y}^{\top}\boldsymbol{K}_{\boldsymbol{\theta}}^{-1}\frac{\partial \boldsymbol{K}_{\boldsymbol{\theta}}}{\partial \theta_i}\boldsymbol{K}_{\boldsymbol{\theta}}^{-1}\boldsymbol{y} - \tfrac{1}{2}\text{tr}\big(\boldsymbol{K}_{\boldsymbol{\theta}}^{-1}\frac{\partial \boldsymbol{K}_{\boldsymbol{\theta}}}{\partial \theta_i}\big)$$

$$= \tfrac{1}{2}\text{tr}\big((\boldsymbol{\alpha}\boldsymbol{\alpha}^{\top} - \boldsymbol{K}_{\boldsymbol{\theta}}^{-1})\frac{\partial \boldsymbol{K}_{\boldsymbol{\theta}}}{\partial \theta_i}\big),$$

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

# Example: Training Data

# Example: Marginal Likelihood Contour



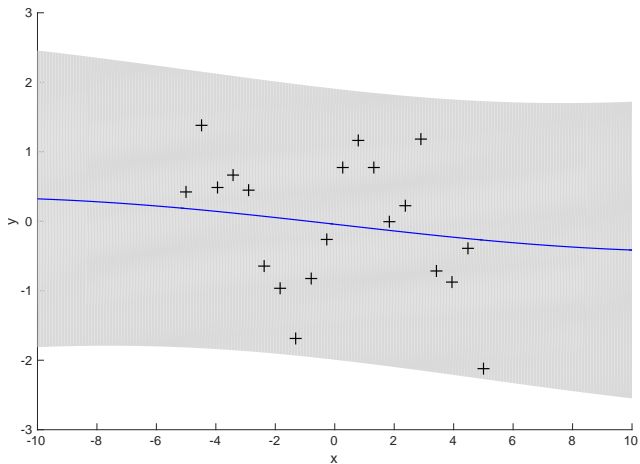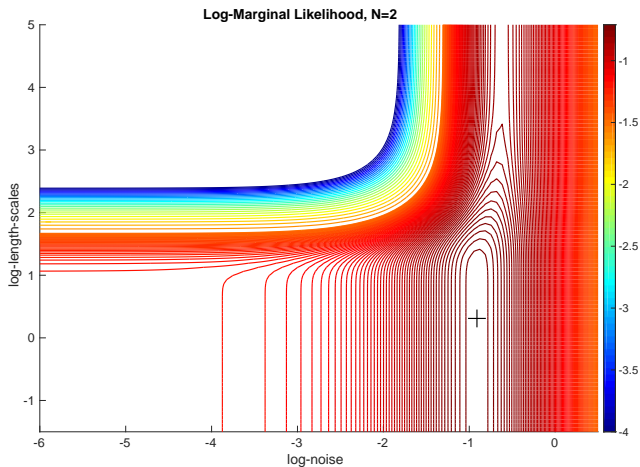Log-Marginal Likelihood, N=20

# Example: Exploring the Modes (1)

# Example: Exploring the Modes (2)

# Marginal Likelihood (1)
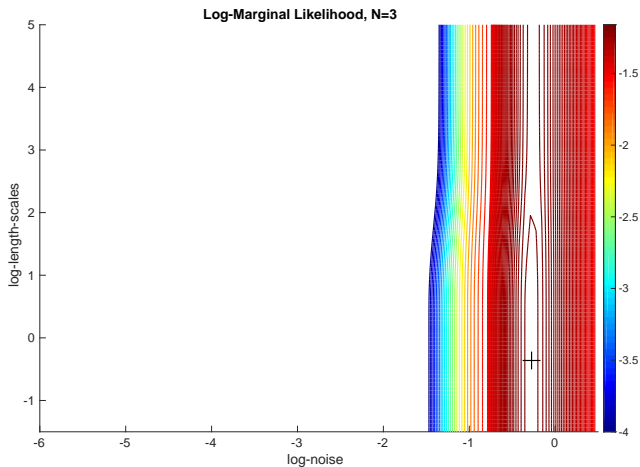


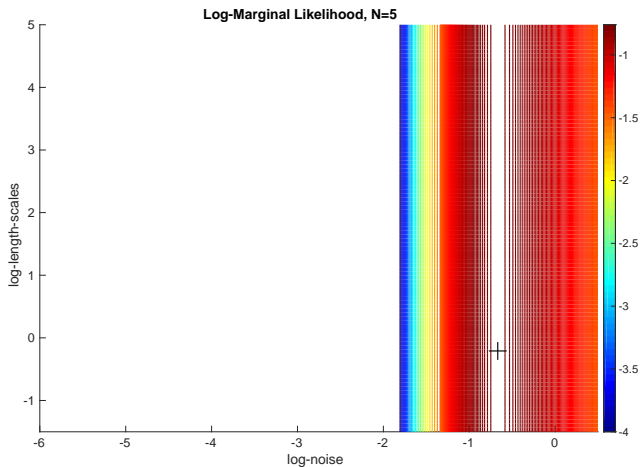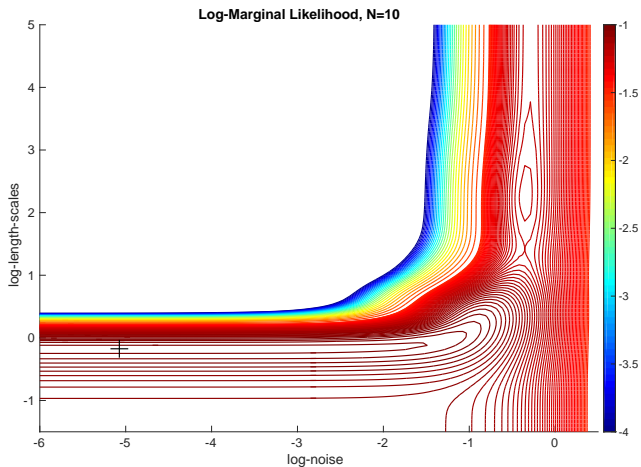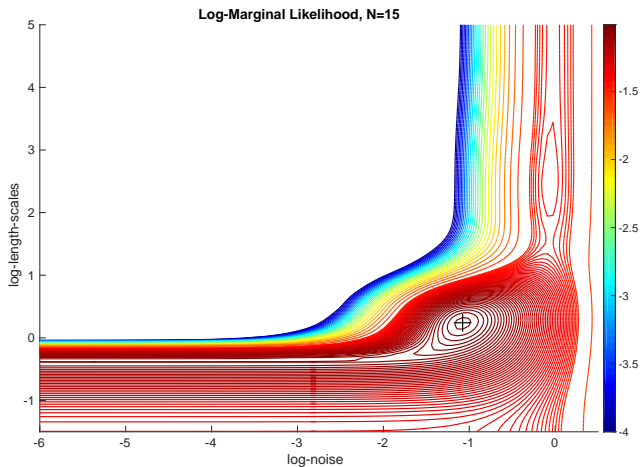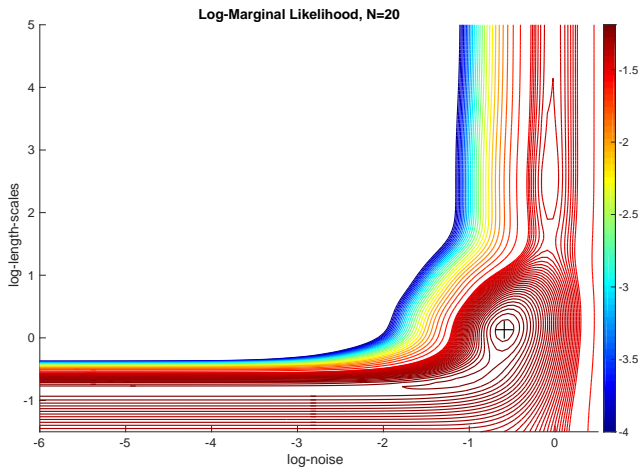**Log-Marginal Likelihood, N=2**

# Marginal Likelihood (2)



Log-Marginal Likelihood, N=3

# Marginal Likelihood (3)

# Marginal Likelihood (4)



**Log-Marginal Likelihood, N=10**
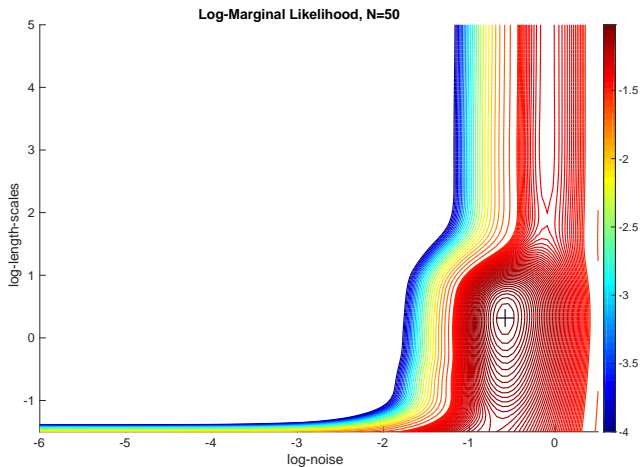
# Marginal Likelihood (5)



Log-Marginal Likelihood, N=15

# Marginal Likelihood (6)



Log-Marginal Likelihood, N=20

# Marginal Likelihood (7)
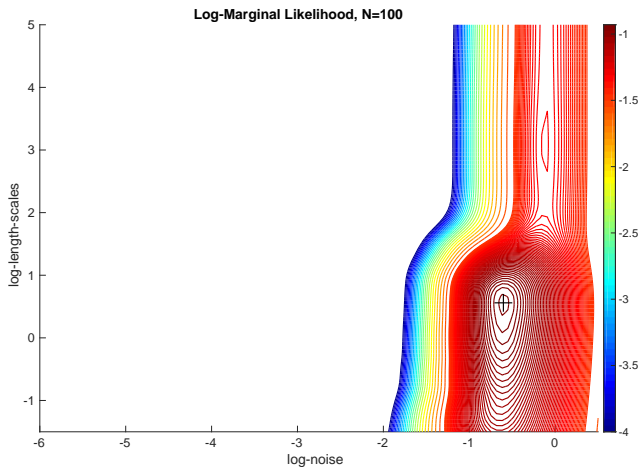


Log-Marginal Likelihood, N=50

# Marginal Likelihood (8)
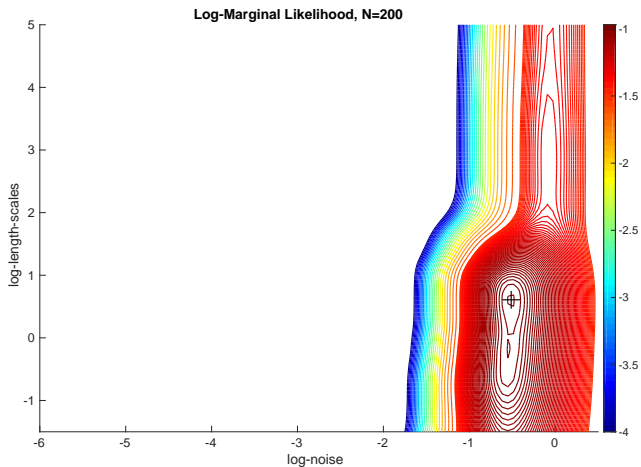


Log-Marginal Likelihood, N=100
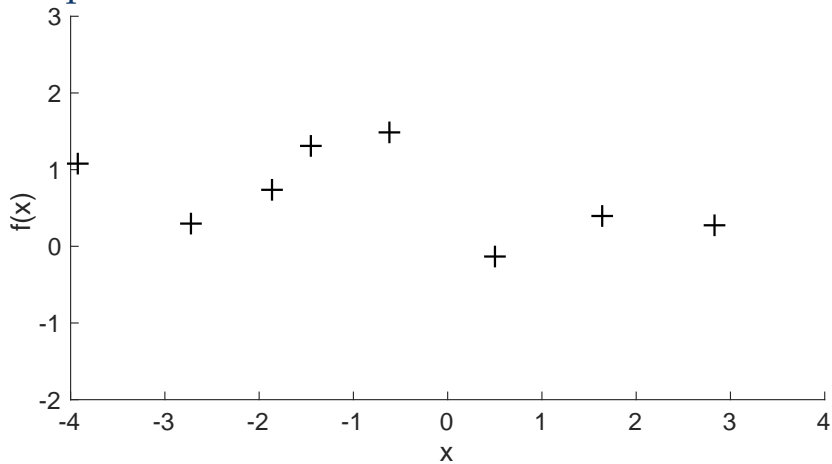
# Marginal Likelihood (9)



Log-Marginal Likelihood, N=200

# 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
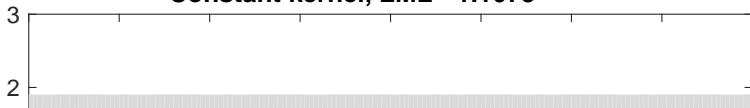  Why can we do not do this easily?

# 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?
‣ Some options:
    ‣ BIC, AIC (see CO-496)
    ‣ Compare marginal likelihood values (assuming a uniform prior on
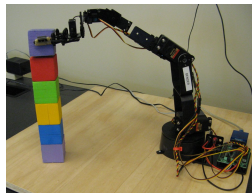      the set of models)
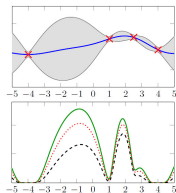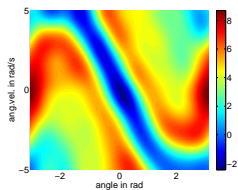
# Example

# 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 $\ell$ to $\approx 0.5$.

‣ 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 $K + \sigma_n^2 I$) by penalizing high signal-to-noise ratios $\sigma_f / \sigma_n$

**Appendix**

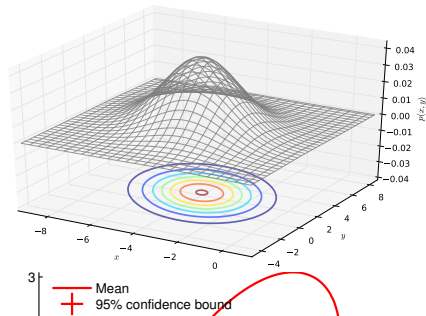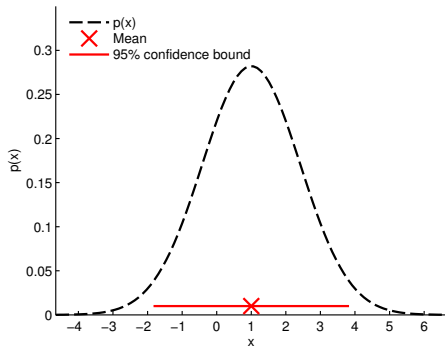# The Gaussian Distribution

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

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

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

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

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

1. Find conditions for $A, b$ to match the mean of $y$
2. Find conditions for $A, b$ to match the covariance of $y$

# 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 $\Sigma$ and mean vector $\mu$.

```
x = randn(D,1);            Sample x ~ N(0, I)
y = chol(Σ)'*x + μ;        Scale x and add offset
```

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

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

# Conditional



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

# Marginal



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

Marginal distribution:

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

‣ 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 $p(\boldsymbol{x}, \tilde{\boldsymbol{x}})$, where $\boldsymbol{x} \in \mathbb{R}^D$ and $\tilde{\boldsymbol{x}} \in \mathbb{R}^k, k \to \infty$ are random variables.
Then

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

where $\boldsymbol{\Sigma}_{\tilde{x}\tilde{x}} \in \mathbb{R}^{k \times k}$ and $\boldsymbol{\Sigma}_{x\tilde{x}} \in \mathbb{R}^{D \times k}$, $k \to \infty$.
However, the marginal remains finite

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

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 $x_{\text{train}}, x_{\text{test}}$
- Then, $x = \{x_{\text{train}}, x_{\text{test}}, x_{\text{other}}\}$
  ($x_{\text{other}}$ plays the role of $\tilde{x}$ from previous slide)

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

$$p(x_{\text{train}}, x_{\text{test}}) = \int p(x_{\text{train}}, x_{\text{test}}, x_{\text{other}}) d\, x_{\text{other}}$$

$$p(x_{\text{test}}|x_{\text{train}}) = \mathcal{N}(\mu_*, \Sigma_*)$$

$$\mu_* = \mu_{\text{test}} + \Sigma_{\text{test,train}} \Sigma_{\text{train}}^{-1}(x_{\text{train}} - \mu_{\text{train}})$$

$$\Sigma_* = \Sigma_{\text{test}} - \Sigma_{\text{test,train}} \Sigma_{\text{train}}^{-1} \Sigma_{\text{train,test}}$$

# Gaussian Process Training: Hierarchical Inference

‣ Level-1 inference (posterior on $f$):

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

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

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

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

# GP as the Limit of an Infinite RBF Network

Consider the universal function approximator

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

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 $\theta_1^2$

▶▶ GP with mean 0 and Gaussian covariance function

# 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] N. A. C. Cressie. *Statistics for Spatial Data*. Wiley-Interscience, 1993.

[3] 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.

[4] M. P. Deisenroth and S. Mohamed. Expectation Propagation in Gaussian Process Dynamical Systems. In *Advances in Neural Information Processing Systems*, pages 2618–2626, 2012.

[5] M. P. Deisenroth and J. W. Ng. Distributed Gaussian Processes. In *Proceedings of the International Conference on Machine Learning*, 2015.

[6] M. P. Deisenroth, C. E. Rasmussen, and J. Peters. Gaussian Process Dynamic Programming. *Neurocomputing*, 72(7–9):1508–1524, March 2009.

[7] 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.

[8] 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.

[9] 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.

[10] 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.

[11] N. Lawrence. Probabilistic Non-linear Principal Component Analysis with Gaussian Process Latent Variable Models. *Journal of Machine Learning Research*, 6:1783–1816, November 2005.

# References II

[12] 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.

[13] M. A. Osborne, R. Garnett, and S. J. Roberts. Gaussian Processes for Global Optimization. In *Proceedings of the International Conference on Learning and Intelligent Optimization*, 2009.

[14] 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.

[15] 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.

[16] 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.

[17] S. Roberts, M. A. Osborne, M. Ebden, S. Reece, N. Gibson, and S. Aigrain. Gaussian Processes for Time Series Modelling. *Philosophical Transactions of the Royal Society (Part A)*, 371(1984), February 2013.

[18] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. de Freitas. Taking the Human out of the Loop: A Review of Bayesian Optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.