

Machine Learning in Robotics

Gaussian Processes

Prof. Dongheui Lee

Institute of Automatic Control Engineering
Technische Universität München

dhlee@tum.de

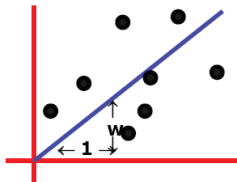
Today's Lecture Outline

- Gaussian Processes: Introduction
- Gaussian Processes for Regression
- Hyper-parameters and their Optimal Selection
- Kernel
- Examples in Robotics

Recall: Linear Regression

Linear regression assumes that expected value of the output given an input is linear.

$$y^{(i)} = wx^{(i)} + \epsilon \quad (1)$$



We just assumed the underlying function $f(x)$ is *linear*. What if the function is cubic, or non-polynomial? Furthermore, if such class model is not known a priori, how can we formulate a hypothesis over the data and maintain general considerations?

Gaussian Processes: definition

Gaussian is widely used to describe **a distribution over vector**:

$$y \sim \mathcal{N}(\mu, \Sigma)$$

Similarly, a Gaussian Process defines **a distribution over functions** $p(f(\mathbf{x}))$, where $f(\mathbf{x})$ is a scalar function mapping some input space \mathcal{X} to \mathbb{R} .

$$f : \mathcal{X} \rightarrow \mathbb{R}$$

where $f(\mathbf{x})$ is a random variable and that it can be an infinite-dimensional quantity.

$$f \sim \mathcal{GP}(m, K)$$

Definition: *A stochastic process is a Gaussian process if for any finite subset $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\} \subset \mathcal{X}$, the marginal distribution $p(f(\mathbf{x}))$ over that finite subset has a multivariate Gaussian distribution*

Gaussian Processes

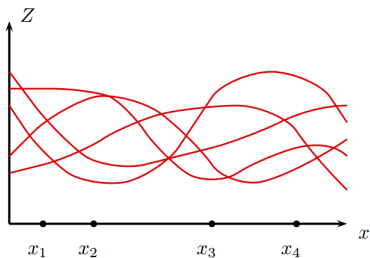
A GP is a Gaussian distribution over functions:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

A GP is parameterized by a mean function $m(\mathbf{x})$, and a covariance function (kernel function $k(\mathbf{x}, \mathbf{x}')$)

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))^T]$$



Gaussian Processes in Practice

A training dataset is given $(\mathbf{x}^{(i)}, f(\mathbf{x}^{(i)}))$ for $i = 1, \dots, n$.

$$f(\mathbf{X}) \sim \mathcal{N}(m(\mathbf{X}), K(\mathbf{X}, \mathbf{X}))$$

$$m(\mathbf{X}) = (m(\mathbf{x}^{(1)}), \dots, m(\mathbf{x}^{(n)}))^T$$

$$K(\mathbf{X}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(n)}) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(n)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(n)}, \mathbf{x}^{(n)}) \end{bmatrix}$$

Kernel Functions

GP use kernels to lift input data in a higher dimensional feature space (*kernel trick*)

- The kernel function $k(\mathbf{x}, \mathbf{x}')$ is a measure of “closeness” between inputs: Close points \rightarrow similar predictions
- The kernel function has to generate a symmetric positive semi-definite covariance matrix for any set of points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$

Common Kernel Functions

- Linear: $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$.
- Squared exponential (SE): $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \sigma_f^2 \exp^{-\frac{1}{2l} \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}$.
- Radial basis function (RBF):

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \sigma_f^2 \exp^{-\frac{1}{2} \sum_{l=1}^d w_l (x_l^{(i)} - x_l^{(j)})^2}, \text{ where } x_l^{(i)} \text{ is the } l\text{-th component of } \mathbf{x}^{(i)}.$$

Gaussian Process Regression

We have a training dataset (X and corresponding observations $f = f(X)$). Given a test set X_* , we want to predict the function output f_* .

The GP defines a joint distribution for $p(f, f_* | X, X_*)$.

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu \\ \mu_* \end{bmatrix}, \begin{bmatrix} K & K_* \\ K_*^T & K_{**} \end{bmatrix} \right)$$

where $\mu = m(X)$, $\mu_* = m(X_*)$, $K = K(X, X)$, $K_* = K(X, X_*)$, $K_{**} = K(X_*, X_*)$.

To infer f_* (or $p(f_* | X, X_*, f)$), the rule for conditional multivariate Gaussians is applied.

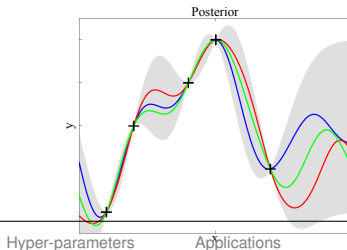
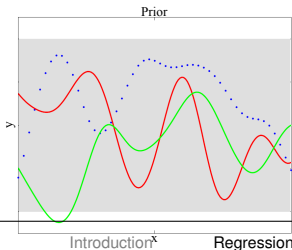
Noiseless GP regression

- **Training Data:** $\mathcal{D} = \{x^{(i)}, f^{(i)} = f(x^{(i)})\}_{i=1}^n$
- **Prior:**

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu \\ \mu_* \end{bmatrix}, \begin{bmatrix} K & K_* \\ K_*^T & K_{**} \end{bmatrix} \right)$$

- **Prediction:** Given a test input x_t ,

$$\begin{aligned} p(f_* | X, X_*, f) &= \mathcal{N}(f_* | \mu^*, \Sigma^*) \\ \mu^* &= \mu_* + K_*^T K^{-1} (f - \mu) \\ \Sigma^* &= K_{**} - K_*^T K^{-1} K_* \end{aligned}$$



Noisy GP regression

- **Training Data:** $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)} = f(\mathbf{x}^{(i)}) + \epsilon\}_{i=1}^n = (\mathbf{X}, \mathbf{y})$

$$f \sim \mathcal{GP}(0, k(\cdot, \cdot))$$

$$\epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

- **Prior:** Given a test input \mathbf{x}_t

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{bmatrix}\right)$$

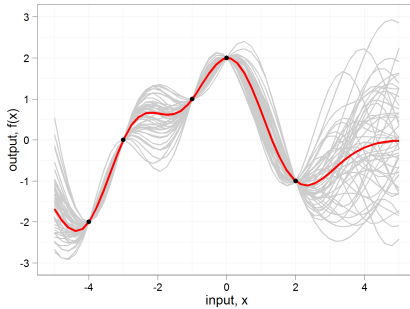
- **Predictive distribution:**

$$p(\mathbf{f}_* | \mathbf{X}, \mathbf{X}_*, \mathbf{y}) = \mathcal{N}(\mathbf{f}_* | \boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

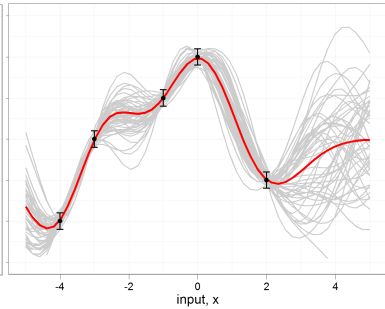
$$\boldsymbol{\mu}^* = \boldsymbol{\mu}_* + \mathbf{K}_*^T [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma}^* = \mathbf{K}_{**} - \mathbf{K}_*^T [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{K}_*$$

Regression Examples



without noise



with noise

The role of the Hyper-parameters

- Kernels are functions of a set of *hyper-parameters* θ
 - Example: Squared Exponential kernel
$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \sigma_f^2 \exp^{-\frac{1}{2l} \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}$$
$$\theta = [\sigma_f, l]$$
- GP regression results strongly depends on the used hyper-parameters

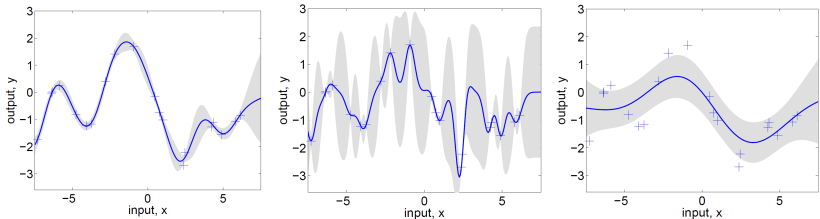


Figure: GP with different hyperparameters. (left) $l = 1$, (middle) $l = 0.3$, (right) $l = 3$

Learning the Hyper-parameters

Hyper-parameters can be learned from observations, which makes GP a parameter-free approach

HOW: Maximize the marginal log-likelihood conditioned on the hyper-parameters θ

1. $p(\mathbf{y}|\mathbf{X}, \theta) = \mathcal{N}(0, \mathbf{K}_y) = \mathcal{N}(0, \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I})$
2. $\log p(\mathbf{y}|\mathbf{X}, \theta) = -\frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{n}{2} \log 2\pi$
3. $\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|\mathbf{X}, \theta) = 0 \leftarrow$ Gradient descent

The terms of the marginal likelihood have a clear meaning:

- $-\mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} / 2$ represents the data-fit
- $-\log |\mathbf{K}_y| / 2$ is the complexity penalty
- $-n \log 2\pi / 2$ is a normalization constant

Properties of GPs

- GP is a nonparametric model.
- GPs are universal approximators. (GP can fit arbitrary functions).
- Probabilistic model: it returns measure of uncertainties
- Minimal assumptions over a distribution of functions (i.e., joint Normal)
- GP for regression and classification.
- Kernel functions: Design of kernel functions for application, Hyperparameters can be learned.

Robotics Applications: PILCO

- Model-free reinforcement learning (RL) algorithms usually requires several rollouts to learn a suitable control policy
 - ▶ Extremely time consuming on real robots
- PILCO (Probabilistic Inference for Learning Control)
 - ▶ Model-based RL
 - ▶ GP for long-term predictions and policy evaluation
 - ▶ Reduced “real” experience to learn the control policy
- PILCO - Working principle
 1. Perform a rollout (real robot) and learn a model of the robot (GP)
 2. Find an optimal policy for the learned model
 3. Apply the policy to the robot
 4. Finish if task learned, otherwise repeat from 2.



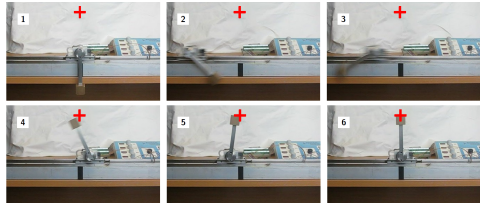
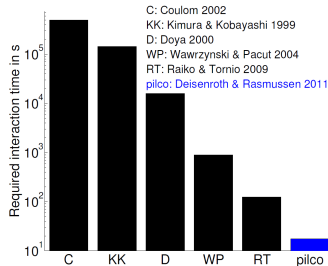
Deisenroth et. al., *Gaussian Processes for Data-Efficient Learning in Robotics and Control*, TPAMI, 2015.



Robotics Applications: PILCO

Task: Swing up a pendulum attached to a cart (cart-pole system, i.e. a cart running on a track and a freely swinging pendulum attached to the cart)

Result: PILCO is able to learn the task significantly faster than state-of-the-art model-based RL approaches



Robotics Applications: LMDS

Locally Modulated Dynamical Systems (LMDS): incrementally refine a set of predefined robotic tasks

- Predefined tasks are stable dynamical systems $\dot{x} = f(x)$, where x is the robot position
- The user can provide kinesthetic demonstrations to modify the task execution
 - ▶ Learn (GP) a matrix $M(x) = (1 - \alpha(x))R(x)$ from user demonstrations
 - ▶ Modify the original task $\dot{x} = M(x)f(x)$ (do not affect the equilibrium points)



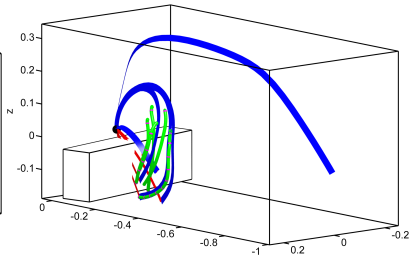
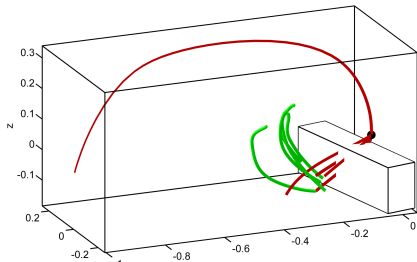
Kronander et. al., *Incremental Motion Learning with Locally Modulated Dynamical Systems*, RAS, 2015.



Robotics Applications: LMDS

Task: insert plates into slots in a dish rack.

- The original dynamical achieve the task from different starting positions
- The user provides kinesthetic demonstrations to prevent the robot to hit the rack and the task is correctly executed



Reference

- Rasmussen and Williams, Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning) - Chap. 1, 2, 3