

RL for Statisticians: an Incomplete Introduction to Model-Based  
Reinforcement Learning

Beatrice Cantoni

Spring 2022

Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>RL; Model-based versus Model-free</b>	<b>2</b>
<b>3</b>	<b>PILCO</b>	<b>3</b>
3.1	Notation and setting . . . . .	3
3.2	Modeling the underlying dynamics using a GP . . . . .	4
3.3	From the GP to the expected return . . . . .	5
3.4	From the expected return to the optimal policy . . . . .	6
<b>4</b>	<b>Mountain car implementation</b>	<b>6</b>
<b>5</b>	<b>Further applications and role in the literature</b>	<b>8</b>
<b>6</b>	<b>Final comments</b>	<b>8</b>

# 1 Introduction

This report was written as final project for the class SDS383D in Spring 2022. Its aim is to review and understand one of the most well known approaches in model-based Reinforcement Learning (RL) field, named PILCO by its inventors [1]. The report is designed to be useful for people with graduate-level statistical knowledge but without deep knowledge of RL literature and techniques. A public library, available online, was also used for simulations. Some difficulties encountered in dealing with the necessary packages and environments will be mentioned. A drawback encountered was the unavailability of certain packages for Windows, as well as some recent updates of important packages such as Tensorflow.

## 2 RL; Model-based versus Model-free

Reinforcement learning [5] is a class of learning problems in which an agent (or controller) interacts with a stochastic dynamic and incompletely known environment with the goal of finding an action-selection strategy or policy to optimize some measure of its long-term performance. The interaction is often modeled conventionally as a *Markov decision process* (MDP). In particular, RL proposes that good policies are inferred by observing rewards (dependent on the actions taken) and adjusting the associated policy that maps states to actions and thence to future rewards. This approach is loosely seen as lying between supervised and unsupervised learning. One big classification within RL models is the distinction between value function and policy search models. Firstly, value function methods approximate a function that maps states (or state-action pairs) to expected returns. Secondly, policy search approaches use a parametrised policy and, by sampling trajectories from the system, estimate the gradient of a reward (or cost) function with respect to the policy parameters and propose a new value.

**Model-based vs Model-free** Traditionally, RL algorithms have been categorized as being either model-based or model-free. In the former category, the agent uses the collected data to first build a model of the domain's dynamics and then uses this model to optimize its policy. In the latter case, the agent directly learns an optimal (or good) action-selection strategy from the collected data. There is some evidence that model-free methods are usually more flexible, whereas model-based methods offer more data efficiency. Generally speaking, model-based methods are more promising to efficiently extract valuable information from available data than model-free methods. On the other hand, model based methods are not widely used in learning from scratch since they suffer from model bias, i.e., they inherently assume that the learned dynamics model sufficiently resembles the real environment. Model bias is especially an issue when only a few samples and no informative prior knowledge about the task to be learned is available.

### 3 PILCO

Deisenroth and Rasmussen [1] proposed PILCO, a **model-based, policy search** method based on a probabilistic approach for learning the underlying dynamics of the model. **Its main advantage is the reduced model bias**, which is one of the key problems of model-based reinforcement learning. The key for such an improvement lies in the probabilistic approach itself, since it allows for explicitly incorporating model uncertainty into long-term planning. When such an intrinsic uncertainty is not in place, approximations of  $p(\mathbf{x}_t)$  and the associated long term predictions risk being essentially arbitrary. This is even more harmful when few data are available. By contrast, a probabilistic function approximator, obtained using a Gaussian Process, as we will see, places a posterior distribution over the transition function, thereby expressing the level of uncertainty about the model.

#### 3.1 Notation and setting

The key elements are:

- A dynamic system

$$\mathbf{x}_t = f(\mathbf{x}_{t-1}, \mathbf{u}_{t-1}),$$

where  $\mathbf{x} \in \mathbb{R}^D$  is a continuous-valued state,  $\mathbf{u}$  are the controls, and  $f$  encodes the unknown transition dynamic. For  $\mathbf{x}_0$ , we are assuming

$$\mathbf{x}_0 \sim N(\mu_0, \Sigma_0)$$

- A *policy/controller*  $\pi$

$$\mathbf{x} \mapsto \pi(\mathbf{x}) = \mathbf{u}$$

is a function parametrised by  $\theta$ .

- The expected return of following a certain  $\pi$  for  $T$  steps is

$$J^\pi(\theta) = \sum_{t=0}^T E_{\mathbf{x}_t}[c(\mathbf{x}_t)],$$

where  $c(\mathbf{x})$  is the cost of being in state  $\mathbf{x}$  at time  $t$ , which encodes some information about a target state  $\mathbf{x}_{target}$ .

The goal is to find a deterministic *policy/controller*  $\pi$  that minimises the *expected return*.

**Examples** For a better understanding of this last element of the PILCO strategy, I will give two examples, coming from the applications mentioned by the authors in the paper.

**Robot Unicycle** One of the implementations described by the authors is the one of a robotic unicycle system. The robot consists of a wheel, a frame, and a flywheel mounted perpendicularly to the frame. The goal of the experiment was to ride the unicycle or, in other words, to train it to learn how to from falling. In this case, the input vector is 2 dimensional: a first torque applied directly on the wheel mimics a human rider using pedals and a second torque applied laterally to the flywheel. The last element needed to specify the setting is the selected controller. In this case, the authors opted for a linear controller  $\pi(\mathbf{x}, \theta) = \mathbf{A}\mathbf{x} + \mathbf{b}$ ; hence in this example  $\theta = \{\mathbf{A}, \mathbf{b}\}$ .

**Cart-Pole Swing-up** Again from the original paper, PILCO was also applied to learning to control

a real cartpole system. The system consisted of a cart with mass 0.7 kg running on a track and a freely swinging pendulum with mass 0.325 kg attached to the cart. In this case, the state of the system is the position of the cart, the velocity of the cart, the angle of the pendulum, and the angular velocity. A horizontal force  $\mathbf{u}$  could be applied to the cart. The objective was to learn a controller capable to swing the pendulum up and to balance it in the inverted position in the middle of the track. As the authors claim, a linear controller cannot be used in this setting. Rather, the learned state feedback controller was a nonlinear Radial Basis Function (RBF) network:

$$\pi(\mathbf{x}, \theta) = \sum_{i=1}^n w_i \phi_i(\mathbf{x})$$

$$\phi_i(\mathbf{x}) = \exp \left( -\frac{1}{2} (\mathbf{x} - \mu_i)^\top \mathbf{\Lambda}^{-1} (\mathbf{x} - \mu_i) \right)$$

**Cart-double pendulum swing-up** Similarly, the authors applied the PILCO learning strategy to the physics of the Cart-double Pendulum Swing-up. The cart-double pendulum system consists of a cart running on a track and a freely swinging two-link pendulum attached to it. The state of the system is the position, the velocity of the cart, the angles and the angular velocities of both attached pendulums. In this case,  $\mathbf{u}$  corresponds to horizontal forces applied to the cart. Similarly to the Cart-pole Swing-up problem, a linear controller is not capable of solving this problem. A standard control approach to solving the cart-double pendulum task is to design two separate controllers, one for the swing up and one linear controller for the balancing task.

### 3.2 Modeling the underlying dynamics using a GP

As said, PILCO's underlying dynamics are implemented using a GP. In particular:

- the touples

$$\tilde{\mathbf{x}}_{t-1} := [\mathbf{x}_{t-1}, \mathbf{u}_{t-1}] \in \mathbb{R}^{D+F}$$

are the training inputs. For future reference and for a better understanding of the following sections, recall that

$$\mathbf{x}_t = f(\mathbf{x}_{t-1}, \mathbf{u}_{t-1}). \quad (1)$$

- the targets or output are the differences  $\Delta_t = \mathbf{x}_t - \mathbf{x}_{t-1}$ .

Since we are using a GP, we have:

$$\Delta_t = \mathbf{x}_t - \mathbf{x}_{t-1} + \epsilon \in \mathbb{R}^D \quad (2)$$

with  $\epsilon \sim N(0, \Sigma_\epsilon)$ . The authors also assume  $\Sigma_\epsilon = \text{diag}([\sigma_{\epsilon_1}, \dots, \sigma_{\epsilon_2}])$ .

Going through the same derivations that we saw in class, one obtains that:

$$m_f(\tilde{\mathbf{x}}_*) = \mathbb{E}_f[\Delta_*] = \mathbf{k}_*^\top (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y} = \mathbf{k}_*^\top \beta, \quad (3)$$

$$\sigma_f^2(\Delta_*) = \text{var}_f[\Delta_*] = k_{**} - \mathbf{k}_*^\top (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{k}_*, \quad (4)$$

respectively, where  $\mathbf{k}_* := k(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_*)$ ,  $k_{**} := k(\tilde{\mathbf{x}}_*, \tilde{\mathbf{x}}_*)$ ,  $\beta := (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y}$ , and  $\mathbf{K}$  being the Gram matrix with entries  $K_{ij} = k(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)$ .

**Mean function and Covariance function** The prior mean function is chosen to be  $m \equiv 0$ . We have seen in class that the selection of the covariance function represents a crucial point when

modeling a process using a GP. In this case, the authors use the squared exponential kernel (SE), defined as:

$$k(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = \alpha^2 \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}} - \tilde{\mathbf{x}}')' \Lambda^{-1}(\tilde{\mathbf{x}} - \tilde{\mathbf{x}}')\right),$$

where we can think of  $\alpha^2$  as the variance of the latent function  $f$ . We have  $\Lambda := \text{diag}([l_1^2, \dots, l_D^2])$ , so that the kernel function is parameterised by the hyperparameters  $l_i, \alpha^2, \Sigma_\epsilon$ . With the SE function,  $l_i$  have the role of a length-scale. Those hyperparameters are learned by *evidence maximization*, which, using a more statistical vocabulary, corresponds to a maximum likelihood strategy. As we commented in class, the GP with this covariance function is very smooth. This is due to the fact that this covariance function is infinitely differentiable. More details about it can be found in [4]. While such strong smoothness assumptions are unrealistic for modelling many physical processes, the squared exponential is probably the most widely-used kernel within the kernel machines field and turns out to not be harmful for models such PILCO.

### 3.3 From the GP to the expected return

Recall that our goal is to find the optimal  $\pi$ , which minimizes the *expected return* that we defined in previous sections. One can easily see how, in order to evaluate it, we need to first obtain an expression for  $p(\mathbf{x}_t)$ . This section is dedicated to explaining how  $p(\mathbf{x}_t)$  is obtained in the model, starting from the GP structure we introduced in the previous section.

First note that, using (1) and (2), we can say that using a GP to model  $\Delta_t$  implies:

$$p(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{u}_{t-1}) = \mathcal{N}(\mathbf{x}_t \mid \mu_t, \Sigma_t) \quad (5)$$

$$\mu_t = \mathbf{x}_{t-1} + \mathbb{E}_f[\Delta_t], \quad (6)$$

$$\Sigma_t = \text{var}_f[\Delta_t] \quad (7)$$

Now it is time for a further (trivial) observation: in order to obtain  $p(\mathbf{x}_t)$ , we need to get rid of  $\mathbf{u}_t$  and  $f$  in expression (5), (6), (7). This goal is achieved by constructing an approximation for both the joint distribution  $p(\tilde{\mathbf{x}}_{t-1}) = p(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$  and the predictive  $p(\Delta_t)$ . Both are approximated using Gaussian distributions with carefully selected means and covariances. For further reference and notation, we can note that those approximations allow us to write

$$p(\mathbf{u}_t) \approx N(\mu_u, \Sigma_u), \quad p(\Delta_t) \approx N(\mu_\Delta, \Sigma_\Delta)$$

The details are not trivial, and a more comprehensive explanation is given in the appendix. At this point, using properties of Gaussians (hence combining the GP model with the approximations for the joint with  $\mathbf{u}_t$  as well as the Gaussian approximation for  $p(\Delta)$ ), one can write

$$p(\mathbf{x}_t) \approx N(\mu_t, \Sigma_t)$$

with

$$\mu_t = \mu_{t-1} + \mu_\Delta \quad (8)$$

$$\Sigma_t = \Sigma_{t-1} + \Sigma_\Delta + \text{cov}[\mathbf{x}_{t-1}, \Delta_t] + \text{cov}[\Delta_t, \mathbf{x}_{t-1}] \quad (9)$$

$$\text{cov}[\mathbf{x}_{t-1}, \Delta_t] = \text{cov}[\mathbf{x}_{t-1}, \mathbf{u}_{t-1}] \Sigma_u^{-1} \text{cov}[\mathbf{u}_{t-1}, \Delta_t] \quad (10)$$

Details about how to derive  $\text{cov}[\mathbf{x}_t, \Delta_t]$  can be found in Deisenroth 2010 [2].

Now that we obtained an approximation for  $p(\mathbf{x}_t)$ , we can evaluate the expected return  $J^\pi$  by computing the expected values

$$\mathbb{E}_{\mathbf{x}_t}[c(\mathbf{x}_t)] = \int c(\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{x}_t$$

In particular, the authors assumed  $c$  to be modeled as a squared exponential subtracted from unity, i.e.

$$c(\mathbf{x}) = 1 - \exp(-\|\mathbf{x} - \mathbf{x}_{target}\|^2/\sigma_c^2) \in [0, 1]$$

where  $\sigma_c^2$  is a terms that controls the width of  $c$ . This choice for  $c$  implies that the expectations can be found analytically.

### 3.4 From the expected return to the optimal policy

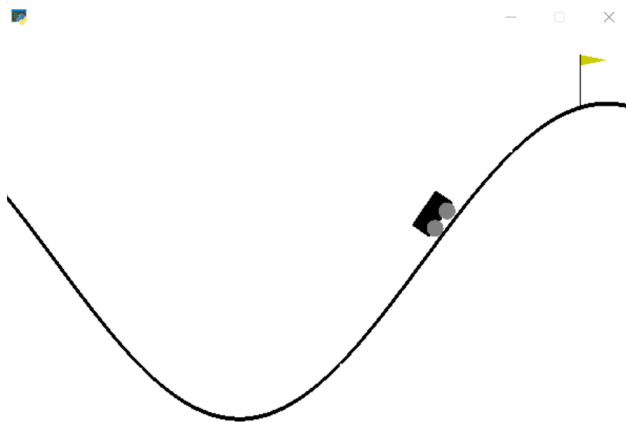
We are now left with the last bite of the PILCO strategy. Recall that examples of policy functions (or controllers) were given earlier in this model section. As said, the controller  $\pi$  can be thought as a function parameterised by  $\theta$ . Hence, finding the optimal policy boils down to finding the optimal  $\theta$ . This is done by analytically computing the gradients of the expected return  $J^\pi$  with respect to the policy parameter  $\theta$ . Note that the parameter space will be different for different applications, since the controllers will have different functional forms. In particular, the derivative  $dJ^\pi/d\theta$  is obtained by repeated application of the chain rule. Details depend on the specific parametrization of the policy  $\pi$ .

## 4 Mountain car implementation

Implementing PILCO from scratch would have been incredibly satisfactory, but also an incredibly unfeasible task for the time available in preparing this project. Still, I was interested in having the chance to look at the mechanics of the model in more operational terms. Hence, I decided to run some experiments based on the code posted online by Nikitas Rontsis and Kyriakos Polymenakos on their GitHub repositories. Tensorflow is used to optimize the parameters in  $\pi$  (that is, in the explanation of the model I just gave, to compute the gradients), while the package GPflow is used to train the GP. The code provides the uses with different PILCO tools: a PILCO basic skeleton which combines the training of the GP and the optimization of the parameters, as well as different specifications of controllers. The user can combine the different pieces and construct a new experiment in the environments offered by OpenAI Gym. Within the OpenAI Gym environments, the most advanced and complex (in terms of both graphics and model structure) are the ones offered in the subgroup Mujoco. Installing those toolkits can be tedious sometimes, but the satisfaction of seeing them in action is worth the effort. Here, a quite simple environment is presented as an example. More elaborate environments can be found in the Mujoco library, presenting 3D rendering and complex objects structures. Unfortunately, due to the last changes in policies and in the packages to download, I had some issues with the more complex Mujoco environments, that I plan to fix soon in order to explore more advanced simulations. More explanation on further directions as well as on software problems encountered are given in the Final Comments section and in the code section in the Appendix.

### The mountain car environment

In this environment, there is a car starting in between two hills. The goal is for the car to reach the top of the hill on the right. The car does not have enough engine power to reach the top of the hill by just accelerating forward. To win, the car must build momentum by swinging itself backwards and forwards until it has enough speed to reach the flag.



The setting is quite simple. In this case,  $\mathbf{x}_t$  corresponds to the car position and velocity so that, in our notation,  $D = 2$ . The car only has 3 possible actions ( $\mathbf{u}$ ) at every state: it can either accelerate forwards, accelerate backwards, or do nothing. Again, the goal is to learn what actions to take at each state to help it reach the flag. Hence, in our notation, the flag will represent  $\mathbf{x}_{target}$ . Following the PILCO strategy I just presented, the car will at every state look at all of its possible actions (accelerate forward, backward, or do not accelerate). It can then calculate the expected value of the future states and select the best action.

As we have seen, a key selection is the choice of the functional form of the policy function  $\pi$ . In this case, similarly to what the authors suggest to use for the Cart-pole swing-up environment, I used a RBF controller (see above in section 3.1). The code is written in Python, and I found its structure very pedagogical. Few comments on it follow:

- Recall that the car needs to learn the strategy on its own. More specifically, we can quote PILCO's inventors: *"Increasing data efficiency requires either having informative prior knowledge or extracting more information from available data. In this paper, we do not assume that any expert knowledge is available (e.g., in terms of demonstrations or differential equations for the dynamics). Instead, we elicit a general policy-search framework for data-efficient learning from scratch"*. In the code, this is achieved by running an initial random rollout to generate a dataset.
- In the actual environment, after the first initial rollout, the variables (hence also the environment reference measures) are normalised for computational stability.
- The packages GP flow and Tensorflow play a key role in training the GP process and optimizing the parameters faster and more efficiently from a computational point of view.
- Notably, the code is modular. This makes it possible to try out the results when choosing different policy functions.
- Importantly, and more on this will be said in the "Further comments on code" section in the Appendix, some recent modifications to OpenAI Gym and Tensorflow create some further

incompatibles when running this code. These problems are in addition to the already tricky usage of Mujoco on Windows machines.

## 5 Further applications and role in the literature

The examples given in the paper were presented earlier in section 4.1. Importantly, given its flexible structure, PILCO can be employed in both more standard benchmark problems and also high-dimensional control problems. Many other applications can be found in the literature. Indeed, after the first seminal paper about PILCO was published, PILCO rapidly became one of the most used strategies for model-based reinforcement learning. Some videos of those implementations can be found at the corresponding YouTube channel: <https://www.youtube.com/user/PilcoLearner>.

An example of usage of PILCO in the literature can be given by the work of Rosenthal et al. [3], who made use of the PILCO framework for designing an algorithm for *safe* policy search. Their work was motivated by the fact that in many real application domains a primary concern is that of safety, particularly the avoidance of specific states or sets of states which are considered dangerous for the system or more generally undesirable (for example, the avoidance of obstacles). Using the probabilistic setting provided by the GP, one can define these constraints *a priori* and then estimate the risk of violating the constraints before any candidate policy is implemented in the actual system. This example, which is one of the many new models that have been built on the PILCO framework, shows once again the potential that statistics, which brings tools to quantify uncertainty, can bring to the field.

## 6 Final comments

The main purpose of this project was to take a first step in understanding one possible way in which statistical tools can be useful for Reinforcing Learning strategies. Also, the goal was to look at the role played by some of the tools we had the chance to look closer at in class. In this case, this tool was the Gaussian Process used to model the underlying evolution of the system.

My feeling is that this project exposed me to learn many new things, some of which are, I believe, really important if I aim to conduct research in the field of RL. In the first place, I had the chance to look closer at the OpenAI Gym, with its really cool environments but also its pitfalls in terms of installing and permit issues. I plan to be able to fix some limitations on my computer and to be able to explore more Mujoco environments in future projects. I also learned that when dealing with Reinforcement Learning research, a good knowledge of tools such as TensorFlow is a prerequisite for making the computation feasible. While I cannot claim to have coded PILCO from scratch, I did some experiments based on the code provided by Nikitas Rontsis and Kyriakos Polymenakos on their websites.

Importantly, this project introduced me to the great potential that statistics can have as a tool for robotics and machine learning techniques as a general category: model uncertainty. As seen with PILCO, this can reduce model bias and provide more efficient procedures. I believe that at least two directions can be taken making use of what learned from this project. On one side, one can think about further applications of the PILCO strategy and about how to export its building blocks to other settings. Moreover, one can ask himself the following question: can we make improvements



to the PILCO setting from a statistical point of view? Can we make the characteristics of the GP (its covariance function, for example) more specific for different environments? On the other side, which is maybe the more open and intellectually challenging one, we can think about other ways in which probabilistic and statistical tools can operate in bringing uncertainty evaluation in the RL field. For example, is it possible to evaluate the uncertainty of well established model-free tools, such as Q-learning? How random are the results obtained with heuristic well know techniques? Can agents be more prone to adapt to new environments by using probabilistic tools to improve their flexibility?

## References

- [1] Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472. Citeseer, 2011.
- [2] Marc Peter Deisenroth. *Efficient reinforcement learning using Gaussian processes*, volume 9. KIT Scientific Publishing, 2010.
- [3] Kyriakos Polymenakos, Alessandro Abate, and Stephen Roberts. Safe policy search with gaussian process models. *arXiv preprint arXiv:1712.05556*, 2017.
- [4] Michael L Stein. *Interpolation of spatial data: some theory for kriging*. Springer Science & Business Media, 1999.
- [5] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

## Appendix

### Approximating the input joint distribution

Note that, at the stage given by equations (3), (4), (5), a joint distribution  $p(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$  is required for predicting  $\mathbf{x}_t$  from  $\mathbf{x}_{t-1}$ . Also, note that those one-step-predictions are crucial for evaluating the expected return  $J$ , since the state distributions  $p(\mathbf{x}_1, \dots, \mathbf{x}_T)$  are obtained by cascading one step predictions. Finally, recall  $\mathbf{u}_{t-1} = \pi(\mathbf{x}_{t-1}, \theta)$ ; hence, the distribution associated computations will depend on the parametrization of the policy  $\pi$ . The authors compute the desired joint distribution as follows. First, they compute the mean  $\mu_u$  and the covariance  $\Sigma_u$  of the predictive control distribution  $p(\mathbf{u}_{t-1})$  by integrating out the state. Subsequently, the cross-covariance  $\text{cov}[\mathbf{x}_{t-1}, \mathbf{u}_{t-1}]$  is computed. Finally, they approximate the joint state-control distribution  $p(\tilde{\mathbf{x}}_{t-1}) = p(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$  with a Gaussian with the correct mean and covariance.

### Approximating the output predictive distribution, $p(\Lambda_t)$

We now simply report the reasoning followed in the paper for completeness of this report.

**MEAN OF THE PREDICTION** Recall that  $D$  is the dimension of  $\mathbf{x}_t$ . Then, following the law of iterated expectations and for target dimensions  $a = 1, \dots, D$ , one has

$$\mu_\Delta^a = \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} [\mathbb{E}_f [f(\tilde{\mathbf{x}}_{t-1}) \mid \tilde{\mathbf{x}}_{t-1}]] = \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} [m_f(\tilde{\mathbf{x}}_{t-1})] \quad (11)$$

$$= \int m_f(\tilde{\mathbf{x}}_{t-1}) \mathcal{N}(\tilde{\mathbf{x}}_{t-1} \mid \tilde{\mu}_{t-1}, \tilde{\Sigma}_{t-1}) d\tilde{\mathbf{x}}_{t-1} \quad (12)$$

$$= \beta_a^\top \mathbf{q}_a \quad (13)$$

with  $\beta_a = (\mathbf{K}_a + \sigma_{\varepsilon_a}^2)^{-1} \mathbf{y}_a$  and  $\mathbf{q}_a = [q_{a1}, \dots, q_{an}]^\top$ . With  $m_f$  given in equation(3), the entries of  $\mathbf{q}_a \in \mathbb{R}^n$  are

$$\begin{aligned} q_{a_i} &= \int k_a(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_{t-1}) \mathcal{N}(\tilde{\mathbf{x}}_{t-1} \mid \tilde{\mu}_{t-1}, \tilde{\Sigma}_{t-1}) d\tilde{\mathbf{x}}_{t-1} \\ &= \frac{\alpha_a^2}{\sqrt{|\tilde{\Sigma}_{t-1} \Lambda_a^{-1} + \mathbf{I}|}} \exp\left(-\frac{1}{2} \nu_i^\top (\tilde{\Sigma}_{t-1} + \Lambda_a)^{-1} \nu_i\right), \\ \nu_i &:= (\tilde{\mathbf{x}}_i - \tilde{\mu}_{t-1}). \end{aligned}$$

Where  $\nu_i$  is the difference between the training input  $\tilde{\mathbf{x}}_i$  and the mean of the "test" input distribution  $p(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$ .

**COVARIANCE MATRIX OF THE PREDICTION** To compute the predictive covariance matrix  $\Sigma_\Delta \in \mathbb{R}^{D \times D}$  the authors distinguish between diagonal elements and off-diagonal elements. Using the law of iterated variances, one obtains for target dimensions  $a, b = 1, \dots, D$

$$\sigma_{aa}^2 = \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} [\text{var}_f [\Delta_a \mid \tilde{\mathbf{x}}_{t-1}]] + \mathbb{E}_{f, \tilde{\mathbf{x}}_{t-1}} [\Delta_a^2] - (\mu_\Delta^a)^2 \quad (14)$$

$$\sigma_{ab}^2 = \mathbb{E}_{f, \tilde{\mathbf{x}}_{t-1}} [\Delta_a \Delta_b] - \mu_\Delta^a \mu_\Delta^b, \quad a \neq b \quad (15)$$

respectively, where  $\mu_\Delta^a$  is known from equation (13). The off-diagonal terms do not contain the additional term  $\mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} [\text{cov}_f [\Delta_a, \Delta_b \mid \tilde{\mathbf{x}}_{t-1}]]$  because of the conditional independence assumption of the GP models: different target dimensions do not covary for given  $\tilde{\mathbf{x}}_{t-1}$ .

First, the terms that are common to both the diagonal and off-diagonal entries are computed.

With the Gaussian approximation  $\mathcal{N}(\tilde{\mathbf{x}}_{t-1} \mid \tilde{\mu}_{t-1}, \tilde{\Sigma}_{t-1})$  of  $p(\tilde{\mathbf{x}}_{t-1})$  and the law of iterated expectations and using (3)

$$\begin{aligned}\mathbb{E}_{f, \tilde{\mathbf{x}}_{t-1}}[\Delta_a \Delta_b] &= \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}}[\mathbb{E}_f[\Delta_a \mid \tilde{\mathbf{x}}_{t-1}] \mathbb{E}_f[\Delta_b \mid \tilde{\mathbf{x}}_{t-1}]] \\ &= \int m_f^a(\tilde{\mathbf{x}}_{t-1}) m_f^b(\tilde{\mathbf{x}}_{t-1}) p(\tilde{\mathbf{x}}_{t-1}) d\tilde{\mathbf{x}}_{t-1}\end{aligned}$$

due to the conditional independence of  $\Delta_a$  and  $\Delta_b$  given  $\tilde{\mathbf{x}}_{t-1}$ . Using again the definition of the mean function  $m_f$  in (3), one obtains

$$\begin{aligned}\mathbb{E}_{f, \tilde{\mathbf{x}}_{t-1}}[\Delta_a \Delta_b] &= \beta_a^\top \mathbf{Q} \beta_b \\ \mathbf{Q} &:= \int k_a(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_{t-1}) k_b(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_{t-1})^\top p(\tilde{\mathbf{x}}_{t-1}) d\tilde{\mathbf{x}}_{t-1}\end{aligned}$$

Now, one can make use of standard results from Gaussian multiplications and integration, and obtain the entries  $Q_{ij}$  of  $\mathbf{Q} \in \mathbb{R}^{n \times n}$   $Q_{ij} = \frac{k_a(\tilde{\mathbf{x}}_i, \tilde{\mu}_{t-1}) k_b(\tilde{\mathbf{x}}_j, \tilde{\mu}_{t-1})}{\sqrt{|\mathbf{R}|}} \exp\left(\frac{1}{2} \mathbf{z}_{ij}^\top \mathbf{R}^{-1} \tilde{\Sigma}_{t-1} \mathbf{z}_{ij}\right)$ , where we defined  $\mathbf{R} := \tilde{\Sigma}_{t-1} (\Lambda_a^{-1} + \Lambda_b^{-1}) + \mathbf{I}$  and  $\mathbf{z}_{ij} := \Lambda_a^{-1} \nu_i + \Lambda_b^{-1} \nu_j$ . Hence, the off-diagonal entries of  $\Sigma_\Delta$  are fully determined by the other results we just obtained.

Finally, from equation (14), we see that the diagonal entries of  $\Sigma_\Delta$  contain an additional term

$$\mathbb{E}_{\tilde{\mathbf{x}}_{t-1}}[\text{var}_f[\Delta_a \mid \tilde{\mathbf{x}}_{t-1}]] = \alpha_a^2 - \text{tr}\left((\mathbf{K}_a + \sigma_{\varepsilon_a}^2 \mathbf{I})^{-1} \mathbf{Q}\right)$$

## Further comments on code

In order to run the code for the examples, we need to execute the following operations. First, OpenAI Gym 0.15.3 and mujoco-py 2.0.2.7. have to be installed manually. Keep in mind that the creators of Mujoco have only recently changed the usage policy. These changes make those environments more accessible, but the installation trickier in a sense. Then, the PILCO package can be installed by running:

```
git clone https://github.com/nrntsis/PILCO && cd PILCO
python setup.py develop
```

The attached code can be used to run some experiments in the MountainCar OpenAI Gym environment.

Unfortunately, recent changes in the Tensorflow package and the the OpenAI Gym environments create some incompatibilities that create serious issues when running code for PILCO. Fixing those problems would require quite some time. Few "wins" on my side were the following:

- I managed to download Mujoco and OpenAI Gym on Windows when the installation for Windows was not available (which was until two months ago, when I started reading the material for this project). New changes in the following weeks made this effort useless.
- I recently managed to have Mujoco running using the administrator account on my computer.

I believe that new adjustments to the packages will be made in the next months that will make fixing the issues more feasible. On my side, if I will have the chance to do research on related topics, I plan to delve more into the mechanics of those tools. As said previously, this is a pitfall of my project. On the other side, I guess this is a lesson in terms of how research doesn't always go as smooth as we would hope it to go.

```

1  #import relevant packages
2  import numpy as np
3  import gym
4  from gpflow import set_trainable
5  np.random.seed(0)
6
7
8  #import models from reference code
9  from pilco.models import PILCO #import main structure
10 from pilco.controllers import RbfController #import selected controller
11 from pilco.rewards import ExponentialReward #import type of reward used for the
    ↪ cost function
12 from utils import rollout, Normalised_Env
13
14 SUBS = 5
15 T = 25
16 env = gym.make('MountainCarContinuous-v0')
17
18 # Initial random rollouts to generate a dataset
19 X1,Y1, _, _ = rollout(env=env, pilco=None, random=True, timesteps=T, SUBS=SUBS,
    ↪ render=True)
20 for i in range(1,5):
21     X1_, Y1_,_,_ = rollout(env=env, pilco=None, random=True, timesteps=T,
    ↪ SUBS=SUBS, render=True)
22     X1 = np.vstack((X1, X1_))
23     Y1 = np.vstack((Y1, Y1_))
24 env.close()
25
26 # A new (normalised) environment is created
27 env = Normalised_Env('MountainCarContinuous-v0', np.mean(X1[:, :2], 0),
    ↪ np.std(X1[:, :2], 0))
28
29 # results of previous rollouts are normalised and used as initial dataset
30 X = np.zeros(X1.shape)
31 X[:, :2] = np.divide(X1[:, :2] - np.mean(X1[:, :2], 0), np.std(X1[:, :2], 0))
32 X[:, 2] = X1[:, -1] # control inputs are not normalised
33 Y = np.divide(Y1 , np.std(X1[:, :2], 0))
34
35 #various setup
36 #dimension set up
37 state_dim = Y.shape[1]
38 control_dim = X.shape[1] - state_dim

```

```

39  #prior setup for GP
40  m_init = np.transpose(X[0,:-1,None])
41  S_init = 0.5 * np.eye(state_dim)
42  #controller setup
43  controller = RbfController(state_dim=state_dim, control_dim=control_dim,
    ↪ num_basis_functions=25)
44  #cost/reward type setup
45  R = ExponentialReward(state_dim=state_dim,
46                        t=np.divide([0.5,0.0] - env.m, env.std),
47                        W=np.diag([0.5,0.1])
48                        )
49
50  pilco = PILCO((X, Y), controller=controller, horizon=T, reward=R, m_init=m_init,
    ↪ S_init=S_init)
51
52  best_r = 0
53  all_Rs = np.zeros((X.shape[0], 1))
54  for i in range(len(all_Rs)):
55      all_Rs[i,0] = R.compute_reward(X[i,None,:-1], 0.001 * np.eye(state_dim))[0]
56
57  ep_rewards = np.zeros((len(X)//T,1))
58
59  for i in range(len(ep_rewards)):
60      ep_rewards[i] = sum(all_Rs[i * T: i*T + T])
61
62  for model in pilco.mgpr.models:
63      model.likelihood.variance.assign(0.05)
64      set_trainable(model.likelihood.variance, False)
65
66  r_new = np.zeros((T, 1))
67  for rollouts in range(5):
68      pilco.optimize_models()
69      pilco.optimize_policy(maxiter=100, restarts=3)
70      import pdb; pdb.set_trace()
71      X_new, Y_new,_,_ = rollout(env=env, pilco=pilco, timesteps=T, SUBS=SUBS,
    ↪ render=True)
72
73  for i in range(len(X_new)):
74      r_new[:, 0] = R.compute_reward(X_new[i,None,:-1], 0.001 *
    ↪ np.eye(state_dim))[0]
75
76  total_r = sum(r_new)
    _, _, r = pilco.predict(m_init, S_init, T)

```

```

77
78     print("Total ", total_r, " Predicted: ", r)
79     X = np.vstack((X, X_new)); Y = np.vstack((Y, Y_new));
80     all_Rs = np.vstack((all_Rs, r_new)); ep_rewards = np.vstack((ep_rewards,
    ↪     np.reshape(total_r,(1,1))))
81     pilco.mgpr.set_data((X, Y))

```