

PILCO Code Documentation v0.9

Marc Peter Deisenroth
Imperial College London, UK
Technische Universität Darmstadt, Germany
m.deisenroth@imperial.ac.uk

Andrew McHutchon
University of Cambridge, UK
ajm257@cam.ac.uk

Joe Hall
University of Cambridge, UK
jah215@cam.ac.uk

Carl Edward Rasmussen
University of Cambridge, UK
cer54@cam.ac.uk

March 1, 2018

Contents

1	Introduction	1
1.1	Intended Use	1
1.2	Software Design and Implementation	2
1.2.1	Model Learning	2
1.2.2	Policy Learning	2
1.2.3	Policy Application	3
1.3	User Interface by Example	4
1.4	Quick Start	4
2	Software Package Overview	5
2.1	Main Modules	5
2.1.1	applyController	6
2.1.2	trainDynModel	7
2.1.3	learnPolicy	7
2.2	Working with a Real Robot	8
3	Important Function Interfaces	10
3.1	GP Predictions	10
3.1.1	Input Arguments	11
3.1.2	Output Arguments	11
3.2	Controller	12
3.2.1	Interface	12
3.3	Cost Functions	13
3.3.1	Interface for Scenario-specific Cost Functions	13
3.3.2	General Building Blocks	14
4	How to Create Your Own Scenario	19
4.1	Necessary Files	19
4.2	ODE Dynamics	19
4.3	Scenario-specific Settings	20
4.3.1	Adding Paths	20
4.3.2	Indices	21
4.3.3	General Settings	22
4.3.4	Plant Structure	23
4.3.5	Policy Structure	24
4.3.6	Cost Function Structure	25
4.3.7	GP Dynamics Model Structure	27

4.3.8	Optimization Parameters (Policy Learning)	27
4.3.9	Plotting Parameters	28
4.3.10	Allocating Variables	28
4.4	Cost Function	28
4.4.1	Interface	28
4.5	Visualization	31
4.6	Main Function	32
4.6.1	Screen Prints and Visualization	33
5	Implemented Scenarios	38
5.1	Pendulum Swing-up	38
5.2	Double Pendulum Swing-up with a Single Actuator (Pendubot)	39
5.3	Double Pendulum Swing-up with Two Actuators	41
5.4	Cart-Pole Swing-up	41
5.5	Cart-Double Pendulum Swing-up	42
5.6	Unicycling	44
5.6.1	Method	45
5.6.2	Wheel FBD	46
5.6.3	Frame FBD	46
5.6.4	Turntable FBD	48
5.6.5	Eliminating Internal Forces	48
6	Testing and Debugging	50
6.1	Gradient Checks for the Controller Function	50
6.1.1	Interface	50
6.2	Gradient Checks for the Cost Function	51
6.2.1	Interface	51
6.3	Gradient Checks for the GP Prediction Function	52
6.3.1	Interface	52
6.4	Gradient Checks for the State Propagation Function	53
6.4.1	Interface	53
6.5	Gradient Checks for Policy Evaluation	53
6.5.1	Interface	53
7	Code and Auto-generated Documentation of the Main Functions	55
7.1	Base Directory	55
7.1.1	applyController.m	55
7.1.2	propagate.m	56
7.1.3	rollout.m	58
7.1.4	trainDynModel.m	60
7.1.5	value.m	61
7.2	Control Directory	62
7.2.1	concat.m	62
7.2.2	congp.m	64
7.2.3	conlin.m	66
7.3	GP Directory	68
7.3.1	train.m	68
7.3.2	hypCurb.m	69

7.3.3	fitc.m	71
7.3.4	gp0.m	73
7.3.5	gp1.m	75
7.3.6	gp2.m	78

Abstract

We describe a Matlab package for the PILCO policy search framework for data-efficient reinforcement learning. This package implements the PILCO learning framework in multiple different scenarios with continuous states and actions: pendulum swing-up, cart-pole swing-up, double-pendulum swing-up (with either one or two actuators), cart-double-pendulum swing-up, and unicycling. Results from some of these scenarios have been presented previously in [3, 4]. The high-level steps of the PILCO algorithm, which are also implemented in this software package, are the following: Learn a Gaussian process (GP) model of the system dynamics, perform deterministic approximate inference for policy evaluation, update the policy parameters using exact gradient information, apply the learned controller to the system. The software package provides an interface that allows for setting up novel tasks without the need to be familiar with the intricate details of model learning, policy evaluation and improvement.

Chapter 1

Introduction

Reinforcement learning (RL) is a general paradigm for learning (optimal) policies for stochastic sequential decision making processes [10]. In many practical engineering applications, such as robotics and control, RL methods are difficult to apply: First, the state and action spaces are often continuous valued and high dimensional. Second, the number of interactions that can be performed with a real system is practically limited. Therefore, learning methods that efficiently extract valuable information from available data are important. Policy search methods have been playing an increasingly important role in robotics as they consider a simplified RL problem and search (optimal) policies in a constrained policy space [1, 3], typically in an episodic set-up, i.e., a finite-horizon set-up with a fixed initial state (distribution).

We present the PILCO software package for data-efficient policy search that allows to learn (non)linear controllers with hundreds of parameters for high-dimensional systems. A key element is a learned probabilistic GP dynamics model. Uncertainty about the learned dynamics model (expressed by the GP posterior) is explicitly taken into account for multiple-step ahead predictions, policy evaluation, and policy improvement.

1.1 Intended Use

The intended use of this software package is to provide a relatively simple interface for practitioners who want to solve RL problems with continuous states and actions efficiently, i.e., without the need of excessively sized data sets. As PILCO is very data efficient, it has been successfully applied to robots to learn policies from scratch, i.e., without the need to provide demonstrations or other “informative” prior knowledge [3, 4]. In this document, we intentionally hide the involved model learning and inference mechanisms to make the code more accessible. Details about inference and model learning can be found in [2].

This software package is *not* ideal for getting familiar with classical RL scenarios and algorithms (e.g., Q-learning, SARSA, TD-learning), which typically involve discrete states and actions. For this purpose, we refer to existing RL software packages, such as RLGlue, CLSquare¹, PIQLE², RL Toolbox³, LibPG⁴, or RLPy⁵.

¹<http://www.ni.uos.de/index.php?id=70>

²<http://piqle.sourceforge.net/>

³<http://www.igi.tugraz.at/ril-toolbox/>

⁴<http://code.google.com/p/libpg/>

⁵<http://acl.mit.edu/RLPy/>

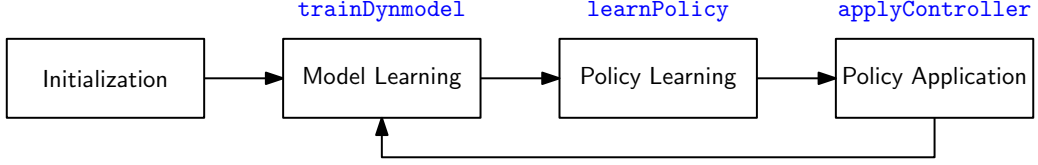


Figure 1.1: Main modules: After an initialization, the first module is responsible for training a GP from available data. The second module is used for policy learning, which consists of policy evaluation and improvement. The third module is responsible for applying the learned policy to the system, which can be either a simulated system or a real system, such as a robot. The collected data from this application is used for updating the model, and the cycle starts from the beginning.

1.2 Software Design and Implementation

The high-level modules 1) Model learning, 2) Policy learning, 3) Policy application are summarized in the following.

1.2.1 Model Learning

The forward model learned is a non-parametric, probabilistic Gaussian process [7]. The non-parametric property of the GP does not require an explicit task-dependent parametrization of the dynamics of the system. The probabilistic property of the GP reduces the effect of model errors.

Inputs to the GP are state-action pairs $(\mathbf{x}_t, \mathbf{u}_t)$, where t is a time index. Training targets are either successor states \mathbf{x}_{t+1} or differences $\Delta_t = \mathbf{x}_{t+1} - \mathbf{x}_t$.

By default, a full GP model is trained by evidence maximization, where we penalize high signal-to-noise ratios in order to maintain numerical stability. There is an option to switch to sparse GPs in case the number of data points exceeds a particular threshold. We implemented the FITC/SPGP sparse GP method proposed by [8] for GP training and predictions at uncertain inputs.

1.2.2 Policy Learning

For policy learning, PILCO uses the learned GP forward model to compute approximate long-term predictions $p(\mathbf{x}_1|\pi), \dots, p(\mathbf{x}_T|\pi)$ for a given controller π . To do so, we follow the analytic moment-matching approach proposed by [6], and approximate all $p(\mathbf{x}_t|\pi)$ by Gaussians $\mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$.

Policy Evaluation. Once the long-term predictions $p(\mathbf{x}_1|\pi), \dots, p(\mathbf{x}_T|\pi)$ are computed, the expected long-term cost

$$J^\pi = \sum_{t=1}^T \mathbb{E}[c(\mathbf{x}_t)|\pi], \quad p(\mathbf{x}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0), \quad (1.1)$$

can be computed analytically for many cost functions c , e.g., polynomials, trigonometric functions, or Gaussian-shaped functions. In our implementation, we typically use a Gaussian-shaped cost function

$$c(\mathbf{x}) = 1 - \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}_{\text{target}}\|_{\mathbf{W}}^2 / \sigma_c^2\right),$$

where $\|\mathbf{x} - \mathbf{x}_{\text{target}}\|_{\mathbf{W}}^2$ is the Mahalanobis distance between \mathbf{x} and $\mathbf{x}_{\text{target}}$, weighted by \mathbf{W} , σ_c is a scaling factor, and $\mathbf{x}_{\text{target}}$ is a target state.

Policy Improvement. To improve the policy, we use gradient-based Quasi-Newton optimization methods, such as BFGS. The required gradients of the expected long-term cost J^π with respect to the policy parameters are computed analytically, which allows for learning low-level policies with hundreds of parameters [3].

1.2.3 Policy Application

This module takes care of applying the learned controller to the (simulated) system. At each time step t , the learned controller π computes the corresponding control signal \mathbf{u}_t from the current state \mathbf{x}_t . An ODE solver is used to determine the corresponding successor state \mathbf{x}_{t+1} . The module returns a trajectory of state-action pairs, from which the training inputs and targets for the GP model can be extracted.

If the controller is applied to a real system, such as a robot, this module is not necessary. Instead, state estimation and control computation need to be performed on the robot directly [4].

1.3 User Interface by Example

```
1 % 0. Initialization
2 settings; % load scenario-specific settings
3
4 for jj = 1:J % Initial J random rollouts
5     [xx, yy, realCost{jj}, latent{jj}] = ...
6         rollout(gaussian(mu0, S0), struct('maxU', policy.maxU), H, plant, cost);
7     x = [x; xx]; y = [y; yy]; % augment training sets for dynamics model
8 end
9
10 % Controlled learning (N iterations)
11 for j = 1:N
12     % 1. Train (GP) dynamics model
13     trainDynModel;
14
15     % 2. Learn policy
16     learnPolicy;
17
18     % 3. Apply controller to system
19     applyController;
20 end
```

In line 2, a scenario-specific settings script is executed. Here, we have to define the policy structure (e.g., parametrization, torque limits), the cost function, the dynamics model (e.g., definition of training inputs/targets), some details about the system/plant (e.g., sampling frequency) that are needed for the ODE solver, and some general parameters (e.g., prediction horizon, discount factor).

In lines 4–8, we create an initial set of trajectory rollouts by applying random actions to the system, starting from a state \mathbf{x}_0 sampled from $p(\mathbf{x}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \mathbf{S}_0)$. The `rollout` function in line 6 takes care of this. For each trajectory, the training inputs and targets are collected in the matrices \mathbf{x} and \mathbf{y} , respectively.

In lines 11–20, the three main modules are executed iteratively. First, the GP dynamics model is learned in the `trainDynModel` script, using the current training data \mathbf{x}, \mathbf{y} (line 13). Second, the policy is learned in the `learnPolicy` script, which updates the policy parameters using analytic policy evaluation and policy gradients. The third module, i.e., the application of the learned controller to the system, is encapsulated in the `applyController` script, which also augments the current data set \mathbf{x}, \mathbf{y} for training the GP model with the new trajectory.

1.4 Quick Start

If you want to try it out without diving into the details, navigate to `<pilco_root>/scenarios/cartPole` and execute the `cartPole_learn` script.

Chapter 2

Software Package Overview

This software package implements the PILCO reinforcement learning framework [3]. The package contains the following directories

- **base**: Root directory. Contains all other directories.
- **control**: Directory that implements several controllers.
- **doc**: Documentation
- **gp**: Everything that has to do with Gaussian processes (training, predictions, sparse GPs etc.)
- **loss**: Several immediate cost functions
- **scenarios**: Different scenarios. Each scenario is packaged in a separate directory with all scenario-specific files
- **util**: Utility files
- **test**: Test functions (derivatives etc.)

2.1 Main Modules

The main modules of the PILCO framework and their interplay are visualized in Figure 1.1. Each module is implemented in a separate script and can be found in `<pilco_root>/base`.

Let us have a look at the high-level functionality of the three main modules in Figure 1.1:

1. `applyController`
 - (a) determine start state
 - (b) generate rollout
 - i. compute control signal $\pi(\mathbf{x}_t)$
 - ii. simulate dynamics (or apply control to real robot)
 - iii. transition to state \mathbf{x}_{t+1}
2. `trainDynModel`
3. `learnPolicy`

- (a) call gradient-based non-convex optimizer **minimize**: minimize value with respect to policy parameters θ
 - i. **propagated**: compute successor state distribution $p(\mathbf{x}_{t+1})$ and gradients $\partial p(\mathbf{x}_{t+1})/\partial \theta$ with respect to the policy parameters and gradients $\partial p(\mathbf{x}_{t+1})/\partial p(\mathbf{x}_t)$ with respect to the previous state distribution $p(\mathbf{x}_t)$.
 - A. trigonometric augmentation of the state distribution $p(\mathbf{x}_t)$
 - B. compute distribution of preliminary (unsquashed) policy $p(\tilde{\pi}(\mathbf{x}_t))$
 - C. compute distribution of squashed (limited-amplitude) policy $p(\pi(\mathbf{x}_t)) = p(\mathbf{u}_t)$
 - D. determine successor state distribution $p(\mathbf{x}_{t+1})$ using GP prediction (**gp***)
 - ii. **cost.fcn**: Scenario-specific function that computes the expected (immediate) cost $\mathbb{E}_{\mathbf{x}}[c(\mathbf{x})]$ and its partial derivatives $\partial \mathbb{E}_{\mathbf{x}}[c(\mathbf{x})]/\partial p(\mathbf{x})$

2.1.1 applyController

```

1 % 1. Generate trajectory rollout given the current policy
2 if isfield(plant, 'constraint'), HH = maxH; else HH = H; end
3 [xx, yy, realCost{j+J}, latent{j}] = ...
4   rollout(gaussian(mu0, S0), policy, HH, plant, cost);
5 disp(xx); % display states of observed trajectory
6 x = [x; xx]; y = [y; yy]; % augment training set
7 if plotting.verbosity > 0
8   if ~ishandle(3); figure(3); else set(0, 'CurrentFigure', 3); end
9   hold on; plot(1:length(realCost{J+j}), realCost{J+j}, 'r'); drawnow;
10 end
11
12 % 2. Make many rollouts to test the controller quality
13 if plotting.verbosity > 1
14   lat = cell(1,10);
15   for i=1:10
16     [~,~,~,lat{i}] = rollout(gaussian(mu0, S0), policy, HH, plant, cost);
17   end
18
19   if ~ishandle(4); figure(4); else set(0, 'CurrentFigure', 4); end; clf(4);
20
21   ldyno = length(dyno);
22   for i=1:ldyno % plot the rollouts on top of predicted error bars
23
24     subplot(ceil(ldyno/sqrt(ldyno)), ceil(sqrt(ldyno)), i); hold on;
25     errorbar(0:length(M{j})(i,:)-1, M{j}(i,:), ...
26       2*sqrt(squeeze(Sigma{j})(i,i,:))) );
27     for ii=1:10
28       plot(0:size(lat{ii}(:,dyno(i)),1)-1, lat{ii}(:,dyno(i)), 'r');
29     end
30     plot(0:size(latent{j}(:,dyno(i)),1)-1, latent{j}(:,dyno(i)), 'g');
31     axis tight
32   end
33   drawnow;
34 end
35
36 % 3. Save data
37 filename = [basename num2str(j) '_H' num2str(H)]; save(filename);

```

The script `applyController` executes the following high-level steps:

1. Generate a trajectory rollout by applying the current policy to the system (lines 1–4). The initial state is sampled from $p(\mathbf{x}_0) = \mathcal{N}(\mu_0, S_0)$, see line 4. This trajectory is used to augment the GP training set (line 6).
2. (optional) Generate more trajectories with different start states $\mathbf{x}_0 \sim p(\mathbf{x}_0)$ and plot a sample distribution of the trajectory distribution (lines 12– 33).
3. Save the entire workspace (line 36).

2.1.2 trainDynModel

```

1 % 1. Train GP dynamics model
2 Du = length(policy.maxU); Da = length(plant.angi); % no. of ctrl and angles
3 xaug = [x(:,dyno) x(:,end-Du-2*Da+1:end-Du)]; % x augmented with angles
4 dynmodel.inputs = [xaug(:,dyni) x(:,end-Du+1:end)]; % use dyni and ctrl
5 dynmodel.targets = y(:,dyno);
6 dynmodel.targets(:,difi) = dynmodel.targets(:,difi) - x(:,dyno(difi));
7
8 dynmodel = dynmodel.train(dynmodel, plant, trainOpt); % train dynamics GP
9
10 % display some hyperparameters
11 Xh = dynmodel.hyp;
12 % noise standard deviations
13 disp(['Learned noise std: ' num2str(exp(Xh(end,:)))]);
14 % signal-to-noise ratios (values > 500 can cause numerical problems)
15 disp(['SNRs : ' num2str(exp(Xh(end-1,:)-Xh(end,:)))]);

```

The script that takes care of training the GP executes the following high-level steps:

1. Extract states and controls from **x**-matrix (lines 2–3)
2. Define the training inputs and targets of the GP (lines 4–6)
3. Train the GP (line 8)
4. Display GP hyper-parameters, the learned noise hyper-parameters, and the signal-to-noise ratios (lines 10–15). This information is very valuable for debugging purposes.

2.1.3 learnPolicy

```

1 % 1. Update the policy
2 opt.fh = 1;
3 [policy.p fX3] = minimize(policy.p, 'value', opt, mu0Sim, S0Sim, ...
4   dynmodel, policy, plant, cost, H);
5
6 % (optional) Plot overall optimization progress
7 if exist('plotting', 'var') && isfield(plotting, 'verbosity') ...
8   && plotting.verbosity > 1
9   if ~ishandle(2); figure(2); else set(0, 'CurrentFigure', 2); end
10  hold on; plot(fX3); drawnow;
11  xlabel('line search iteration'); ylabel('function value')
12 end
13

```

```

14 % 2. Predict trajectory from p(x0) and compute cost trajectory
15 [M{j} Sigma{j}] = pred(policy, plant, dynmodel, mu0Sim(:,1), SOSim, H);
16 [fantasy.mean{j} fantasy.std{j}] = ...
17   calcCost(cost, M{j}, Sigma{j}); % predict cost trajectory
18
19 % (optional) Plot predicted immediate costs (as a function of the time steps)
20 if exist('plotting', 'var') && isfield(plotting, 'verbosity') ...
21   && plotting.verbosity > 0
22   if ~ishandle(3); figure(3); else set(0, 'CurrentFigure', 3); end
23   clf(3); errorbar(0:H, fantasy.mean{j}, 2*fantasy.std{j}); drawnow;
24   xlabel('time step'); ylabel('immediate cost');
25 end

```

learnPolicy

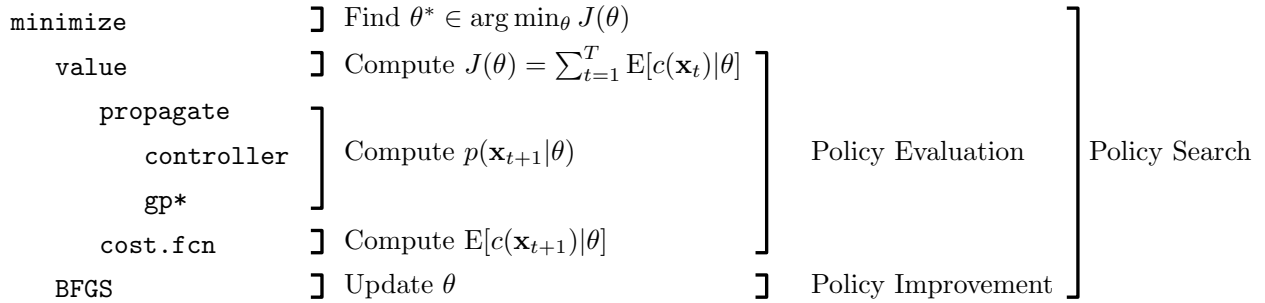


Figure 2.1: Functions being called from `learnPolicy.m` for learning the policy.

1. Learn the policy by calling `minimize`. Figure 2.1 depicts the functions that are called by `learnPolicy` in order to perform the policy search to find a good parameter set θ^* .
2. (optional) Plot overall optimization progress.
3. Long-term prediction of a state trajectory from $p(\mathbf{x}_0)$ using the learned policy (line 15) by calling `pred`. This prediction is equivalent to the last predicted trajectory during policy learning, i.e., the predicted state trajectory that belongs to the learned controller.
4. The predicted state trajectory is used to compute the corresponding distribution over immediate costs (lines 16–17) by calling `calcCost`.
5. (optional) Plot the predicted immediate cost distribution as a function of the time steps (lines 19–25).

2.2 Working with a Real Robot

When you want to apply PILCO to a learning controller parameters for a real robot, only a few modifications are necessary. As policy learning is not real-time anyway, it does not make too much sense performing it on the robot directly. Therefore, the robot only needs to know about the learned policy, but nothing about the learned dynamics model.

Here is a list of modifications:

- An ODE does not need to be specified for simulating the system.

- All trajectory rollouts are executed directly on the robot.
- The module `applyController` needs to take care of generating a trajectory on the robot.
- For generating a trajectory using the robot, the probably least coding extensive approach is the following:
 1. Learn the dynamics model and the policy.
 2. Save the learned policy parameters in a file.
 3. Transfer the parameters to your robot
 4. Write a controller function in whatever programming language the robot needs.
 5. When the controller is applied, just map the measured state through the policy to obtain the desired control signal.
 6. Save the recorded trajectory in a file and make it available to PILCO and save them in `xx`, `yy`.

Here is a high-level code-snippet that explains the main steps.

```

1 % 0. Initialization
2 settings; % load scenario-specific settings
3
4 applyController_on_robot; % collect data from robot
5
6 % Controlled learning (N iterations)
7 for j = 1:N
8   % 1. Train (GP) dynamics model
9   trainDynModel;
10
11   % 2. Learn policy
12   learnPolicy;
13
14   % 3. Apply controller to system
15   applyController_on_robot;
16 end

```

We successfully applied this procedure on different hardware platforms [4, 3].

Chapter 3

Important Function Interfaces

The PILCO software package relies on several high-level functionalities with unified interfaces:

- Predicting with GPs when the test input is Gaussian distributed. We have implemented several versions of GPs (including sparse GPs), which perform these predictions. The generic interface is detailed in the following.
- Controller functions. With a unified interface, it is straightforward to swap between controllers in a learning scenario. We discuss the generic interface in this chapter.
- Cost functions. With this software package, we ship implementations of several cost functions. The interface of them is discussed in this chapter.

3.1 GP Predictions

Table 3.1 gives an overview of all implemented functions that are related to predicting with GPs at a Gaussian distributed test input $\mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$. We assume that the input dimension is D and the

Table 3.1: Overview of functions for GP predictions with Gaussian distributed test inputs.

	[M S V]	[dMdm dSdm dVdm dMds dSds dVds]	[dMdP dSdP dVdP]	sparse	prob. GP
gp0	✓			✗	✓
gp0d	✓			✗	✓
gp1	✓			✓	✓
gp1d	✓			✓	✓
gp2	✓			✗	✗
gp2d	✓			✗	✗

predictive dimension is E . All functions are in the directory `<pilco_root>/gp/`.

The convention in the function name is that a “d” indicates that derivatives are computed. For instance, **gp0d** computes the same function values as **gp0**, but it additionally computes some derivatives. We have three different categories of functions for GP predictions:

- **gp0**: The underlying model is a full probabilistic GP model. This model is used for implementing the standard GP dynamics model.

- **gp1**: The underlying model is a sparse GP model. In particular, we use the SPGP/FITC approach by Snelson and Ghahramani [8]. This model is used for implementing the GP dynamics when the data set is too large.
- **gp2**: The underlying model is a full “deterministic” GP model. The model differs from the full probabilistic model (**gp0**) by ignoring the posterior uncertainty about the underlying function. The model essentially consists of the mean function only. This makes it functionally equivalent to a radial-basis-function (RBF) network. This kind of model is used for implementing nonlinear controllers.

3.1.1 Input Arguments

For all functions **gp***, the input arguments are the following:

1. **gpmodel**: Structure containing all relevant information
 - **.hyp**: log-hyper-parameters in a $(D + 2) \times E$ matrix (D log-length scales, 1 log-signal-standard deviation, 1 log-noise-standard deviation per predictive dimension)
 - **.inputs**: training inputs in an $n \times D$ matrix
 - **.targets**: training targets in an $n \times E$ matrix
2. **m**: Mean of the state distribution $p(\mathbf{x})$, $(D \times 1)$
3. **s**: Covariance matrix of the state distribution $p(\mathbf{x})$, $(D \times D)$

3.1.2 Output Arguments

The **gp*** functions can be used to compute the mean and the covariance of the joint distribution $p(\mathbf{x}, f(\mathbf{x}))$, where $f \sim \mathcal{GP}$ and $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{s})$.

All functions **gp*** predict the mean **M** and the covariance **S** of $p(f(\mathbf{x}))$ as well as $\mathbf{V} = \mathbf{s}^{-1} \text{cov}[\mathbf{x}, f(\mathbf{x})]$. Note that **gp1*** compute these values using sparse GPs and **gp2*** use only the mean function of the GP, i.e., the posterior uncertainty about f is discarded.

For policy learning, we require from the *dynamics model* the following derivatives:

- **dMdm**: $\partial M / \partial m \in \mathbb{R}^{E \times D}$ The derivative of the mean of the prediction with respect to the mean of the input distribution.
- **dSdm**: $\partial S / \partial m \in \mathbb{R}^{E^2 \times D}$ The derivative of the covariance of the prediction with respect to the mean of the input distribution.
- **dVdm**: $\partial V / \partial m \in \mathbb{R}^{DE \times D}$ The derivative of V with respect to the mean of the input distribution.
- **dMds**: $\partial M / \partial s \in \mathbb{R}^{E \times D^2}$ The derivative of the mean of the prediction with respect to the covariance of the input distribution.
- **dSds**: $\partial S / \partial s \in \mathbb{R}^{E^2 \times D^2}$ The derivative of the covariance of the prediction with respect to the covariance of the input distribution.
- **dVds**: $\partial V / \partial s \in \mathbb{R}^{DE \times D^2}$ The derivative of V with respect to the covariance of the input distribution.

As `gp0d` and `gp1d` are the functions used to propagate uncertainties through a GP dynamics model, they all compute these derivatives, see Table 3.1.

When we use `gp2*` as a convenient implementation of an RBF network *controller*, we additionally require the gradients of `M`, `S`, `V` with respect to the “parameters” of the GP, which are abbreviated by `P` in Table 3.1. These parameters comprise the training inputs, the training targets, and the log-hyper-parameters:

- `dMdP` = $\{\partial M/\partial X, \partial M/\partial y, \partial M/\partial \theta\}$: The derivative of the mean prediction with respect to the training inputs X , the training targets y , and the log-hyper-parameters θ .
- `dSdP` = $\{\partial S/\partial X, \partial S/\partial y, \partial S/\partial \theta\}$: The derivative of the covariance of the prediction with respect to the training inputs X , the training targets y , and the log-hyper-parameters θ .
- `dVdP` = $\{\partial V/\partial X, \partial V/\partial y, \partial V/\partial \theta\}$: The derivative of V with respect to the training inputs X , the training targets y , and the log-hyper-parameters θ .

3.2 Controller

The control directory is located at `<pilco_root>/control`. The controllers compute the (unconstrained) control signals $\tilde{\pi}(\mathbf{x})$.

The generic function call is as follows, where `controller` is a generic name¹:

```
1 function [M, S, V, dMdm, dSdm, dCdm, dMds, dSds, dVds, dMdP, dSdP, dVdP] ...
2         = controller(policy, m, s)
```

3.2.1 Interface

Let us explain the interface in more detail

3.2.1.1 Input Arguments

All controllers expect the following inputs

1. `policy`: A struct with the following fields
 - `policy.fcn`: A function handle to `controller`. This is not needed by the controller function itself, but by other functions that call `controller`.
 - `policy.p`: The policy parameters. Everything that is in this field is considered a free parameter and optimized during policy learning.
 - `policy.<>`: Other arguments the controller function requires.
2. `m`: $\mathbb{E}[\mathbf{x}] \in \mathbb{R}^D$ The mean of the state distribution $p(\mathbf{x})$.
3. `s`: $\mathbb{V}[\mathbf{x}] \in \mathbb{R}^{D \times D}$ The covariance matrix of the state distribution $p(\mathbf{x})$.

¹We have implemented two controller functions: `conlin` and `congp`

3.2.1.2 Output Arguments

All controller functions are expected to compute

1. **M**: $\mathbb{E}[\tilde{\pi}(\mathbf{x})] \in \mathbb{R}^F$ The mean of the predicted (unconstrained) control signal
2. **S**: $\mathbb{V}[\tilde{\pi}(\mathbf{x})] \in \mathbb{R}^{F \times F}$ The covariance matrix of the predicted (unconstrained) control signal
3. **V**: $\mathbb{V}(\mathbf{x})^{-1} \text{cov}[\mathbf{x}, \tilde{\pi}(\mathbf{x})]$ The cross-covariance between the (input) state \mathbf{x} and the control signal $\tilde{\pi}(\mathbf{x})$, pre-multiplied with $\mathbb{V}(\mathbf{x})^{-1}$, the inverse of the covariance matrix of $p(\mathbf{x})$. We do not compute $\text{cov}[\mathbf{x}, \tilde{\pi}(\mathbf{x})]$ because of numerical reasons.
4. Gradients. The gradients of all output arguments with respect to all input arguments are computed:
 - **dMdm**: $\partial M / \partial m \in \mathbb{R}^{F \times D}$ The derivative of the mean of the predicted control with respect to the mean of the state distribution.
 - **dSdm**: $\partial S / \partial m \in \mathbb{R}^{F^2 \times D}$ The derivative of the covariance of the predicted control with respect to the mean of the state distribution.
 - **dVdm**: $\partial V / \partial m \in \mathbb{R}^{DF \times D}$ The derivative of V with respect to the mean of the state distribution.
 - **dMds**: $\partial M / \partial s \in \mathbb{R}^{F \times D^2}$ The derivative of the mean of the predicted control with respect to the covariance of the state distribution.
 - **dSds**: $\partial S / \partial s \in \mathbb{R}^{F^2 \times D^2}$ The derivative of the covariance of the predicted control with respect to the covariance of the state distribution.
 - **dVds**: $\partial V / \partial s \in \mathbb{R}^{DF \times D^2}$ The derivative of V with respect to the covariance of the state distribution.
 - **dMdθ**: $\partial M / \partial \theta \in \mathbb{R}^{F \times |\theta|}$ The derivative of the mean of the predicted control with respect to the policy parameters θ .
 - **dSdθ**: $\partial S / \partial \theta \in \mathbb{R}^{F^2 \times |\theta|}$ The derivative of the covariance of the predicted control with respect to the policy parameters θ .
 - **dVdθ**: $\partial V / \partial \theta \in \mathbb{R}^{DF \times |\theta|}$ The derivative of V with respect to the policy parameters θ .

3.3 Cost Functions

Any (generic) cost function is supposed to compute the expected (immediate) cost $\mathbb{E}[c(\mathbf{x})]$ and the corresponding variance $\mathbb{V}[c(\mathbf{x})]$ for a Gaussian distributed state $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{S})$.

Cost functions have to be written for each scenario. Example cost functions can be found in `<pilco_root>/scenarios/*`.

3.3.1 Interface for Scenario-specific Cost Functions

```
1 function [L, dLdm, dLds, S] = loss(cost, m, s)
```

Input Arguments

<code>cost</code>	cost structure	
<code>.p</code>	parameters that are required to compute the cost, e.g., length of pendulum	[P x 1]
<code>.expl</code>	(optional) exploration parameter	
<code>.target</code>	target state	[D x 1]
<code>m</code>	mean of state distribution	[D x 1]
<code>s</code>	covariance matrix for the state distribution	[D x D]

We only listed typical fields of the `cost` structure. It is possible to add more information. `cost.expl` allows for UCB-type exploration, in which case the returned cost L should be computed according to

$$L(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[c(\mathbf{x})] + \kappa \sqrt{\mathbb{V}_{\mathbf{x}}[c(\mathbf{x})]}, \quad (3.1)$$

where κ is an exploration parameter stored in `cost.expl`. Exploration is encouraged for $\kappa < 0$ and discouraged for $\kappa > 0$. By default, exploration is disabled, i.e., `cost.expl=0`. A target state can be passed in via `cost.target`, but could also be hard-coded in the cost function.

Output Arguments

<code>L</code>	expected cost	[1 x 1]
<code>dLdm</code>	derivative of expected cost wrt. state mean vector	[1 x D]
<code>dLds</code>	derivative of expected cost wrt. state covariance matrix	[1 x D^2]
<code>S</code>	variance of cost	[1 x 1]

Note that the expected cost $L = \mathbb{E}[c(\mathbf{x})]$ can take care of UCB-type exploration, see Equation (3.1). The gradients of L with respect to the mean (`dLdm`) and covariance (`dLds`) of the input distribution are required for policy learning.

3.3.2 General Building Blocks

We have implemented some generic building blocks that can be called by the scenario-specific cost functions. In the following, we detail the computation of a saturating cost function `<pilco.root>/loss/lossSat.m` and a quadratic cost function `<pilco.root>/loss/lossQuad.m`.

3.3.2.1 Saturating Cost

`lossSat` computes the expectation and variance of a saturating cost

$$1 - \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{z})^\top \mathbf{W}(\mathbf{x} - \mathbf{z})\right) \in [0, 1]$$

and their derivatives, where $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{S})$, and a is a normalizing constant. The matrix \mathbf{W} is never inverted and plays the role of a precision matrix. Moreover, it is straightforward to eliminate the influence of state variables in the cost function by setting the corresponding values in \mathbf{W} to 0.

```
1 function [L, dLdm, dLds, S, dSdm, dSds, C, dCdm, dCds] = lossSat(cost, m, s)
```

Input arguments:

cost		
.z:	target state	[D x 1]
.W:	weight matrix	[D x D]
m	mean of input distribution	[D x 1]
s	covariance matrix of input distribution	[D x D]

Output arguments:

L	expected loss	[1 x 1]
dLdm	derivative of L wrt input mean	[1 x D]
dLds	derivative of L wrt input covariance	[1 x D^2]
S	variance of loss	[1 x 1]
dSdm	derivative of S wrt input mean	[1 x D]
dSds	derivative of S wrt input covariance	[1 x D^2]
C	inv(S) times input-output covariance	[D x 1]
dCdm	derivative of C wrt input mean	[D x D]
dCds	derivative of C wrt input covariance	[D x D^2]

Implementation

```

2 % some precomputations
3 D = length(m); % get state dimension
4
5 % set some defaults if necessary
6 if isfield(cost, 'W'); W = cost.W; else W = eye(D); end
7 if isfield(cost, 'z'); z = cost.z; else z = zeros(D,1); end
8
9 SW = s*W;
10 iSpW = W/(eye(D)+SW);

```

In lines 6–7, we check whether the weight matrix \mathbf{W} and the state \mathbf{z} exist. Their default values are \mathbf{I} and $\mathbf{0}$, respectively. Lines 9–10 do some pre-computations of matrices that will be frequently used afterwards.

```

11 % 1. Expected cost
12 L = -exp(-(m-z)'*iSpW*(m-z)/2)/sqrt(det(eye(D)+SW)); % in interval [-1,0]
13
14 % 1a. derivatives of expected cost
15 if nargout > 1
16     dLdm = -L*(m-z)'*iSpW; % wrt input mean
17     dLds = L*(iSpW*(m-z)*(m-z)'-eye(D))*iSpW/2; % wrt input covariance matrix
18 end

```

In lines 11–18, we compute the expected cost $L = \mathbb{E}[c(\mathbf{x})]$ and its derivatives with respect to the mean and the covariance matrix of the input distribution. A detailed derivation can be found in [2]. Note that at the moment, $L \in [-1, 0]$ (line 12).

```

19 % 2. Variance of cost
20 if nargout > 3
21     i2SpW = W/(eye(D)+2*SW);
22     r2 = exp(-(m-z)'*i2SpW*(m-z))/sqrt(det(eye(D)+2*SW));

```

```

23 S = r2 - L^2;
24 if S < 1e-12; S=0; end % for numerical reasons
25 end
26
27 % 2a. derivatives of variance of cost
28 if nargout > 4
29 % wrt input mean
30 dSdm = -2*r2*(m-z)'*iSpW-2*L*dLdm;
31 % wrt input covariance matrix
32 dSds = r2*(2*iSpW*(m-z)*(m-z)'-eye(D))*iSpW-2*L*dLds;
33 end

```

In lines 19–33, we compute the variance $V[c(\mathbf{x})]$ of the cost and its derivatives with respect to the mean and the covariance matrix of the input distribution. A detailed derivation can be found in [2]. If the variance $V[c(\mathbf{x})] < 10^{-12}$, we set it to 0 for numerical reasons (line 24).

```

34 % 3. inv(s)*cov(x,L)
35 if nargout > 6
36 t = W*z - iSpW*(SW*z+m);
37 C = L*t;
38 dCdm = t*dLdm - L*iSpW;
39 dCds = -L*(bsxfun(@times,iSpW,permute(t,[3,2,1])) + ...
40 bsxfun(@times,permute(iSpW,[1,3,2]),t'))/2;
41 dCds = bsxfun(@times,t,dLds(:)') + reshape(dCds,D,D^2);
42 end

```

If required, lines 34–42 compute $\mathbf{S}^{-1}\text{cov}[\mathbf{x}, c(\mathbf{x})]$ and the corresponding derivatives with respect to the mean and the covariance of the (Gaussian) state distribution $p(\mathbf{x})$.

```

43 L = 1+L; % bring cost to the interval [0,1]

```

Line 43 brings the expected cost L to the interval $[0, 1]$.

3.3.2.2 Quadratic Cost

`lossQuad` computes the expectation and variance of a quadratic cost

$$c(\mathbf{x}) = (\mathbf{x} - \mathbf{z})^\top \mathbf{W}(\mathbf{x} - \mathbf{z})$$

and their derivatives with respect to the mean and covariance matrix of the (Gaussian) input distribution $p(\mathbf{x})$.

```

1 function [L, dLdm, dLds, S, dSdm, dSds, C, dCdm, dCds] = lossQuad(cost, m, S)

```

Input arguments

<code>cost</code>		
<code>.z:</code>	target state	[D x 1]
<code>.W:</code>	weight matrix	[D x D]
<code>m</code>	mean of input distribution	[D x 1]
<code>s</code>	covariance matrix of input distribution	[D x D]

Output arguments

L	expected loss	$\begin{bmatrix} 1 & x & 1 \end{bmatrix}$
dLdm	derivative of L wrt input mean	$\begin{bmatrix} 1 & x & D \end{bmatrix}$
dLds	derivative of L wrt input covariance	$\begin{bmatrix} 1 & x & D^2 \end{bmatrix}$
S	variance of loss	$\begin{bmatrix} 1 & x & 1 \end{bmatrix}$
dSdm	derivative of S wrt input mean	$\begin{bmatrix} 1 & x & D \end{bmatrix}$
dSds	derivative of S wrt input covariance	$\begin{bmatrix} 1 & x & D^2 \end{bmatrix}$
C	inv(S) times input-output covariance	$\begin{bmatrix} D & x & 1 \end{bmatrix}$
dCdm	derivative of C wrt input mean	$\begin{bmatrix} D & x & D \end{bmatrix}$
dCds	derivative of C wrt input covariance	$\begin{bmatrix} D & x & D^2 \end{bmatrix}$

Implementation

```

2 D = length(m); % get state dimension
3
4 % set some defaults if necessary
5 if isfield(cost, 'W'); W = cost.W; else W = eye(D); end
6 if isfield(cost, 'z'); z = cost.z; else z = zeros(D,1); end

```

In lines 5–6, we check whether the weight matrix \mathbf{W} and the state \mathbf{z} exist. Their default values are \mathbf{I} and $\mathbf{0}$, respectively.

```

7 % 1. expected cost
8 L = S(:)'*W(:) + (z-m)'*W*(z-m);
9
10 % 1a. derivatives of expected cost
11 if nargout > 1
12     dLdm = 2*(m-z)'*W; % wrt input mean
13     dLds = W';          % wrt input covariance matrix
14 end

```

In lines 7–14, we compute the expected cost $L = \mathbb{E}[c(\mathbf{x})]$ and its derivatives with respect to the mean and the covariance matrix of the input distribution. A detailed derivation can be found in [2].

```

15 % 2. variance of cost
16 if nargout > 3
17     S = trace(W*S*(W+W')*S) + (z-m)'*(W+W')*S*(W+W')*(z-m);
18     if S < 1e-12; S = 0; end % for numerical reasons
19 end
20
21 % 2a. derivatives of variance of cost
22 if nargout > 4
23     % wrt input mean
24     dSdm = -(2*(W+W')*S*(W+W)*(z-m))';
25     % wrt input covariance matrix
26     dSds = W'*S'*(W+W')+(W+W')'*S'*W' + (W+W')*(z-m)*((W+W')*(z-m))';
27 end

```

In lines 15–27, we compute the variance $\mathbb{V}[c(\mathbf{x})]$ of the cost and its derivatives with respect to the mean and the covariance matrix of the input distribution. A detailed derivation can be found in [2]. If the variance $\mathbb{V}[c(\mathbf{x})] < 10^{-12}$, we set it to 0 for numerical reasons (line 18).

```

28 % 3. inv(s) times IO covariance with derivatives
29 if nargout > 6
30     C = 2*W*(m-z);
31     dCdm = 2*W;
32     dCds = zeros(D,D^2);
33 end

```

If required, lines 28–33 compute $\mathbf{S}^{-1}\text{cov}[\mathbf{x}, c(\mathbf{x})]$ and the corresponding derivatives with respect to the mean and the covariance of the (Gaussian) state distribution $p(\mathbf{x})$.

Chapter 4

How to Create Your Own Scenario

In this chapter, we explain in sufficient detail how to set up a new scenario by going step-by-step through the cart-pole scenario, which can be found in `<pilco_root>/scenarios/cartPole`.

4.1 Necessary Files

For each scenario, we need the following set of files, which are specific to this scenario. In the cart-pole case, these files are the following:

- `settings_cp.m`: A file that contains scenario-specific settings and initializations
- `loss_cp.m`: A cost function
- `dynamics_cp.m`: A file that implements the ODE, which governs the dynamics.¹
- `learn_cp.m`: A file that puts everything together
- (optional) visualization

4.2 ODE Dynamics

In the following, we briefly describe the interface and the functionality of the cart-pole dynamics. The PILCO code assumes by default that the dynamics model is given by an ODE that is solved numerically using ODE45 (see `<pilco_root>/util/simulate.m` for more details).

```
1 function dz = dynamics_cp(t, z, f)
```

Input arguments:

<code>t</code>	current time step (called from ODE solver)	
<code>z</code>	state	[4 x 1]
<code>f</code>	(optional): force $f(t)$	

The input arguments are as follows:

- `t`: The current time.

¹When working with a real robot, this file is not needed.

- **z**: The state. It is assumed that the state z is given as follows:

$$z = [x, \dot{x}, \dot{\theta}, \theta],$$

where x is the position of the cart (given in m), \dot{x} is the cart's velocity (in m/s), $\dot{\theta}$ is the pendulum's angular velocity in rad/s, and θ is the pendulum's angle in rad. For the angle θ , we chose 0 rad to be the angle when the pendulum hangs downward.

- **f**: The applied force to the cart.

Output arguments:

dz if 3 input arguments: state derivative wrt time
 if only 2 input arguments: total mechanical energy

The function returns either \dot{z} , if three input arguments are given, or the total mechanical energy with two input arguments. The total mechanical energy can be used to verify that the system preserves energy (here, the friction coefficient in line 5 needs to be set to 0).

```
2 l = 0.5; % [m]      length of pendulum
3 m = 0.5; % [kg]     mass of pendulum
4 M = 0.5; % [kg]     mass of cart
5 b = 0.1; % [N/m/s]  coefficient of friction between cart and ground
6 g = 9.82; % [m/s^2] acceleration of gravity
```

In lines 2–6, the parameters of the cart-pole system are defined: the length of the pendulum l (line 2), the mass of the pendulum m (line 3), the mass of the cart M (line 4), the coefficient of friction between cart and ground b (line 5), and the acceleration of gravity g (line 6).

```
7 if nargin==3
8     dz = zeros(4,1);
9     dz(1) = z(2);
10    dz(2) = ( 2*m*l*z(3)^2*sin(z(4)) + 3*m*g*sin(z(4))*cos(z(4)) ...
11             + 4*f(t) - 4*b*z(2) ) / ( 4*(M+m) - 3*m*cos(z(4))^2 );
12    dz(3) = (-3*m*l*z(3)^2*sin(z(4))*cos(z(4)) - 6*(M+m)*g*sin(z(4)) ...
13            - 6*(f(t)-b*z(2))*cos(z(4)) ) / ( 4*l*(m+M) - 3*m*l*cos(z(4))^2 );
14    dz(4) = z(3);
15 else
16    dz = (M+m)*z(2)^2/2 + 1/6*m*l^2*z(3)^2 + m*l*(z(2)*z(3)-g)*cos(z(4))/2;
17 end
```

Lines 7–17 compute either \dot{z} (lines 8–14) or the total mechanical energy (line 16). A principled derivation of the system dynamics and the mechanical energy can be found in [2].

4.3 Scenario-specific Settings

4.3.1 Adding Paths

```
1 rand('seed',1); randn('seed',1); format short; format compact;
2 % include some paths
```

```

3 try
4     rd = '../.. /';
5     addpath([rd 'base'],[rd 'util'],[rd 'gp'],[rd 'control'],[rd 'loss']);
6 catch
7 end

```

First, we include (relative) paths to the directories required for learning and initialize the random seed to make the experiments reproducible. The setting of the random seeds (line 1) will cause some warnings in newer versions of MATLAB, but it is backwards compatible to MATLAB 2007.

4.3.2 Indices

```

8 % 1. Define state and important indices
9
10 % 1a. Full state representation (including all augmentations)
11 %
12 % 1 x          cart position
13 % 2 v          cart velocity
14 % 3 dtheta     angular velocity
15 % 4 theta      angle of the pendulum
16 % 5 sin(theta) complex representation ...
17 % 6 cos(theta) of theta
18 % 7 u          force applied to cart
19 %
20
21 % 1b. Important indices
22 % odei indices for the ode solver
23 % augi indices for variables augmented to the ode variables
24 % dyno indices for the output from the dynamics model and indices to loss
25 % angi indices for variables treated as angles (using sin/cos representation)
26 % dyni indices for inputs to the dynamics model
27 % poli indices for the inputs to the policy
28 % difi indices for training targets that are differences (rather than values)
29
30 odei = [1 2 3 4];           % variables for the ode solver
31 augi = [];                 % variables to be augmented
32 dyno = [1 2 3 4];          % variables to be predicted (and known to loss)
33 angi = [4];               % angle variables
34 dyni = [1 2 3 5 6];        % variables that serve as inputs to the dynamics GP
35 poli = [1 2 3 5 6];        % variables that serve as inputs to the policy
36 difi = [1 2 3 4];          % variables that are learned via differences

```

We now define important state indices that are required by the code that does the actual learning. We assume that the state is given as

$$\mathbf{x} = [x, \dot{x}, \dot{\theta}, \theta]^\top,$$

where x is the position of the cart, \dot{x} the corresponding velocity, and $\dot{\theta}$ and θ are the angular velocity and the angle of the pendulum.

The ODE-solver requires to know what parts of the state are used for the forward dynamics. These indices are captured by `odei` (line 30).

The predictive dimensions of the dynamics GP model are stored in `dyno` (line 32).

The indices in `angi` (line 33) indicate which variables are angles. We represent these angle variables in the complex plane

$$\theta \mapsto [\sin \theta, \cos \theta]^\top$$

to be able to exploit the wrap-around condition $\theta \equiv \theta + 2k\pi$, $k \in \mathbb{Z}$. With this augmentation, we define the auxiliary state vector, i.e., the state vector augmented with the complex representation of the angle, as

$$\mathbf{x} = [x, \dot{x}, \dot{\theta}, \theta, \sin \theta, \cos \theta]^\top. \quad (4.1)$$

The `dyni` indices (line 34) describe which variables from the auxiliary state vector in Equation (4.1) are used as the training inputs for the GP dynamics model. Note that we use the complex representation $[\sin \theta, \cos \theta]$ instead of θ , i.e., we no longer need θ in the inputs of the dynamics GP.

The `poli` indices (line 35) describe which state variables from the auxiliary state vector in Equation (4.1) are used as inputs to the policy.

The `difi` indices (line 36) are a subset of `dyno` and contain the indices of the state variables for which the GP training targets are differences

$$\Delta_t = \mathbf{x}_{t+1} - \mathbf{x}_t$$

instead of \mathbf{x}_{t+1} . Using differences as training targets encodes an implicit prior mean function $m(\mathbf{x}) = \mathbf{x}$. This means that when leaving the training data, the GP predictions do not fall back to 0 but they remain constant. A practical implication is that learning differences Δ_t generalizes better across different parts of the state space. From a learning point of view, training a GP on differences is much simpler than training it on absolute values: The function to be learned does not vary so much, i.e., we do not need so many data points in the end. From a robotics point of view, robot dynamics are typically relative to the current state, they do not so much depend on absolute coordinates.

4.3.3 General Settings

```

37 % 2. Set up the scenario
38 dt = 0.10; % [s] sampling time
39 T = 4.0; % [s] initial prediction horizon time
40 H = ceil(T/dt); % prediction steps (optimization horizon)
41 mu0 = [0 0 0 0]'; % initial state mean
42 S0 = diag([0.1 0.1 0.1 0.1].^2); % initial state covariance
43 N = 15; % number controller optimizations
44 J = 1; % initial J trajectories of length H
45 K = 1; % no. of initial states for which we optimize
46 nc = 100; % number of controller basis functions

```

`dt` is the sampling time, i.e., $1/dt$ is the sampling frequency.

`T` is the length of the prediction horizon in seconds.

`H` = T/dt is the length of the prediction horizon in time steps.

`mu0` is the mean of the distribution $p(\mathbf{x}_0)$ of the initial state. Here, `mu0` = `0` encodes that the cart is in the middle of the track (with zero velocity) with the pendulum hanging down (with zero angular velocity).

`S0` is the covariance matrix of $p(\mathbf{x}_0)$.

`N` is the number of times the loop in Figure 1.1 is executed.

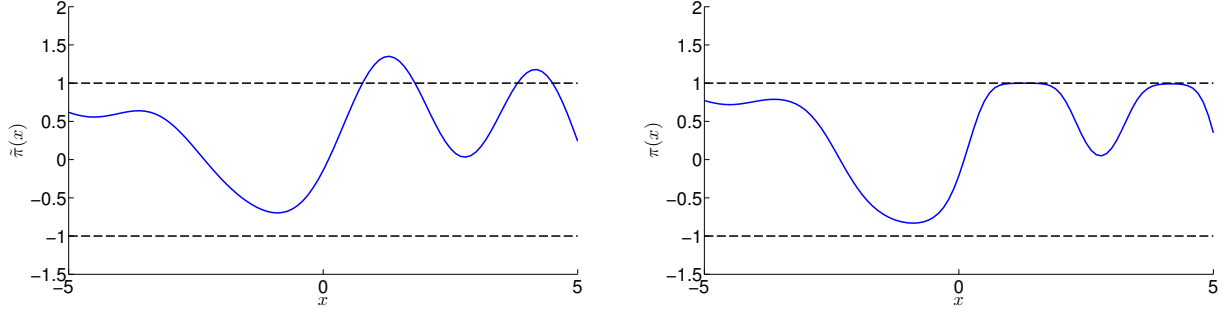


Figure 4.1: Preliminary policy $\tilde{\pi}$ and squashed policy π . The squashing function ensures that the control signals $\mathbf{u} = \pi(\mathbf{x})$ do not exceed the values $\pm \mathbf{u}_{\max}$.

J is the number of initial random rollouts, i.e., rollouts with a random policy. These rollouts are used to collect an initial data set for training the first GP dynamics model. Usually, we set this to 1.

K is the number of initial states for which the policy is learned. The code can manage initial state distributions of the form

$$p(\mathbf{x}_0) \propto \sum_{i=1}^K \mathcal{N}(\boldsymbol{\mu}_0^{(i)}, \boldsymbol{\Sigma}_0), \quad (4.2)$$

which corresponds to a distributions with different means $\boldsymbol{\mu}_0^{(i)}$ but shared covariance matrices $\boldsymbol{\Sigma}_0$.

nc is the number of basis functions of the policy. In this scenario, we use a nonlinear policy of the form

$$\pi(\mathbf{x}) = u_{\max} \sigma \tilde{\pi}(\mathbf{x}) \quad (4.3)$$

$$\tilde{\pi}(\mathbf{x}) = \sum_{i=1}^{\text{nc}} w_i \exp \left(-\frac{1}{2} (\mathbf{x} - \mathbf{c}_i)^\top \mathbf{W} (\mathbf{x} - \mathbf{c}_i) \right), \quad (4.4)$$

where σ is a squashing function, which maps its argument to the interval $[-1, 1]$, \mathbf{c}_i are the locations of the Gaussian-shaped basis functions, and \mathbf{W} is a (shared) weight matrix.

4.3.4 Plant Structure

```

47 % 3. Plant structure
48 plant.dynamics = @dynamics_cp;           % dynamics ode function
49 plant.noise = diag(ones(1,4)*0.01.^2);   % measurement noise
50 plant.dt = dt;
51 plant.ctrl = @zoh;                       % controler is zero order hold
52 plant.odei = odei;
53 plant.augi = augi;
54 plant.angi = angi;
55 plant.poli = poli;
56 plant.dyno = dyno;
57 plant.dyni = dyni;
58 plant.difi = difi;
59 plant.prop = @propagated;

```

`plant.dynamics` requires a function handle to the function that implements the ODE for simulating the system.

`plant.noise` contains the measurement noise covariance matrix. We assume the noise is zero-mean Gaussian. The noise is added to the (latent) state during a trajectory rollout (see `base/rollout.m`).

`plant.dt` is the sampling time, i.e., $1/dt$ is the sampling frequency.

`plant.ctrl` is the controller to be applied. Here, `@zoh` implements a zero-order-hold controller. Other options are `@foh` (first-order-hold) and `@lag` (first-order-lag). For more information, have a look at `base/simulate.m`.

`plant.odei-plant.difi` copy the indices defined earlier into the `plant` structure.

`plant.prop` requires a function handle to the function that computes $p(\mathbf{x}_{t+1})$ from $p(\mathbf{x}_t)$, i.e., a one-step (probabilistic) prediction. In this software package, we implemented a fairly generic function called `propagated`, which computes the predictive state distribution $p(\mathbf{x}_{t+1})$ and the partial derivatives that are required for gradient-based policy search. For details about the gradients, we refer to [3].

4.3.5 Policy Structure

```

60 % 4. Policy structure
61 policy.fcn = @(policy,m,s)conCat(@congp,@gSat,policy,m,s);% controller
62                                     % representation
63 policy.maxU = 10;                    % max. amplitude of
64                                     % control
65 [mm ss cc] = gTrig(mu0, S0, plant.angi); % represent angles
66 mm = [mu0; mm]; cc = S0*cc; ss = [S0 cc; cc' ss]; % in complex plane
67 policy.p.inputs = gaussian(mm(poli), ss(poli,poli), nc)'; % init. location of
68                                     % basis functions
69 policy.p.targets = 0.1*randn(nc, length(policy.maxU)); % init. policy targets
70                                     % (close to zero)
71 policy.p.hyp = log([1 1 1 0.7 0.7 1 0.01])'; % initialize policy
72                                     % hyper-parameters

```

The policy we use in this example is the nonlinear policy given in Equations (4.3)–(4.4). The policy function handle is stored in `policy.fcn` (line 61). In this particular example, the policy is a concatenation of two functions: the RBF controller (`@congp`, see `ctrl/congp.m`), which is parametrized as the mean of a GP with a squared exponential covariance function, and a squashing function σ (`@gSat`, see `<pilco_root>/util/gSat.m`), defined as

$$\sigma(x) = u_{\max} \frac{9 \sin x + \sin(3x)}{8}, \quad (4.5)$$

which is the third-order Fourier series expansion of a trapezoidal wave, normalized to the interval $[-u_{\max}, u_{\max}]$.

`policy.maxU` (line 63) defines the maximum force value u_{\max} in Equation (4.3). We assume that $u \in [-u_{\max}, u_{\max}]$.

In lines 65–66, we augment the original state by $[\sin \theta, \cos \theta]$ by means of `gTrig.m`, where the indices of the angles θ are stored in `plant.angi`. We compute a Gaussian approximation to the joint distribution $p(\mathbf{x}, \sin \theta, \cos \theta)$. The representation of angles θ by means of $[\sin \theta, \cos \theta]$ avoids discontinuities and automatically takes care of the “wrap-around condition”, i.e., $\theta \equiv \theta + 2k\pi, k \in \mathbb{Z}$.

The following lines are used to initialize the policy parameters. The policy parameters are generally stored in `policy.p`. We distinguish between three types of policy parameters:

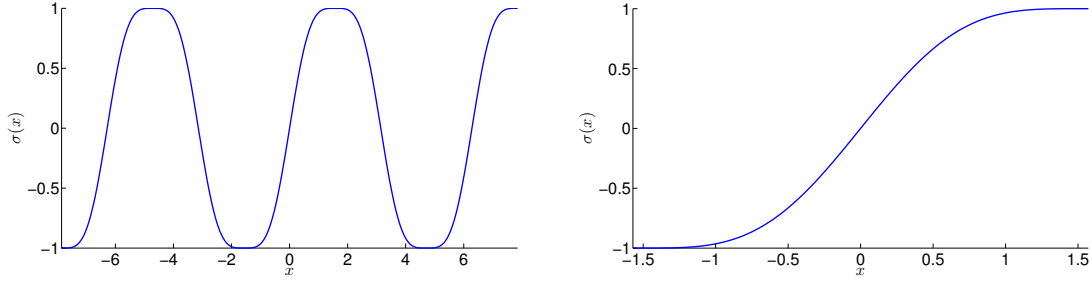


Figure 4.2: Squashing function.

- **policy.p.inputs:** These values play the role of the training inputs of a GP and correspond to the centers \mathbf{c}_i of the policy in Equation (4.4). We sample the initial locations of the centers from the initial state distribution $p(\mathbf{x}_0)$. We compute this initial state distribution in lines 65–66, where we account for possible angle representations (`plant.angi`) of the state. If `plant.angi` is empty, lines 65–66 do not do anything and `mm=mu0` and `ss=S0`.
- **policy.p.targets:** These values play the role of GP training targets are initialized to values close to zero.
- **policy.p.hyp:** These values play the role of the GP log-hyper-parameters: log-length-scales, log-signal-standard-deviation, and log-noise-standard-deviation. We initialize the policy hyper-parameters as follows:
 - Length-scales: The length-scales weight the dimensions of the state that are passed to the policy (see `poli`). In our case these are: $x, \dot{x}, \dot{\theta}, \sin \theta, \cos \theta$. We initialize the first three length-scales to 1. These values largely depend on the scale of the input data. In our example, the cart position and velocity are measured in m and m/s, respectively, the angular velocity is measured in rad/s. The last two length-scales scale trigonometric values, i.e., $\sin(\theta)$ and $\cos(\theta)$. Since these trigonometric functions map their argument into the interval $[-1, 1]$, we choose a length-scale of 0.7, which is somewhat smaller than unity.
 - Signal variance: We set the signal variance of the controller $\tilde{\pi}$ to 1. Note that we squash $\tilde{\pi}$ through σ , see Equation (4.3). To exploit the full potential of the squashing function σ , it is sufficient to cover the domain $[-\pi/2, \pi/2]$. Therefore, we initialize the signal variance to 1.
 - Noise variance: The noise variance is only important only important as a relative factor to the signal variance. This ratio essentially determines how smooth the policy is. We initialize the noise variance to 0.01^2 .

4.3.6 Cost Function Structure

In the following, we set up a structure for the immediate cost function.

```

73 % 5. Set up the cost structure
74 cost.fcn = @loss_cp;           % cost function
75 cost.gamma = 1;               % discount factor
76 cost.p = 0.5;                 % length of pendulum
77 cost.width = 0.25;            % cost function width

```

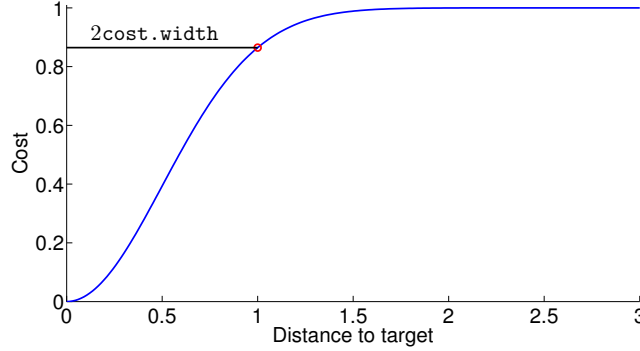


Figure 4.3: In the saturating cost function, see Equation (4.6), `cost.width` determines how far away from the target a cost $c < 1$ can be attained.

```

78 cost.expl = 0.0; % exploration parameter (UCB)
79 cost.angle = plant.angi; % index of angle (for cost function)
80 cost.target = [0 0 0 pi]'; % target state

```

In line 74, we store the function handle `@loss_cp` in `cost.fcn`. The function `loss_cp` implements a saturating cost function (an unnormalized Gaussian subtracted from 1) with spread σ_c , i.e.,

$$c(\mathbf{x}) = 1 - \exp\left(-\frac{1}{2\sigma_c^2}\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2\right) \in [0, 1] \quad (4.6)$$

where $\mathbf{x}_{\text{target}}$ is a target state.

We set the discount factor `cost.gamma` to 1 (line 75) as we look at a finite-horizon problem.

The following parameters are specific to the cost function:

- `cost.p` (line 76) is the length of the pendulum. This length is needed to compute the Euclidean distance of the tip of the pendulum from the desired location in the inverted position.
- `cost.width` (line 77) is the spread/width σ_c of the cost function. Looking at the cost function in (4.6) and the target state $[x, \dot{x}, \theta, \dot{\theta}] = [0, *, *, 2k\pi + \pi]$, $k \in \mathbb{Z}$, the factor 0.25 essentially encodes that the pendulum has to be above horizontal, such that a cost substantially different from 1 is incurred (the length of the pendulum is 0.5 m). Figure 4.3 illustrates a simplified scenario with $\sigma_c = 1/2$. For $c \approx 1$, it can get difficult to obtain any useful gradients. As a rule of thumb, one can set `cost.width` = $\|\boldsymbol{\mu}_0 - \mathbf{x}_{\text{target}}\|/10$.
- `cost.expl` (line 78) is a UCB-type exploration parameter. Negative values encourage exploration, positive values encourage the policy staying in regions with good predictive performance. We set the value to 0 in order to disable any kind of additional exploration or exploitation.
- `cost.angle` (line 79) tells the cost function, which indices of the state are angles. In the cost function, we also represent angles in the complex plane.
- `cost.target` (line 80) defines the target state $\mathbf{x}_{\text{target}}$. Here, the target state is defined as the cart being in the middle of the track and the pendulum upright, without any velocity or angular velocity.

4.3.7 GP Dynamics Model Structure

In the following, we set up the structure `dynmodel` for the GP dynamics model.

```
81 % 6. Dynamics model structure
82 dynmodel.fcn = @gp1d;           % function for GP predictions
83 dynmodel.train = @train;        % function to train dynamics model
84 dynmodel.induce = zeros(300,0,1); % shared inducing inputs (sparse GP)
85 trainOpt = [300 500];          % defines the max. number of line searches
86                               % when training the GP dynamics models
87                               % trainOpt(1): full GP,
88                               % trainOpt(2): sparse GP (FITC)
```

We generally assume that the model uses a squared exponential covariance, a Gaussian likelihood, and a zero prior mean. Therefore, these parameters are not explicitly specified here.

`dynmodel.fcn` (line 82) contains a function handle to `gp1d`, which can predict with (sparse) GPs at uncertain inputs. If the GP is not sparse but full, `gp1d` calls `gp0d`, which implements GP predictions at uncertain inputs with the full GP model.

`dynmodel.train` (line 83) contains a function handle to `train`, which is responsible for GP training.

`dynmodel.induce` (line 84) is optional and tells us when to switch from full GPs to sparse GPs.

`dynmodel.induce` is a tensor of the form $a \times b \times c$, where a is the number of inducing inputs², $b = 0$ tells us that there are no inducing inputs yet, and c is either 1 or the number of predictive dimensions, which corresponds to the number of indices stored in `dyno`. For $c = 1$, the inducing inputs are shared among all GPs. Otherwise, sets of inducing inputs are separately learned for each predictive dimension.

`trainOpt` (line 85) contains the number of line searches for GP hyper-parameter training used by the full GP (first number) and the sparse GP (second parameter).

4.3.8 Optimization Parameters (Policy Learning)

In the following lines, we define (optional) parameters for policy learning. Generally, we use an adapted version of `minimize.m`³, a non-convex gradient-based optimizer. These optional parameters are stored in a structure `opt`.

```
89 % 7. Parameters for policy optimization
90 opt.length = 150;           % max. number of line searches
91 opt.MFEPLS = 30;           % max. number of function evaluations
92                             % per line search
93 opt.verbosity = 1;          % verbosity: specifies how much
94                             % information is displayed during
95                             % policy learning. Options: 0-3
```

`opt.length` (line 90) sets the maximum number of line searches after which the optimizer returns the best parameter set so far.

`opt.MFEPLS` (line 91) is the maximum number of function evaluations per line search. Either the line search succeeds by finding a parameter set with a gradient close to 0 or it does not succeed and aborts after `opt.MFEPLS` many function (and gradient) evaluations.

`opt.verbosity` (line 93) regulates the verbosity of the optimization procedure. Verbosity ranges from 0 (no information displayed) to 3 (visualize the line search and the computed gradients⁴).

²If the size of the training set exceeds a , the full GP automatically switches to its sparse approximation.

³<http://www.gaussianprocess.org/gpml/code/matlab/util/minimize.m>

⁴This is great for debugging.

4.3.9 Plotting Parameters

```

96 % 8. Plotting verbosity
97 plotting.verbosity = 1;           % 0: no plots
98                                % 1: some plots
99                                % 2: all plots

```

`plotting.verbosity` (line 97) is an easy way of controlling how much information is visualized during policy learning.

4.3.10 Allocating Variables

In lines 100-103, we simply initialize a few array that are used to store data during the learning process.

```

100 % 9. Some initializations
101 x = []; y = [];
102 fantasy.mean = cell(1,N); fantasy.std = cell(1,N);
103 realCost = cell(1,N); M = cell(N,1); Sigma = cell(N,1);

```

4.4 Cost Function

In the following, we describe how to define a cost function for the cart-pole swing-up scenario. The cost function is stored in `loss_cp.m` and implements the cost

$$c(\mathbf{x}) = \frac{1}{\#\text{cost.cw}} \sum_{i=1}^{\#\text{cost.cw}} \left(1 - \exp \left(- \frac{1}{2(\sigma_c^{(i)})^2} \|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2 \right) \right) \quad (4.7)$$

which is a generalized version of Equation (4.6). In particular, `cost.cw` can be an array of different widths $\sigma_c^{(i)}$, which is used to compute a cost function $c(\mathbf{x})$ as a mixture of cost functions with different widths.

The mean and the variance of the cost $c(\mathbf{x})$ are computed by averaging over the Gaussian state distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{S})$ with mean \mathbf{m} and covariance matrix \mathbf{S} . Derivatives of the expected cost and the cost variance with respect to the mean and the covariance of the input distribution are computed when desired.

4.4.1 Interface

```

1 function [L, dLdm, dLds, S] = loss_cp(cost, m, s)

```

Input arguments:

<code>cost</code>	cost structure	
<code>.p</code>	length of pendulum	[1 x 1]
<code>.width</code>	array of widths of the cost (summed together)	
<code>.expl</code>	(optional) exploration parameter	
<code>.angle</code>	(optional) array of angle indices	

.target	target state	[D x 1]
m	mean of state distribution	[D x 1]
s	covariance matrix for the state distribution	[D x D]

Output arguments:

L	expected cost	[1 x 1]
dLdm	derivative of expected cost wrt. state mean vector	[1 x D]
dLds	derivative of expected cost wrt. state covariance matrix	[1 x D^2]
S	variance of cost	[1 x 1]

```

2 if isfield(cost, 'width'); cw = cost.width; else cw = 1; end
3 if ~isfield(cost, 'expl') || isempty(cost.expl); b = 0; else b = cost.expl; end

```

In lines 2–3, we check whether a scaling factor (array) and an exploration parameter exist. Default values are 1 (no scaling) and 0 (no exploration), respectively.

```

4 % 1. Some precomputations
5 D0 = size(s,2); % state dimension
6 D1 = D0 + 2*length(cost.angle); % state dimension (with sin/cos)
7
8 M = zeros(D1,1); M(1:D0) = m; S = zeros(D1); S(1:D0,1:D0) = s;
9 Mdm = [eye(D0); zeros(D1-D0,D0)]; Sdm = zeros(D1*D1,D0);
10 Mds = zeros(D1,D0*D0); Sds = kron(Mdm,Mdm);

```

In line 5, the dimension of the state is determined.

In line 6, the dimension of the state is augmented to account for potential angles in the state, which require a representation on the unit circle via $\sin \theta$ and $\cos \theta$. Therefore, the (fully) augmented state variable is then given as

$$\mathbf{j} = [x, \dot{x}, \dot{\theta}, \theta, \sin \theta, \cos \theta]^\top. \quad (4.8)$$

Lines 8–10 initialize the output arguments to 0.

In the following lines, the distance $\mathbf{x} - \mathbf{x}_{\text{target}}$ is computed.

```

11 % 2. Define static penalty as distance from target setpoint
12 e11 = cost.p; % pendulum length
13 Q = zeros(D1); Q([1 D0+1],[1 D0+1]) = [1 e11]'*[1 e11]; Q(D0+2,D0+2) = e11^2;

```

Line 12 stores the pendulum length in `e11`.

In line 12, the matrix \mathbf{Q} is computed, such that $(\mathbf{j} - \mathbf{j}_{\text{target}})^\top \mathbf{Q} (\mathbf{j} - \mathbf{j}_{\text{target}})$ is the squared Euclidean distance between the tip of the pendulum in the current state and the target state. For $\mathbf{x}_{\text{target}} = [0, *, *, \pi]^\top$, i.e., the pendulum is balanced in the inverted position in the middle of the track, the Euclidean distance is given as

$$\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2 = x^2 + 2xl \sin \theta + 2l^2 + 2l^2 \cos \theta = (\mathbf{j} - \mathbf{j}_{\text{target}})^\top \mathbf{Q} (\mathbf{j} - \mathbf{j}_{\text{target}}), \quad (4.9)$$

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 & 0 & 0 & l & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ l & 0 & 0 & 0 & l^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & l^2 \end{bmatrix}. \quad (4.10)$$

Note that at this point only the \mathbf{Q} -matrix is determined.

```

14 % 3. Trigonometric augmentation
15 if D1-D0 > 0
16     % augment target
17     target = [cost.target(:); gTrig(cost.target(:), 0*s, cost.angle)];
18
19     % augment state
20     i = 1:D0; k = D0+1:D1;
21     [M(k) S(k,k) C mdm sdm Cdm mds sds Cds] = gTrig(M(i), S(i,i), cost.angle);
22
23     % compute derivatives (for augmentation)
24     X = reshape(1:D1*D1, [D1 D1]); XT = X'; % vectorized indices
25     I=0*X; I(i,i)=1; ii=X(I==1)'; I=0*X; I(k,k)=1; kk=X(I==1)';
26     I=0*X; I(i,k)=1; ik=X(I==1)'; ki=XT(I==1)';
27
28     Mdm(k,:) = mdm*Mdm(i,:) + mds*Sdm(ii,:); % chainrule
29     Mds(k,:) = mdm*Mds(i,:) + mds*Sds(ii,:);
30     Sdm(kk,:) = sdm*Mdm(i,:) + sds*Sdm(ii,:);
31     Sds(kk,:) = sdm*Mds(i,:) + sds*Sds(ii,:);
32     dCdm = Cdm*Mdm(i,:) + Cds*Sdm(ii,:);
33     dCds = Cdm*Mds(i,:) + Cds*Sds(ii,:);
34
35     S(i,k) = S(i,i)*C; S(k,i) = S(i,k)'; % off-diagonal
36     SS = kron(eye(length(k)), S(i,i)); CC = kron(C', eye(length(i)));
37     Sdm(ik,:) = SS*dCdm + CC*Sdm(ii,:); Sdm(ki,:) = Sdm(ik,:);
38     Sds(ik,:) = SS*dCds + CC*Sds(ii,:); Sds(ki,:) = Sds(ik,:);
39 end

```

This block is only executed if angles are present (the check is performed in line 15).

First (line 17), the target state $\mathbf{x}_{\text{target}}$ is augmented to

$$\mathbf{j}_{\text{target}} = [\mathbf{x}_{\text{target}}, \dot{\mathbf{x}}_{\text{target}}, \dot{\theta}_{\text{target}}, \theta_{\text{target}}, \sin(\theta_{\text{target}}), \cos(\theta_{\text{target}})]^T.$$

In line 21, the state distribution $p(\mathbf{x})$ is probabilistically augmented to $p(\mathbf{j})$, where \mathbf{j} is given in Equation (4.1). Note that $p(\mathbf{j})$ cannot be computed analytically. Instead, we compute the mean and covariance of $p(\mathbf{j})$ exactly and approximate $p(\mathbf{j})$ by a Gaussian. This probabilistic augmentation and the corresponding derivatives with respect to the mean and covariance of the state distribution are computed by `<pilco_root>/util/gTrig.m`.

Lines 28–38 compute the derivatives of the mean and covariance of $p(\mathbf{j})$ and the cross-covariance $\text{cov}[\mathbf{x}, \mathbf{j}]$ with respect to the mean and covariance of $p(\mathbf{x})$ using the chain rule.

```

40 % 4. Calculate loss
41 L = 0; dLdm = zeros(1,D0); dLds = zeros(1,D0*D0); S = 0;
42 for i = 1:length(cw) % scale mixture of immediate costs
43     cost.z = target; cost.W = Q/cw(i)^2;
44     [r rdM rdS s2 s2dM s2dS] = lossSat(cost, M, S);
45
46     L = L + r; S = S + s2;
47     dLdm = dLdm + rdM(:)'*Mdm + rdS(:)'*Sdm;
48     dLds = dLds + rdM(:)'*Mds + rdS(:)'*Sds;
49
50     if (b~=0 || ~isempty(b)) && abs(s2)>1e-12
51         L = L + b*sqrt(s2);
52         dLdm = dLdm + b/sqrt(s2) * ( s2dM(:)'*Mdm + s2dS(:)'*Sdm )/2;

```

```

53     dLds = dLds + b/sqrt(s2) * ( s2dM(:)'*Mds + s2dS(:)'*Sds )/2;
54     end
55 end
56
57 % normalize
58 n = length(cw); L = L/n; dLdm = dLdm/n; dLds = dLds/n; S = S/n;

```

After all the pre-computations, in lines 40–58, the expected cost for Equation (4.7) is finally computed: For all widths of the cost structure (line 42), we compute the mean and the variance of the saturating cost in Equation (4.6), including the derivatives with respect to the mean and the covariance of $p(j)$, see line 44. For these computations, the function `loss/lossSat.m` is called.

Lines 47–48 compute the derivatives of the expected cost and the variance of the cost with respect to the mean and covariance matrix of the state distribution $p(\mathbf{x})$ by applying the chain rule.

If exploration is desired (line 50), we add $\kappa\sqrt{V[c(\mathbf{x})]}$ to the $E[c(\mathbf{x})]$ to allow for UCB-type exploration, see Equation (3.1)

4.5 Visualization

The following lines of code display a sequence of images (video) of a cart-pole trajectory and can be found in `<pilco_root>/scenarios/cartPole/draw_rollout_cp.m`.

```

1 % Loop over states in trajectory (= time steps)
2 for r = 1:size(xx,1)
3     if exist('j','var') && ~isempty(M{j})
4         draw_cp(latent{j}(r,1), latent{j}(r,4), latent{j}(r,end), cost, ...
5             ['trial # ' num2str(j+J) ', T=' num2str(H*dt) ' sec'], ...
6             ['total experience (after this trial): ' num2str(dt*size(x,1)) ...
7             ' sec'], M{j}(:,r), Sigma{j}(:, :, r));
8     else
9         draw_cp(latent{jj}(r,1), latent{jj}(r,4), latent{jj}(r,end), cost, ...
10            ['(random) trial # ' num2str(jj) ', T=' num2str(H*dt) ' sec'], ...
11            ['total experience (after this trial): ' num2str(dt*size(x,1)) ...
12            ' sec'])
13     end
14     pause(dt);
15 end

```

At each time step r of the most recent trajectory (stored in `xx`), an image of the current cart-pole state is drawn by repeatedly calling `draw_cp`. We distinguish between two modes (`if-else` statement): Either we draw a trajectory *after* having learned a policy (`if` statement), or we draw a trajectory from a random rollout, i.e., before a policy is learned (`else` statement). In the first case, `draw_cp` also visualizes the long-term predictive means and variances of the tip of the pendulum, which are not given otherwise.

The `draw_cp` function plots the cart-pole system with reward, applied force, and predictive uncertainty of the tip of the pendulum. We just describe the interface in the following. The code is available at `<pilco_root>/scenarios/cartPole/draw_cp.m`.

```

1 function draw_cp(x, theta, force, cost, text1, text2, M, S)

```

Input arguments:

x	position of the cart
theta	angle of pendulum
force	force applied to cart
cost	cost structure
.fcfn	function handle (it is assumed to use saturating cost)
.<>	other fields that are passed to cost
M	(optional) mean of state
S	(optional) covariance of state
text1	(optional) text field 1
text2	(optional) text field 2

4.6 Main Function

The main function executes the following high-level steps.

1. Load scenario-specific setting
2. Create J initial trajectories by applying random controls
3. Controlled learning:
 - (a) Train dynamics model
 - (b) Learn policy
 - (c) Apply policy to system

```

1 % 1. Initialization
2 clear all; close all;
3 settings_cp; % load scenario-specific settings
4 basename = 'cartPole-'; % filename used for saving data

```

Lines 1–4 load scenario-specific settings (line 3) and define a basename for data that is stored throughout the execution.

```

5 % 2. Initial J random rollouts
6 for jj = 1:J
7     [xx, yy, realCost{jj}, latent{jj}] = ...
8     rollout(gaussian(mu0, S0), struct('maxU', policy.maxU), H, plant, cost);
9     x = [x; xx]; y = [y; yy]; % augment training sets for dynamics model
10    if plotting.verbosity > 0; % visualization of trajectory
11        if ~ishandle(1); figure(1); else set(0, 'CurrentFigure', 1); end; clf(1);
12        draw_rollout_cp;
13    end
14
15 end
16
17 mu0Sim(odei,:) = mu0; S0Sim(odei,odei) = S0;
18 mu0Sim = mu0Sim(dyno); S0Sim = S0Sim(dyno,dyno);

```

To kick off learning, we need to create an initial (small) data set that can be learned for learning the first GP dynamics model. For this, we generate J trajectories of length H by applying random control

signals using `<pilco_root>/base/rollout.m`, (lines 7–8). Generally, the training data for the GP is stored in `x` and `y` (line 9). If desired, the trajectories of the cart-pole system are visualized (lines 10–13). In lines 17–18, we define variables `muOSim` and `SOSim`, which are used subsequently.

```

19 % 3. Controlled learning (N iterations)
20 for j = 1:N
21     trainDynModel; % train (GP) dynamics model
22     learnPolicy; % learn policy
23     applyController; % apply controller to system
24     disp(['controlled trial # ' num2str(j)]);
25     if plotting.verbosity > 0; % visualization of trajectory
26         if ~ishandle(1); figure(1); else set(0, 'CurrentFigure', 1); end; clf(1);
27         draw_rollout_cp;
28     end
29 end

```

The actual learning happens in lines 19–29, where PILCO performs `N` (controlled) iterations of dynamics model learning (line 21), policy search (line 22), and controller application to the system (line 23). If desired, the trajectories of the cart-pole system are visualized (lines 25–28).

4.6.1 Screen Prints and Visualization

When training the GP dynamics model, a typical feedback in the command window looks as follows:

```

Train hyper-parameters of full GP ...
GP 1/4
Initial Function Value 1.853671e+01
linesearch #      31; value -8.942969e+01
GP 2/4
Initial Function Value 3.115190e+01
linesearch #      34; value -4.210842e+01
GP 3/4
Initial Function Value 8.045295e+01
linesearch #      30; value 4.742728e+00
GP 4/4
Initial Function Value -3.771471e+00
linesearch #      37; value -6.971879e+01
Learned noise std: 0.016818    0.016432    0.04385    0.019182
SNRs               : 28.91172    116.4612    112.2714    49.12984

```

In this cart-pole example, we train four GPs (one for each predictive dimension stored in `dyno`). The hyper-parameters are learned after 30–40 line searches. Positive values are generally an indicator that the learned model is not so good at predicting this target dimension. This might be due to sparse data and/or very nonlinear dynamics. At the end, the learned noise standard deviations are displayed, together with the signal-to-noise ratios (SNRs). The SNRs are computed as

$$SNR = \frac{\sqrt{\sigma_f^2}}{\sqrt{\sigma_{\text{noise}}^2}},$$

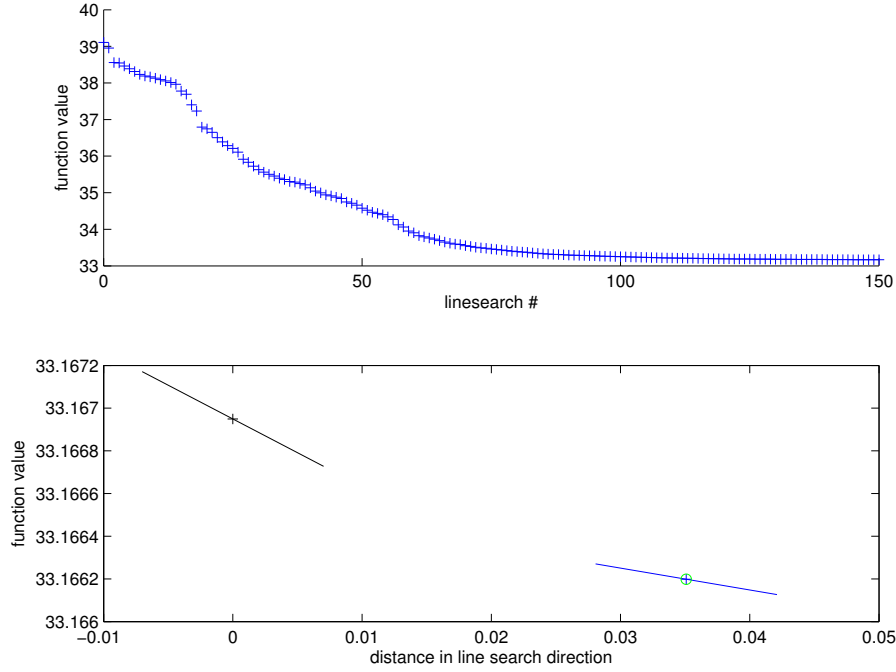


Figure 4.4: Typical display during policy learning: The top subplot shows the overall progress of policy learning as a function of the number of line searches. In this example, the initial policy caused a total expected cost $J(\theta)$ of 39.11; the policy after 150 line search searches caused an expected cost of 33.16. The bottom subplot shows the progress of the current line search as a function of the distance in line search direction. The function values (+) and the corresponding gradients are visualized. For more detailed information about the visualization and the verbosity of the output (`opt.verbosity`), we refer to `<pilco_root>/util/minimize.m`.

where σ_f^2 is the variance of the underlying function and σ_{noise}^2 is the inferred noise variance. High SNRs (> 500) are penalized during GP training in `<pilco_root>/gp/hypCurb.m` to improve numerical stability. If there are SNRs > 500 , it might be worthwhile adding some more noise to the GP training targets.

Figure 4.4 displays a typical screen output during policy optimization, showing the overall progress of policy learning (top subplot) and the progress of the current line search (bottom subplot). The following lines are simultaneously displayed in the command window to show the current progress of learning:

```
Initial Function Value 3.910902e+01
line search #      150;  value 3.316620e+01
```

If `opt.verbosity` < 3 , i.e., we are not interested in permanently displaying the gradients as in Figure 4.4, it is possible to display the overall optimization progress at the end of each policy search by setting `plotting.verbosity` = 2. The corresponding figure is given in Figure 4.5. This figure visually indicates whether the policy search was close to convergence. If it frequently happens that the curve does not flatten out, it might be worthwhile increasing the value of `opt.length` in the scenario-specific settings file (here: `settings_cp.m`).

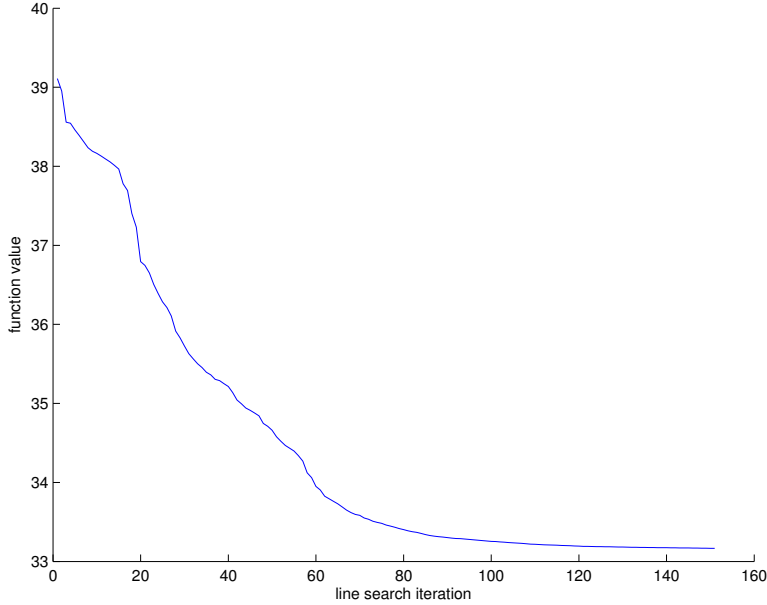


Figure 4.5: Overall progress of a policy search with 150 line searches.

Figure 4.6 shows the predicted and the incurred immediate cost when applying the policy. Initially, the cost is at unity, which means that the initial state is far from the target area. After around 6 time steps, the predictive uncertainty in the cost increases, which means that the predictive state distribution substantially overlaps with the target area. In Figure 4.6 the reason is that the predictive state increases very quickly during the first time steps.

Figure 4.7 shows two sets of example trajectories for each predictive dimension. One set is generated in the early stages of learning (Figure 4.7(a)), the other one is generated in the final stages of learning (Figure 4.7(b)). The green trajectory is stored in `xx` and will be used to augment the training data set for the GP model. The red trajectories are generated to (a) indicate whether the green trajectory is a “typical” trajectory or an outlier, (b) show the quality of the long-term predictive state distributions, whose 95% predictive confidence intervals are indicated by the blue error bars, (c) give an intuition how sensitive the currently learned controller is to different initial states $\mathbf{x}_0 \sim p(\mathbf{x}_0)$. As shown in Figure 4.7(a), in the early stages of learning, the controller is very sensitive to the initial conditions, i.e., the controlled trajectories vary a lot. A good controller is robust to the uncertain initial conditions and leads to very similar controlled trajectories as shown in Figure 4.7(b).

Figure 4.8 shows a snapshot of a visualization of a trajectory of the cart-pole system. The cart runs on an (infinite) track, the pendulum is freely swinging. The cross denotes the target location for balancing the tip of the pendulum, and the blue ellipse represents the 95% confidence bound of a t -step ahead prediction of the location of the tip of the pendulum when applying the learned policy. We also display the applied force (green bar), whose values u_{\max} and $-u_{\max}$ correspond to a green bar being either full to the right or left side. Here, “full” means at the end of the black line that represents the track. Moreover, we show the incurred immediate reward (negative cost) of the current state of the system. The immediate reward is represented by a yellow bar, whose maximum is at the right end of the black bar. We also display the number of total trials (this includes the random initial trials), the prediction horizon T , and the total experience after this trial.

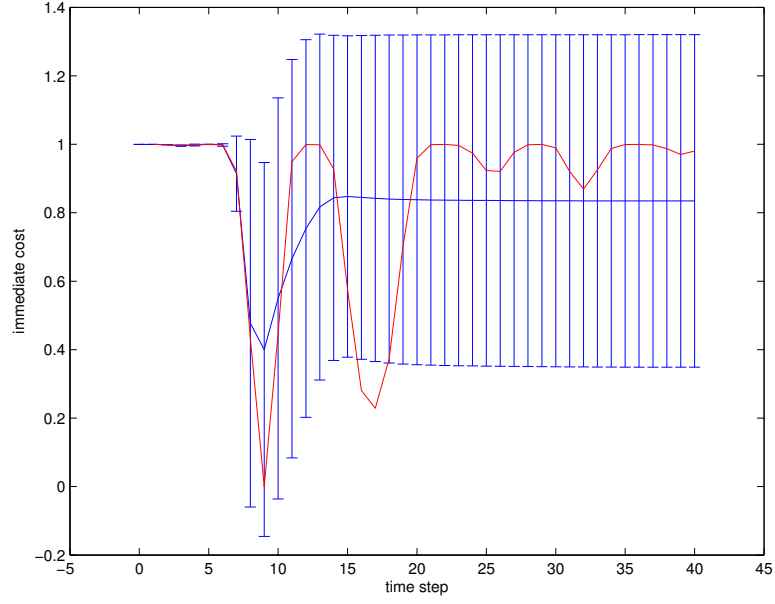


Figure 4.6: Predicted immediate cost (blue error bar) and corresponding incurred cost when applying the policy (red) for each time step of the prediction horizon.

Table 4.1: `plotting.verbosity` Overview.

	Figure 4.5	Figure 4.6	Figure 4.7	Figure 4.8
<code>plotting.verbosity=0</code>	✗	✗	✗	✗
<code>plotting.verbosity=1</code>	✗	✓	✗	✓
<code>plotting.verbosity=2</code>	✓	✓	✓	✓

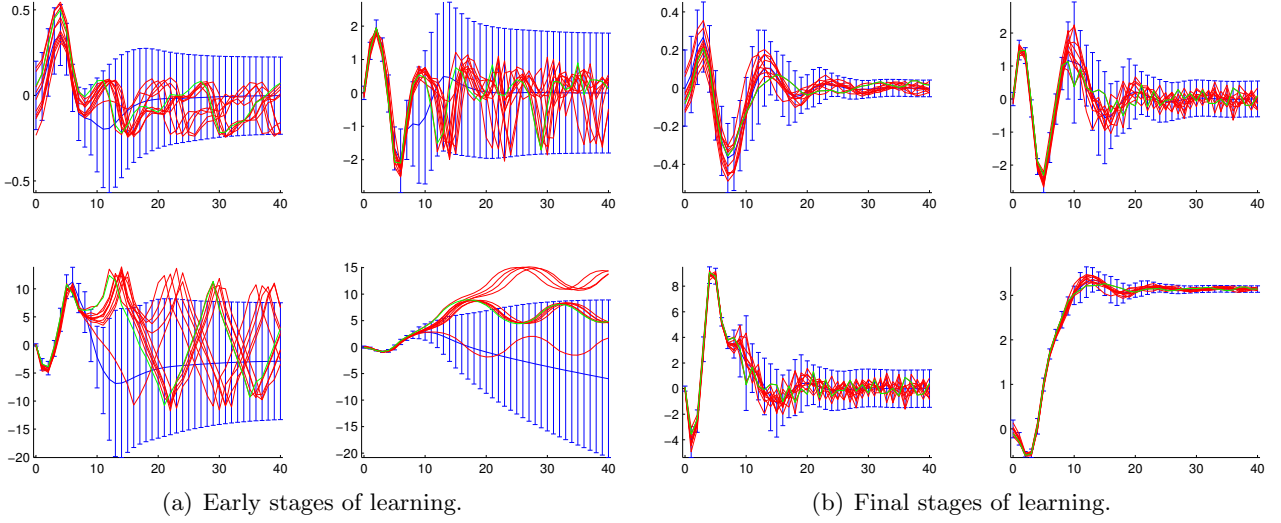


Figure 4.7: Example trajectories. For each predictive dimension (**dyno**), ten state trajectories are displayed, which occurred when the same policy $\pi(\theta^*)$ is applied to the system. Note that the initial state differs as it is sampled from the start state distribution, i.e., $\mathbf{x}_0 \sim p(\mathbf{x}_0)$. The green trajectory will be used in the next learning iteration to update the GP training set; the other (red) trajectories are generated to get an impression of the quality of the long-term predictive state distributions, the 95% confidence intervals of which are shown by the blue error bars.

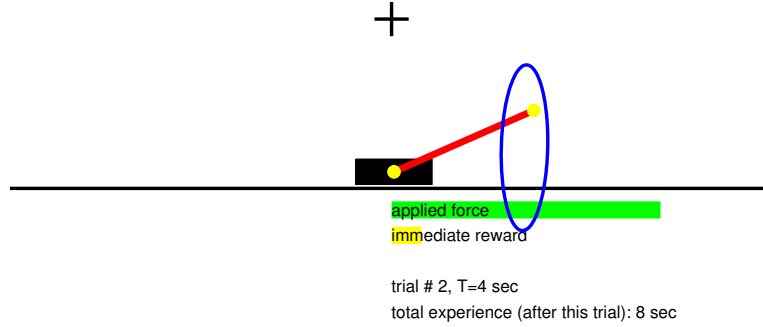


Figure 4.8: Visualization of a trajectory of the cart-pole swing-up during learning. The cart runs on an (infinite) track, the pendulum is freely swinging. The cross denotes the target location for balancing the tip of the pendulum, and the blue ellipse represents the 95% confidence bound of a t -step ahead prediction of the location of the tip of the pendulum when applying the learned policy. We also display the applied force (green bar), whose values u_{\max} and $-u_{\max}$ correspond to a green bar being either full to the right or left side. Here, “full” means at the end of the black line that represents the track. Moreover, we show the incurred immediate reward (negative cost) of the current state of the system. The immediate reward is represented by a yellow bar, whose maximum is at the right end of the black bar. We also display the number of total trials (this includes the random initial trials), the prediction horizon T , and the total experience after this trial.

Chapter 5

Implemented Scenarios

In the following, we introduce the scenarios that are shipped with this software package, and detail the derivation of the corresponding equations of motion, taken from [2]. All scenarios have their own folder in `<pilco.root>/scenarios/`.

Table 5.1: State and control space dimensions in the implemented scenarios.

Task	State Space	Control Space
Pendulum	\mathbb{R}^2	\mathbb{R}
Pendubot	\mathbb{R}^4	\mathbb{R}
Double Pendulum	\mathbb{R}^4	\mathbb{R}^2
Cart-Pole	\mathbb{R}^4	\mathbb{R}
Cart-Double Pendulum	\mathbb{R}^6	\mathbb{R}
Unicycling	\mathbb{R}^{12}	\mathbb{R}^2

Table 5.1 gives an overview of the corresponding state and control dimensions.

5.1 Pendulum Swing-up

The pendulum shown in Figure 5.1 possesses a mass m and a length l . The pendulum angle φ is measured anti-clockwise from hanging down. A torque u can be applied to the pendulum. Typical values are: $m = 1$ kg and $l = 1$ m.

The coordinates x and y of the midpoint of the pendulum are

$$\begin{aligned}x &= \frac{1}{2}l \sin \varphi, \\y &= -\frac{1}{2}l \cos \varphi,\end{aligned}$$

and the squared velocity of the midpoint of the pendulum is

$$v^2 = \dot{x}^2 + \dot{y}^2 = \frac{1}{4}l^2 \dot{\varphi}^2.$$

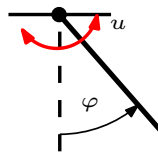


Figure 5.1: Pendulum.

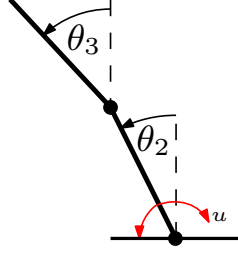


Figure 5.2: Pendubot.

We derive the equations of motion via the system Lagrangian L , which is the difference between kinetic energy T and potential energy V and given by

$$L = T - V = \frac{1}{2}mv^2 + \frac{1}{2}I\dot{\varphi}^2 + \frac{1}{2}mgl \cos \varphi, \quad (5.1)$$

where $g = 9.82 \text{ m/s}^2$ is the acceleration of gravity and $I = \frac{1}{12}ml^2$ is the moment of inertia of a pendulum around the pendulum midpoint.

The equations of motion can generally be derived from a set of equations defined through

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} - \frac{\partial L}{\partial q_i} = Q_i,$$

where Q_i are the non-conservative forces and q_i and \dot{q}_i are the state variables of the system. In our case,

$$\begin{aligned} \frac{\partial L}{\partial \dot{\varphi}} &= \frac{1}{4}ml^2\dot{\varphi} + I\dot{\varphi} \\ \frac{\partial L}{\partial \varphi} &= -\frac{1}{2}mgl \sin \varphi \end{aligned}$$

yield

$$\ddot{\varphi}(\frac{1}{4}ml^2 + I) + \frac{1}{2}mgl \sin \varphi = u - b\dot{\varphi},$$

where b is a friction coefficient. Collecting both variables $\mathbf{z} = [\dot{\varphi}, \varphi]^\top$ the equations of motion can be conveniently expressed as two coupled ordinary differential equations

$$\frac{d\mathbf{z}}{dt} = \begin{bmatrix} u - bz_1 - \frac{1}{2}mgl \sin z_2 \\ \frac{1}{4}ml^2 + I \\ z_1 \end{bmatrix},$$

which can be simulated numerically.

5.2 Double Pendulum Swing-up with a Single Actuator (Pendubot)

The Pendubot in Figure 5.2 is a two-link (mass m_2 and m_3 and length l_2 and l_3 respectively), underactuated robot as described by [9]. The first joint exerts torque, but the second joint cannot. The system has four continuous state variables: two joint positions and two joint velocities. The angles of the joints, θ_2 and θ_3 , are measured anti-clockwise from upright. An applied external torque u controls the first joint. Typical values are: $m_2 = 0.5\text{kg}$, $m_3 = 0.5\text{kg}$, $l_2 = 0.6\text{m}$, $l_3 = 0.6\text{m}$.

The Cartesian coordinates x_2, y_2 and x_3, y_3 of the midpoints of the pendulum elements are

$$\begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} -\frac{1}{2}l_2 \sin \theta_2 \\ \frac{1}{2}l_2 \cos \theta_2 \end{bmatrix}, \quad \begin{bmatrix} x_3 \\ y_3 \end{bmatrix} = \begin{bmatrix} -l_2 \sin \theta_2 - \frac{1}{2}l_3 \sin \theta_3 \\ l_2 \cos \theta_2 + \frac{1}{2}l_3 \cos \theta_3 \end{bmatrix}, \quad (5.2)$$

and the squared velocities of the pendulum midpoints are

$$v_2^2 = \dot{x}_2^2 + \dot{y}_2^2 = \frac{1}{4}l_2^2\dot{\theta}_2^2 \quad (5.3)$$

$$v_3^2 = \dot{x}_3^2 + \dot{y}_3^2 = l_2^2\dot{\theta}_2^2 + \frac{1}{4}l_3^2\dot{\theta}_3^2 + l_2l_3\dot{\theta}_2\dot{\theta}_3 \cos(\theta_2 - \theta_3). \quad (5.4)$$

The system Lagrangian is the difference between the kinematic energy T and the potential energy V and given by

$$L = T - V = \frac{1}{2}m_2v_2^2 + \frac{1}{2}m_3v_3^2 + \frac{1}{2}I_2\dot{\theta}_2^2 + \frac{1}{2}I_3\dot{\theta}_3^2 - m_2gy_2 - m_3gy_3,$$

where the angular moment of inertia around the pendulum midpoint is $I = \frac{1}{12}ml^2$, and $g = 9.82\text{m/s}^2$ is the acceleration of gravity. Using this moment of inertia, we assume that the pendulum is a thin (but rigid) wire. Plugging in the squared velocities (5.3) and (5.4), we obtain

$$\begin{aligned} L = & \frac{1}{8}m_2l_2^2\dot{\theta}_2^2 + \frac{1}{2}m_3(l_2^2\dot{\theta}_2^2 + \frac{1}{4}l_3^2\dot{\theta}_3^2 + l_2l_3\dot{\theta}_2\dot{\theta}_3 \cos(\theta_2 - \theta_3)) \\ & + \frac{1}{2}I_2\dot{\theta}_2^2 + \frac{1}{2}I_3\dot{\theta}_3^2 - \frac{1}{2}m_2gl_2 \cos \theta_2 - m_3g(l_2 \cos \theta_2 + \frac{1}{2}l_3 \cos \theta_3). \end{aligned}$$

The equations of motion are

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} - \frac{\partial L}{\partial q_i} = Q_i, \quad (5.5)$$

where Q_i are the non-conservative forces and q_i and \dot{q}_i are the state variables of the system. In our case,

$$\begin{aligned} \frac{\partial L}{\partial \dot{\theta}_2} &= l_2^2\dot{\theta}_2(\frac{1}{4}m_2 + m_3) + \frac{1}{2}m_3l_2l_3\dot{\theta}_3 \cos(\theta_2 - \theta_3) + I_2\dot{\theta}_2, \\ \frac{\partial L}{\partial \theta_2} &= -\frac{1}{2}m_3l_2l_3\dot{\theta}_2\dot{\theta}_3 \sin(\theta_2 - \theta_3) + (\frac{1}{2}m_2 + m_3)gl_2 \sin \theta_2, \\ \frac{\partial L}{\partial \dot{\theta}_3} &= m_3l_3(\frac{1}{4}l_3\dot{\theta}_3 + \frac{1}{2}l_2\dot{\theta}_2 \cos(\theta_2 - \theta_3)) + I_3\dot{\theta}_3, \\ \frac{\partial L}{\partial \theta_3} &= \frac{1}{2}m_3l_3(l_2\dot{\theta}_2\dot{\theta}_3 \sin(\theta_2 - \theta_3) + g \sin \theta_3) \end{aligned}$$

lead to the equations of motion

$$\begin{aligned} u &= \ddot{\theta}_2(l_2^2(\frac{1}{4}m_2 + m_3) + I_2) + \ddot{\theta}_3\frac{1}{2}m_3l_3l_2 \cos(\theta_2 - \theta_3) \\ &\quad + l_2(\frac{1}{2}m_3l_3\dot{\theta}_3^2 \sin(\theta_2 - \theta_3) - g \sin \theta_2(\frac{1}{2}m_2 + m_3)), \\ 0 &= \ddot{\theta}_2\frac{1}{2}l_2l_3m_3 \cos(\theta_2 - \theta_3) + \ddot{\theta}_3(\frac{1}{4}m_3l_3^2 + I_3) - \frac{1}{2}m_3l_3(l_2\dot{\theta}_2^2 \sin(\theta_2 - \theta_3) + g \sin \theta_3). \end{aligned}$$

To simulate the system numerically, we solve the linear equation system

$$\begin{bmatrix} l_2^2(\frac{1}{4}m_2 + m_3) + I_2 & \frac{1}{2}m_3l_3l_2 \cos(\theta_2 - \theta_3) \\ \frac{1}{2}l_2l_3m_3 \cos(\theta_2 - \theta_3) & \frac{1}{4}m_3l_3^2 + I_3 \end{bmatrix} \begin{bmatrix} \ddot{\theta}_2 \\ \ddot{\theta}_3 \end{bmatrix} = \begin{bmatrix} c_2 \\ c_3 \end{bmatrix}$$

for $\ddot{\theta}_2$ and $\ddot{\theta}_3$, where

$$\begin{bmatrix} c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} -l_2(\frac{1}{2}m_3l_3\dot{\theta}_3^2 \sin(\theta_2 - \theta_3) - g \sin \theta_2(\frac{1}{2}m_2 + m_3)) + u \\ \frac{1}{2}m_3l_3(l_2\dot{\theta}_2^2 \sin(\theta_2 - \theta_3) + g \sin \theta_3) \end{bmatrix}.$$

5.3 Double Pendulum Swing-up with Two Actuators

The dynamics of the double pendulum with two actuators (one at the shoulder, one at the elbow), are derived exactly as described in Section 5.2, with the single modification that we need to take the second control signal into account in the equations of motion

$$\begin{aligned} u_1 &= \ddot{\theta}_2 \left(l_2^2 \left(\frac{1}{4} m_2 + m_3 \right) + I_2 \right) + \ddot{\theta}_3 \frac{1}{2} m_3 l_3 l_2 \cos(\theta_2 - \theta_3) \\ &\quad + l_2 \left(\frac{1}{2} m_3 l_3 \dot{\theta}_3^2 \sin(\theta_2 - \theta_3) - g \sin \theta_2 \left(\frac{1}{2} m_2 + m_3 \right) \right), \\ u_2 &= \ddot{\theta}_2 \frac{1}{2} l_2 l_3 m_3 \cos(\theta_2 - \theta_3) + \ddot{\theta}_3 \left(\frac{1}{4} m_3 l_3^2 + I_3 \right) - \frac{1}{2} m_3 l_3 \left(l_2 \dot{\theta}_2^2 \sin(\theta_2 - \theta_3) + g \sin \theta_3 \right). \end{aligned}$$

To simulate the system numerically, we solve the linear equation system

$$\begin{bmatrix} l_2^2 \left(\frac{1}{4} m_2 + m_3 \right) + I_2 & \frac{1}{2} m_3 l_3 l_2 \cos(\theta_2 - \theta_3) \\ \frac{1}{2} l_2 l_3 m_3 \cos(\theta_2 - \theta_3) & \frac{1}{4} m_3 l_3^2 + I_3 \end{bmatrix} \begin{bmatrix} \ddot{\theta}_2 \\ \ddot{\theta}_3 \end{bmatrix} = \begin{bmatrix} c_2 \\ c_3 \end{bmatrix}$$

for $\ddot{\theta}_2$ and $\ddot{\theta}_3$, where

$$\begin{bmatrix} c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} -l_2 \left(\frac{1}{2} m_3 l_3 \dot{\theta}_3^2 \sin(\theta_2 - \theta_3) - g \sin \theta_2 \left(\frac{1}{2} m_2 + m_3 \right) \right) + u_1 \\ \frac{1}{2} m_3 l_3 \left(l_2 \dot{\theta}_2^2 \sin(\theta_2 - \theta_3) + g \sin \theta_3 \right) + u_2 \end{bmatrix}.$$

5.4 Cart-Pole Swing-up

The cart-pole system (inverted pendulum) shown in Figure 5.3 consists of a cart with mass m_1 and an attached pendulum with mass m_2 and length l , which swings freely in the plane. The pendulum angle θ_2 is measured anticlockwise from hanging down. The cart can move horizontally with an applied external force u and a parameter b , which describes the friction between cart and ground. Typical values are: $m_1 = 0.5 \text{ kg}$, $m_2 = 0.5 \text{ kg}$, $l = 0.6 \text{ m}$ and $b = 0.1 \text{ N/m/s}$.

The position of the cart along the track is denoted by x_1 . The coordinates x_2 and y_2 of the midpoint of the pendulum are

$$\begin{aligned} x_2 &= x_1 + \frac{1}{2} l \sin \theta_2, \\ y_2 &= -\frac{1}{2} l \cos \theta_2, \end{aligned}$$

and the squared velocity of the cart and the midpoint of the pendulum are

$$\begin{aligned} v_1^2 &= \dot{x}_1^2 \\ v_2^2 &= \dot{x}_2^2 + \dot{y}_2^2 = \dot{x}_1^2 + \frac{1}{4} l^2 \dot{\theta}_2^2 + l \dot{x}_1 \dot{\theta}_2 \cos \theta_2, \end{aligned}$$

respectively. We derive the equations of motion via the system Lagrangian L , which is the difference between kinetic energy T and potential energy V and given by

$$L = T - V = \frac{1}{2} m_1 v_1^2 + \frac{1}{2} m_2 v_2^2 + \frac{1}{2} I \dot{\theta}_2^2 - m_2 g y_2, \quad (5.6)$$

where $g = 9.82 \text{ m/s}^2$ is the acceleration of gravity and $I = \frac{1}{12} m l^2$ is the moment of inertia of a pendulum around the pendulum midpoint. Plugging this value for I into the system Lagrangian (5.6), we obtain

$$L = \frac{1}{2} (m_1 + m_2) \dot{x}_1^2 + \frac{1}{6} m_2 l^2 \dot{\theta}_2^2 + \frac{1}{2} m_2 l (\dot{x}_1 \dot{\theta}_2 + g) \cos \theta_2.$$

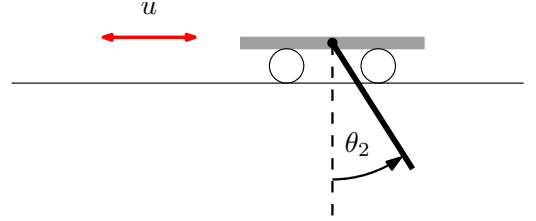


Figure 5.3: Cart-pole system.

The equations of motion can generally be derived from a set of equations defined through

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} - \frac{\partial L}{\partial q_i} = Q_i, \quad (5.7)$$

where Q_i are the non-conservative forces and q_i and \dot{q}_i are the state variables of the system. In our case,

$$\begin{aligned} \frac{\partial L}{\partial \dot{x}_1} &= (m_1 + m_2)\dot{x}_1 + \frac{1}{2}m_2l\dot{\theta}_2 \cos \theta_2, \\ \frac{\partial L}{\partial x_1} &= 0, \\ \frac{\partial L}{\partial \dot{\theta}_2} &= \frac{1}{3}m_2l^2\dot{\theta}_2 + \frac{1}{2}m_2l\dot{x}_1 \cos \theta_2, \\ \frac{\partial L}{\partial \theta_2} &= -\frac{1}{2}m_2l(\dot{x}_1\dot{\theta}_2 + g), \end{aligned}$$

lead to the equations of motion

$$\begin{aligned} (m_1 + m_2)\ddot{x}_1 + \frac{1}{2}m_2l\ddot{\theta}_2 \cos \theta_2 - \frac{1}{2}m_2l\dot{\theta}_2^2 \sin \theta_2 &= u - b\dot{x}_1, \\ 2l\ddot{\theta}_2 + 3\ddot{x}_1 \cos \theta_2 + 3g \sin \theta_2 &= 0. \end{aligned}$$

Collecting the four variables $\mathbf{z} = [x_1, \dot{x}_1, \theta_2, \dot{\theta}_2]^\top$ the equations of motion can be conveniently expressed as four coupled ordinary differential equations

$$\frac{d\mathbf{z}}{dt} = \begin{bmatrix} z_2 \\ \frac{2m_2lz_3^2 \sin z_4 + 3m_2g \sin z_4 \cos z_4 + 4u - 4bz_2}{4(m_1 + m_2) - 3m_2 \cos^2 z_4} \\ \frac{-3m_2lz_3^2 \sin z_4 \cos z_4 - 6(m_1 + m_2)g \sin z_4 - 6(u - bz_2) \cos z_4}{4l(m_1 + m_2) - 3m_2l \cos^2 z_4} \\ z_3 \end{bmatrix},$$

which can be simulated numerically.

5.5 Cart-Double Pendulum Swing-up

The cart-double pendulum dynamic system (see Figure 5.4) consists of a cart with mass m_1 and an attached double pendulum with masses m_2 and m_3 and lengths l_2 and l_3 for the two links, respectively. The double pendulum swings freely in the plane. The angles of the pendulum, θ_2 and θ_3 , are measured anti-clockwise from upright. The cart can move horizontally, with an applied external force u and the coefficient of friction b . Typical values are: $m_1 = 0.5$ kg, $m_2 = 0.5$ kg, $m_3 = 0.5$ kg, $l_2 = 0.6$ m, $l_3 = 0.6$ m, and $b = 0.1$ Ns/m. The coordinates, x_2 , y_2 and x_3 , y_3 of the midpoint of the pendulum elements are

$$\begin{aligned} \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} &= \begin{bmatrix} x_1 - \frac{1}{2}l_2 \sin \theta_2 \\ \frac{1}{2}l_2 \cos \theta_2 \end{bmatrix} \\ \begin{bmatrix} x_3 \\ y_3 \end{bmatrix} &= \begin{bmatrix} x_1 - l_2 \sin \theta_2 - \frac{1}{2}l_3 \sin \theta_3 \\ l_2 \cos \theta_2 + \frac{1}{2}l_3 \cos \theta_3 \end{bmatrix}. \end{aligned}$$

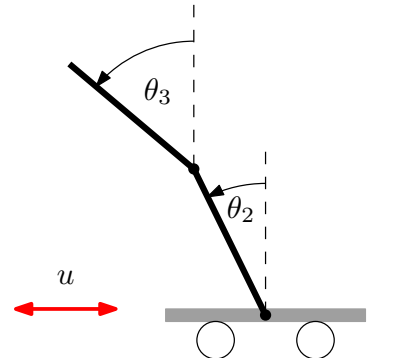


Figure 5.4: Cart-double pendulum.

The squared velocities of the cart and the pendulum midpoints are

$$\begin{aligned} v_1^2 &= \dot{x}_1^2, \\ v_2^2 &= \dot{x}_2^2 + \dot{y}_2^2 = \dot{x}_1^2 - l_2 \dot{x}_1 \dot{\theta}_2 \cos \theta_2 + \frac{1}{4} l_2^2 \dot{\theta}_2^2, \\ v_3^2 &= \dot{x}_3^2 + \dot{y}_3^2 = \dot{x}_1^2 + l_2^2 \dot{\theta}_2^2 + \frac{1}{4} l_3^2 \dot{\theta}_3^2 - 2l_2 \dot{x}_1 \dot{\theta}_2 \cos \theta_2 - l_3 \dot{x}_1 \dot{\theta}_3 \cos \theta_3 + l_2 l_3 \dot{\theta}_2 \dot{\theta}_3 \cos(\theta_2 - \theta_3). \end{aligned}$$

The system Lagrangian is the difference between the kinematic energy T and the potential energy V and given by

$$\begin{aligned} L = T - V &= \frac{1}{2} m_1 v_1^2 + \frac{1}{2} m_2 v_2^2 + \frac{1}{2} m_3 v_3^2 + \frac{1}{2} I_2 \dot{\theta}_2^2 + \frac{1}{2} I_3 \dot{\theta}_3^2 - m_2 g y_2 - m_3 g y_3 \\ &= \frac{1}{2} (m_1 + m_2 + m_3) \dot{x}_1^2 - \frac{1}{2} m_2 l_2 \dot{x}_1 \dot{\theta}_2 \cos(\theta_2) - \frac{1}{2} m_3 (2l_2 \dot{x}_1 \dot{\theta}_2 \cos(\theta_2) + l_3 \dot{x}_1 \dot{\theta}_3 \cos(\theta_3)) \\ &\quad + \frac{1}{8} m_2 l_2^2 \dot{\theta}_2^2 + \frac{1}{2} I_2 \dot{\theta}_2^2 + \frac{1}{2} m_3 (l_2^2 \dot{\theta}_2^2 + \frac{1}{4} l_3^2 \dot{\theta}_3^2 + l_2 l_3 \dot{\theta}_2 \dot{\theta}_3 \cos(\theta_2 - \theta_3)) + \frac{1}{2} I_3 \dot{\theta}_3^2 \\ &\quad - \frac{1}{2} m_2 g l_2 \cos(\theta_2) - m_3 g (l_2 \cos(\theta_2) + \frac{1}{2} l_3 \cos(\theta_3)). \end{aligned}$$

The angular moment of inertia I_j , $j = 2, 3$ around the pendulum midpoint is $I_j = \frac{1}{12} m l_j^2$, and $g = 9.82 \text{ m/s}^2$ is the acceleration of gravity. These moments inertia imply the assumption that the pendulums are thin (but rigid) wires.

The equations of motion are

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} - \frac{\partial L}{\partial q_i} = Q_i, \quad (5.8)$$

where Q_i are the non-conservative forces. We obtain the partial derivatives

$$\begin{aligned} \frac{\partial L}{\partial \dot{x}_1} &= (m_1 + m_2 + m_3) \dot{x}_1 - (\frac{1}{2} m_2 + m_3) l_2 \dot{\theta}_2 \cos \theta_2 - \frac{1}{2} m_3 l_3 \dot{\theta}_3 \cos \theta_3, \\ \frac{\partial L}{\partial x_1} &= 0, \\ \frac{\partial L}{\partial \dot{\theta}_2} &= (m_3 l_2^2 + \frac{1}{4} m_2 l_2^2 + I_2) \dot{\theta}_2 - (\frac{1}{2} m_2 + m_3) l_2 \dot{x}_1 \cos \theta_2 + \frac{1}{2} m_3 l_2 l_3 \dot{\theta}_3 \cos(\theta_2 - \theta_3), \\ \frac{\partial L}{\partial \theta_2} &= (\frac{1}{2} m_2 + m_3) l_2 (\dot{x}_1 \dot{\theta}_2 + g) \sin \theta_2 - \frac{1}{2} m_3 l_2 l_3 \dot{\theta}_2 \dot{\theta}_3 \sin(\theta_2 - \theta_3), \\ \frac{\partial L}{\partial \dot{\theta}_3} &= m_3 l_3 \left[-\frac{1}{2} \dot{x}_1 \cos \theta_3 + \frac{1}{2} l_2 \dot{\theta}_2 \cos(\theta_2 - \theta_3) + \frac{1}{4} l_3 \dot{\theta}_3 \right] + I_3 \dot{\theta}_3, \\ \frac{\partial L}{\partial \theta_3} &= \frac{1}{2} m_3 l_3 \left[(\dot{x}_1 \dot{\theta}_3 + g) \sin \theta_3 + l_2 \dot{\theta}_2 \dot{\theta}_3 \sin(\theta_2 - \theta_3) \right] \end{aligned}$$

leading to the equations of motion

$$\begin{aligned} (m_1 + m_2 + m_3) \ddot{x}_1 + \frac{1}{2} m_2 + m_3) l_2 (\dot{\theta}_2^2 \sin \theta_2 - \ddot{\theta}_2 \cos \theta_2) \\ + \frac{1}{2} m_3 l_3 (\dot{\theta}_3^2 \sin \theta_3 - \ddot{\theta}_3 \cos \theta_3) &= u - b \dot{x}_1 \\ (m_3 l_2^2 + I_2 + \frac{1}{4} m_2 l_2^2) \ddot{\theta}_2 - (\frac{1}{2} m_2 + m_3) l_2 (\ddot{x}_1 \cos \theta_2 + g \sin \theta_2) \\ + \frac{1}{2} m_3 l_2 l_3 [\ddot{\theta}_3 \cos(\theta_2 - \theta_3) + \dot{\theta}_3^2 \sin(\theta_2 - \theta_3)] &= 0 \\ (\frac{1}{4} m_2 l_3^2 + I_3) \ddot{\theta}_3 - \frac{1}{2} m_3 l_3 (\ddot{x}_1 \cos \theta_3 + g \sin \theta_3) \\ + \frac{1}{2} m_3 l_2 l_3 [\ddot{\theta}_2 \cos(\theta_2 - \theta_3) - \dot{\theta}_2^2 \sin(\theta_2 - \theta_3)] &= 0 \end{aligned}$$

These three linear equations in $(\ddot{x}_1, \ddot{\theta}_2, \ddot{\theta}_3)$ can be rewritten as the linear equation system

$$\begin{bmatrix} (m_1 + m_2 + m_3) & -\frac{1}{2} (m_2 + 2m_3) l_2 \cos \theta_2 & -\frac{1}{2} m_3 l_3 \cos \theta_3 \\ -(\frac{1}{2} m_2 + m_3) l_2 \cos \theta_2 & m_3 l_2^2 + I_2 + \frac{1}{4} m_2 l_2^2 & \frac{1}{2} m_3 l_2 l_3 \cos(\theta_2 - \theta_3) \\ -\frac{1}{2} m_3 l_3 \cos \theta_3 & \frac{1}{2} m_3 l_2 l_3 \cos(\theta_2 - \theta_3) & \frac{1}{4} m_2 l_3^2 + I_3 \end{bmatrix} \begin{bmatrix} \ddot{x}_1 \\ \ddot{\theta}_2 \\ \ddot{\theta}_3 \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix},$$

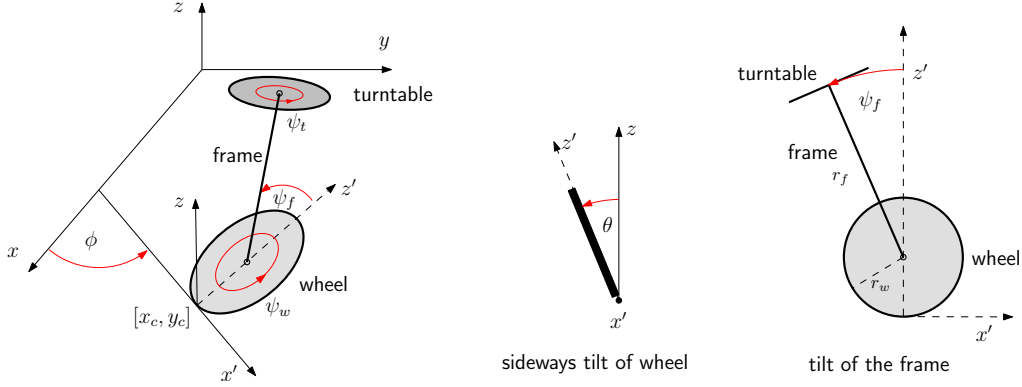


Figure 5.5: Unicycle.

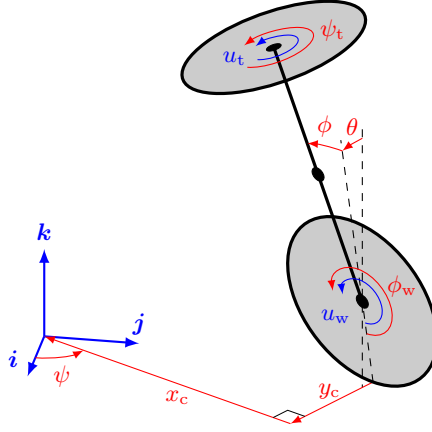


Figure 5.6: The robotic unicycle with state variables: pitch angle ϕ , roll angle θ , yaw angle ψ , wheel angle ϕ_w , turntable angle ψ_t , the associated angular velocities and the location of the global origin (x_c, y_c) . The actions are the wheel motor torque u_w and the turntable motor torque u_t . The global coordinate system is defined by the vectors i, j and k .

where

$$\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} u - b\dot{x}_1 - \frac{1}{2}(m_2 + 2m_3)l_2\dot{\theta}_2^2 \sin \theta_2 - \frac{1}{2}m_3l_3\dot{\theta}_3^2 \sin \theta_3 \\ (\frac{1}{2}m_2 + m_3)l_2g \sin \theta_2 - \frac{1}{2}m_3l_2l_3\dot{\theta}_3^2 \sin(\theta_2 - \theta_3) \\ \frac{1}{2}m_3l_3[g \sin \theta_3 + l_2\dot{\theta}_2^2 \sin(\theta_2 - \theta_3)] \end{bmatrix}.$$

This linear equation system can be solved for $\ddot{x}_1, \ddot{\theta}_2, \ddot{\theta}_3$ and used for numerical simulation.

5.6 Unicycling

A sketch of the unicycle system is shown in Figure 5.5. The equations of motion that govern the unicycle were derived in the MSc thesis by Forster [5]. We shall provide a sketch of the full derivation here, in which we follow the steps taken by Forster [5], Section 3.3.

Constant	Description	Value	Units
m_w	Wheel mass	1.0	kg
r_w	Wheel radius	0.225	m
A_w	Moment of inertia of wheel around \mathbf{i}_w	0.0242	kg m ²
C_w	Moment of inertia of wheel around \mathbf{k}_w	0.0484	kg m ²
m_f	Frame mass	23.5	kg
r_f	Frame centre of mass to wheel	0.54	m
A_f	Moment of inertia of frame around \mathbf{i}_f	0.4248	kg m ²
B_f	Moment of inertia of frame around \mathbf{j}_f	0.4608	kg m ²
C_f	Moment of inertia of frame around \mathbf{k}_f	0.8292	kg m ²
m_t	Turntable mass	10.0	kg
r_t	Frame centre of mass to turntable	0.27	m
A_t	Moment of inertia of turntable around \mathbf{i}_t	1.3	kg m ²
C_t	Moment of inertia of turntable around \mathbf{k}_t	0.2	kg m ²
g	Gravitational acceleration	9.81	m s ⁻²

Table 5.2: Physical constants used for the simulated robotic unicycle. The coordinate systems defined by the \mathbf{i}, \mathbf{j} and \mathbf{k} vectors are shown in Figure 5.7.

5.6.1 Method

The robotic unicycle is shown in Figure 5.6 with global coordinate system defined by the orthonormal vectors \mathbf{i}, \mathbf{j} and \mathbf{k} . The spatial position of the unicycle is fully defined by the pitch angle ϕ , roll angle θ , yaw angle ψ , wheel angle ϕ_w , turntable angle ψ_t and location of the global origin with respect to the body-centred coordinate system (x_c, y_c) . We chose the state vector to be $\mathbf{x} = [\phi, \dot{\phi}, \theta, \dot{\theta}, \psi, \dot{\psi}, \dot{\phi}_w, \dot{\psi}_t, x_c, y_c]^\top \in \mathbb{R}^{10}$ where we exclude ϕ_w and ψ_t since they clearly have no effect on the dynamics. The action vector \mathbf{u} is made up of a wheel motor torque u_w and a turntable motor torque u_t .

Let us start with the coordinates (x_c, y_c) . These are centred on the point of contact with the floor and define the location of the global origin. The coordinate x_c lies parallel to the current direction of travel and y_c is orthogonal to it. These coordinates evolve according to

$$\dot{x}_c = r_w \dot{\phi}_w \cos \psi \quad (5.9)$$

$$\dot{y}_c = r_w \dot{\phi}_w \sin \psi \quad (5.10)$$

where r_w is the wheel radius. The full unicycle model was obtained by analysing the wheel, frame and turntable as individual Free Body Diagrams (FBDs), as depicted in Figure 5.7. Linear momentum and moment of momentum for each FBD were then resolved to yield six scalar equations for each free body. The internal forces were then eliminated to yield five independent scalar equations which govern the evolution of the angular states. A description of the physical constants of the system along with the values we use in this thesis are given in Table 5.2.

5.6.2 Wheel FBD

The wheel coordinate system is defined by the orthonormal vectors $\mathbf{i}_w, \mathbf{j}_w$ and \mathbf{k}_w as shown in Figure 5.7(a). We begin by noting that the angular velocity of the wheel coordinate system is $\Omega_w = \dot{\psi}\mathbf{k} + \dot{\theta}\mathbf{j}_w$. Now noting that the angular velocity of the wheel only differs in the \mathbf{k}_w direction and assuming no slip between wheel and floor we have expressions for the velocity and angular velocity of the wheel

$$\begin{aligned}\mathbf{v}_w &= -(\boldsymbol{\omega}_w \times r_w \mathbf{i}_w) \\ \boldsymbol{\omega}_w &= \Omega_w + \dot{\phi}_w \mathbf{k}_w\end{aligned}$$

From these expressions we can derive the acceleration of the wheel $\dot{\mathbf{v}}_w$ and the rate of change of angular momentum $\dot{\mathbf{h}}_w$ as

$$\begin{aligned}\dot{\mathbf{v}}_w &= \frac{\partial \mathbf{v}_w}{\partial t} + (\Omega_w \times \mathbf{v}_w) \\ \dot{\mathbf{h}}_w &= A_w \frac{\partial \omega_w[1]}{\partial t} \mathbf{i}_w + A_w \frac{\partial \omega_w[2]}{\partial t} \mathbf{j}_w + C_w \frac{\partial \omega_w[3]}{\partial t} \mathbf{k}_w + (\Omega_w \times \mathbf{h}_w)\end{aligned}$$

where angular momentum in the wheel frame of reference is $\mathbf{h}_w = [A_w; A_w; C_w] \circ \boldsymbol{\omega}_w$. Now we consider the forces acting on the wheel free body. These are given by the unknown quantities: axle force F_w , axle torque Q_w & reaction force R and the known quantities: wheel weight W_w & friction torque T . These forces and moments are shown in the right-hand plot of Figure 5.7(a). Note that we actually know the component of the axle torque Q_w in the \mathbf{k}_w direction as it is given by the reaction of the wheel motor on the wheel itself u_w . Resolving the rate of change of linear momentum and the rate of change of angular momentum around the center of mass leads to

$$m_w \dot{\mathbf{v}}_w = R + F_w + W_w \quad (5.11)$$

$$\dot{\mathbf{h}}_w = (r_w \mathbf{i}_w \times R) + Q_w + T \quad (5.12)$$

5.6.3 Frame FBD

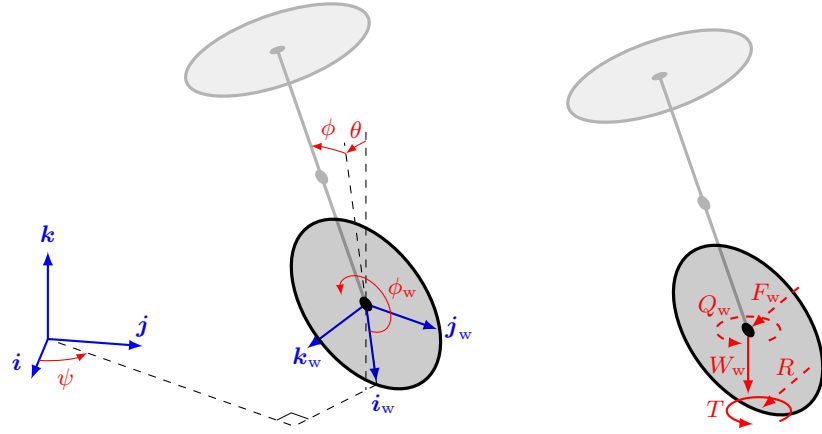
The frame coordinate system is defined by the orthonormal vectors $\mathbf{i}_f, \mathbf{j}_f$ and $\mathbf{k}_f = \mathbf{k}_w$ as shown in Figure 5.7(b). In this case, the angular velocity of the frame $\boldsymbol{\omega}_f$ is given by the angular velocity of the wheel plus an additional spin about the wheel axis and the velocity of the frame \mathbf{v}_f is given by the velocity of the wheel plus an additional rotation about the wheel centre

$$\begin{aligned}\mathbf{v}_f &= \mathbf{v}_w - (\boldsymbol{\omega}_f \times r_f \mathbf{i}_f) \\ \boldsymbol{\omega}_f &= \Omega_w + \dot{\phi} \mathbf{k}_f\end{aligned}$$

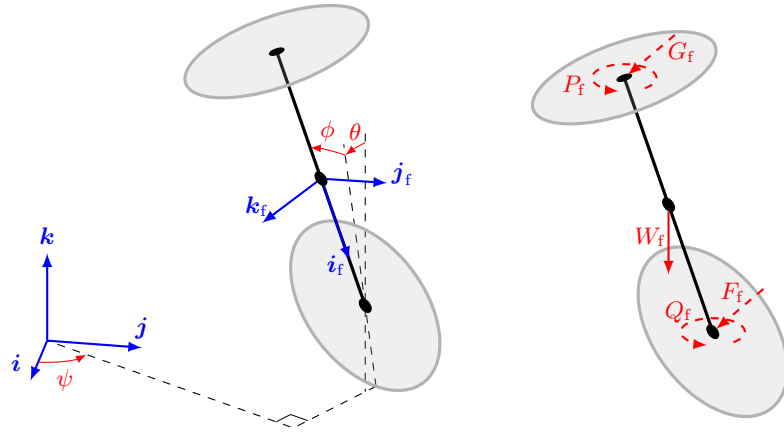
As before, we can now derive the acceleration of the frame $\dot{\mathbf{v}}_f$ and the rate of change of angular momentum $\dot{\mathbf{h}}_f$ as

$$\begin{aligned}\dot{\mathbf{v}}_f &= \frac{\partial \mathbf{v}_f}{\partial t} + (\Omega_f \times \mathbf{v}_f) \\ \dot{\mathbf{h}}_f &= A_f \frac{\partial \omega_f[1]}{\partial t} \mathbf{i}_f + B_f \frac{\partial \omega_f[2]}{\partial t} \mathbf{j}_f + C_f \frac{\partial \omega_f[3]}{\partial t} \mathbf{k}_f + (\Omega_f \times \mathbf{h}_f)\end{aligned}$$

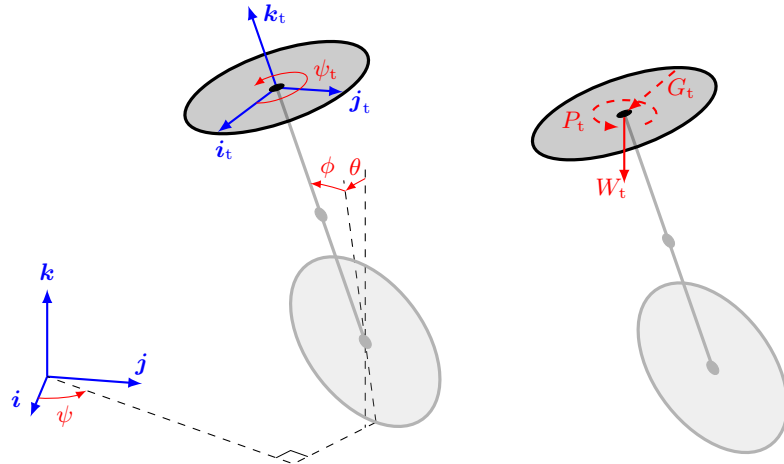
where angular momentum of the frame in this frame of reference is $\mathbf{h}_f = [A_f; B_f; C_f] \circ \boldsymbol{\omega}_f$. The forces and moments acting on the frame are shown on the right in Figure 5.7(b). They consist of the known



(a) The wheel as a free body diagram and coordinate system defined by i_w, j_w and k_w



(b) The frame as a free body diagram and coordinate system defined by i_f, j_f and k_f



(c) The turntable as a free body diagram and coordinate system defined by i_t, j_t and k_t

Figure 5.7: Free body diagrams of the wheel, frame and turntable of the unicycle. The model has the angular state variables pitch ϕ , roll θ , yaw ψ , wheel angle ϕ_w and turntable angle ψ_t . The vectors i, j and k define the global fixed frame of reference. The centres of mass for each component are shown by the black dots. Forces and moments are shown on the right, with unknown quantities as dashed lines.

frame weight W_f and the unknown: wheel axle force $F_f = -F_w$, wheel axle torque $Q_f = -Q_w$, turntable axle force G_f & turntable axle torque P_f . But again we know that the dimension of P_f acting along \mathbf{i}_f is given by the reaction of the frame to the turntable motor u_t . So resolving the rate of change of linear momentum and the rate of change of angular momentum around the centre of mass gives us

$$m_f \dot{\mathbf{v}}_f = F_f + G_f + W_f \quad (5.13)$$

$$\dot{\mathbf{h}}_f = (r_f \mathbf{i}_f \times F_f) - (r_t \mathbf{i}_f \times G_f) + Q_f + P_f \quad (5.14)$$

5.6.4 Turntable FBD

Finally, the turntable coordinate system is defined by the orthonormal vectors $\mathbf{i}_t = \mathbf{k}_f, \mathbf{j}_t = \mathbf{j}_f$ and $\mathbf{k}_t = -\mathbf{i}_f$ as shown in Figure 5.7(c). The velocity of the turntable centre \mathbf{v}_t is equal to the velocity of the wheel plus an additional term caused by rotation about the wheel centre, while the angular velocity $\boldsymbol{\omega}_t$ differs from $\Omega_t = \Omega_f$ only along \mathbf{k}_t

$$\mathbf{v}_t = \mathbf{v}_w + (\Omega_t \times r \mathbf{k}_t)$$

$$\boldsymbol{\omega}_t = \Omega_t + \dot{\psi}_t \mathbf{k}_t$$

Again, we derive the acceleration of the frame $\dot{\mathbf{v}}_t$ and the rate of change of angular momentum $\dot{\mathbf{h}}_t$ as

$$\begin{aligned} \dot{\mathbf{v}}_t &= \frac{\partial \mathbf{v}_f}{\partial t} + (\Omega_f \times \mathbf{v}_f) \\ \dot{\mathbf{h}}_t &= A_t \frac{\partial \omega_t[1]}{\partial t} \mathbf{i}_t + A_t \frac{\partial \omega_t[2]}{\partial t} \mathbf{j}_t + C_t \frac{\partial \omega_t[3]}{\partial t} \mathbf{k}_t + (\Omega_t \times \mathbf{h}_t) \end{aligned}$$

where $\mathbf{h}_t = [A_t; A_t; C_t] \circ \boldsymbol{\omega}_t$. The forces and moments acting on the turntable lead to the last of our equations

$$m_t \dot{\mathbf{v}}_t = G_t + W_t \quad (5.15)$$

$$\dot{\mathbf{h}}_t = P_t \quad (5.16)$$

5.6.5 Eliminating Internal Forces

We now have 18 kinematic relationships given by Equations (5.11)–(5.16) which govern the dynamics of the unicycle. These can be reduced to five expressions by eliminating the 13 scalar unknowns found in the unknown internal forces, F and G , the unknown reaction force R and the partially unknown torques Q and P . The first can be obtained from Equation (5.16) and noting that the component of P_f about \mathbf{k}_t is the reaction to the motor torque u_t

$$\dot{\mathbf{h}}_t^\top \mathbf{k}_t = u_t \quad (5.17)$$

The second can be obtained by first making use of the relationships $G_f = -G_t$ & $P_f = -P_t$ and then rearranging Equation (5.13), Equation (5.15) and Equation (5.16) to get

$$F_f = m_t \dot{\mathbf{v}}_t + m_f \dot{\mathbf{v}}_f - W_t - W_f$$

$$G_f = W_t - m_t \dot{\mathbf{v}}_t$$

$$P_f = -\dot{\mathbf{h}}_t$$

Substituting these into Equation (5.14) and noting that $Q_w = -Q_f$ gives us

$$Q_w = -\dot{\mathbf{h}}_f - \dot{\mathbf{h}}_t - \left(r_f \mathbf{i}_f \times (W_f + W_t - m_f \dot{\mathbf{v}}_f - m_t \dot{\mathbf{v}}_t) \right) - \left(r_t \mathbf{i}_f \times (W_t - m_t \dot{\mathbf{v}}_t) \right)$$

Once again, the component of Q_w about the wheel axis is equal to the torque provided by the wheel motor u_w , therefore $Q_w^\top \mathbf{k}_w = u_w$ and we have our second expression

$$-\left(\dot{\mathbf{h}}_f + \dot{\mathbf{h}}_t + \left(r_f \mathbf{i}_f \times (W_f + W_t - m_f \dot{\mathbf{v}}_f - m_t \dot{\mathbf{v}}_t) \right) + \left(r_t \mathbf{i}_f \times (W_t - m_t \dot{\mathbf{v}}_t) \right) \right)^\top \mathbf{k}_w = u_w \quad (5.18)$$

Finally, Equation (5.11) can be combined with our expression for $F_f = -F_w$ to find the reaction force at the base

$$\mathbf{R} = m_w \dot{\mathbf{v}}_w + m_f \dot{\mathbf{v}}_f + m_t \dot{\mathbf{v}}_t - W_w - W_f - W_t$$

and can be substituted into Equation (5.12) to obtain the following three relationships

$$\begin{aligned} \dot{\mathbf{h}}_w = & \left(r_w \mathbf{i}_w \times (m_w \dot{\mathbf{v}}_w + m_f \dot{\mathbf{v}}_f + m_t \dot{\mathbf{v}}_t - W_w - W_f - W_t) \right) \\ & - \left(\dot{\mathbf{h}}_f + \dot{\mathbf{h}}_t + \left(r_f \mathbf{i}_f \times (W_f + W_t - m_f \dot{\mathbf{v}}_f - m_t \dot{\mathbf{v}}_t) \right) + \left(r_t \mathbf{i}_f \times (W_t - m_t \dot{\mathbf{v}}_t) \right) \right) + T \end{aligned} \quad (5.19)$$

The five expressions contained in Equations (5.17)–(5.19) form the foundation for the model we use. The only unknowns remaining in these relationships are the states: $\phi, \dot{\phi}, \theta, \dot{\theta}, \psi, \dot{\psi}, \dot{\phi}_w, \dot{\psi}_t$, the torque actions: u_w, u_t and the accelerations: $\ddot{\phi}, \ddot{\theta}, \ddot{\psi}, \ddot{\phi}_w, \ddot{\psi}_t$. The equations were then rearranged by Forster, through use of the Symbolic Toolbox in MATLAB, to isolate the five acceleration terms. We use this mapping along with Equations (5.9)–(5.10) to perform simulations of the unicycle.

Chapter 6

Testing and Debugging

A crucial ingredient that is required for PILCO's success is the computation of some gradients. Every important function needs to compute some sort of function value and the derivative of this function value with respect to some of the input values.

We have included some relatively generic¹ gradient test functions in `<pilco_root>/test/`. In all functions, the analytic gradients are compared with finite-difference approximations (see `<pilco_root>/test/checkgrad.m`).

Usually, the test functions can be called with the variables in the MATLAB workspace during the execution of PILCO.

6.1 Gradient Checks for the Controller Function

`conT` checks gradients of the policy (controller) functions.

6.1.1 Interface

```
1 function [dd dy dh] = conT(deriv, policy, m, s, delta)
```

Input arguments:

`deriv` desired derivative. options:

- (i) 'dMdm' - derivative of the mean of the predicted control
 wrt the mean of the input distribution
- (ii) 'dMds' - derivative of the mean of the predicted control
 wrt the variance of the input distribution
- (iii) 'dMdp' - derivative of the mean of the predicted control
 wrt the controller parameters
- (iv) 'dSdm' - derivative of the variance of the predicted control
 wrt the mean of the input distribution
- (v) 'dSds' - derivative of the variance of the predicted control
 wrt the variance of the input distribution
- (vi) 'dSdp' - derivative of the variance of the predicted control
 wrt the controller parameters

¹The interfaces of the most important functions, e.g., controllers, cost functions, GP predictions, are unified.

- (vii) 'dCdm' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the predicted control})$ wrt the mean of the input distribution
- (viii) 'dCds' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the predicted control})$ wrt the variance of the input distribution
- (ix) 'dCdP' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the predicted control})$ wrt the controller parameters

policy policy structure
 .fcn function handle to policy
 .<> other fields that are passed on to the policy
 m mean of the input distribution
 s covariance of the input distribution
 delta (optional) finite difference parameter. Default: 1e-4

Output arguments:

dd relative error of analytical vs. finite difference gradient
 dy analytical gradient
 dh finite difference gradient

6.2 Gradient Checks for the Cost Function

lossT checks gradients of the (immediate) cost functions, which are specific to each scenario.

6.2.1 Interface

```
1 function [dd dy dh] = lossT(deriv, policy, m, s, delta)
```

Input arguments:

deriv desired derivative. options:

- (i) 'dMdm' - derivative of the mean of the predicted cost wrt the mean of the input distribution
- (ii) 'dMds' - derivative of the mean of the predicted cost wrt the variance of the input distribution
- (iii) 'dSdm' - derivative of the variance of the predicted cost wrt the mean of the input distribution
- (iv) 'dSds' - derivative of the variance of the predicted cost wrt the variance of the input distribution
- (v) 'dCdm' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the predicted cost})$ wrt the mean of the input distribution
- (vi) 'dCds' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the predicted cost})$ wrt the variance of the input distribution

cost cost structure
 .fcn function handle to cost
 .<> other fields that are passed on to the cost
 m mean of the input distribution
 s covariance of the input distribution
 delta (optional) finite difference parameter. Default: 1e-4

Output arguments:

dd relative error of analytical vs. finite difference gradient
dy analytical gradient
dh finite difference gradient

6.3 Gradient Checks for the GP Prediction Function

gpT checks gradients of the functions that implement the GP prediction at an uncertain test input, i.e., all gp* functions. These functions can be found in <pilco.root>/gp/.

6.3.1 Interface

```
1 function [dd dy dh] = gpT(deriv, gp, m, s, delta)
```

Input arguments:

deriv desired derivative. options:

- (i) 'dMdm' - derivative of the mean of the GP prediction
 wrt the mean of the input distribution
- (ii) 'dMds' - derivative of the mean of the GP prediction
 wrt the variance of the input distribution
- (iii) 'dMdp' - derivative of the mean of the GP prediction
 wrt the GP parameters
- (iv) 'dSdm' - derivative of the variance of the GP prediction
 wrt the mean of the input distribution
- (v) 'dSds' - derivative of the variance of the GP prediction
 wrt the variance of the input distribution
- (vi) 'dSdp' - derivative of the variance of the GP prediction
 wrt the GP parameters
- (vii) 'dVdm' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the GP prediction})$ wrt the mean of the input distribution
- (viii) 'dVds' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the GP prediction})$ wrt the variance of the input distribution
- (ix) 'dVdp' - derivative of $\text{inv}(s) \cdot (\text{covariance of the input and the GP prediction})$ wrt the GP parameters

gp GP structure

.fcn function handle to the GP function used for predictions at uncertain inputs

.<> other fields that are passed on to the GP function

m mean of the input distribution

s covariance of the input distribution

delta (optional) finite difference parameter. Default: 1e-4

Output arguments:

dd	relative error of analytical vs. finite difference gradient
dy	analytical gradient
dh	finite difference gradient

6.4 Gradient Checks for the State Propagation Function

`propagateT` checks gradients of the function `propagated`, which implements state propagation $p(\mathbf{x}_t) \rightarrow p(\mathbf{x}_{t+1})$. This function can be found in `<pilco_root>/base/`.

6.4.1 Interface

```
1 [dd dy dh] = propagateT(deriv, plant, dynmodel, policy, m, s, delta)
```

Input arguments:

m	mean of the state distribution at time t	[D x 1]
s	covariance of the state distribution at time t	[D x D]
plant	plant structure	
dynmodel	dynamics model structure	
policy	policy structure	

Output arguments:

Mnext	predicted mean at time t+1	[E x 1]
Snext	predicted covariance at time t+1	[E x E]
dMdm	output mean wrt input mean	[E x D]
dMds	output mean wrt input covariance matrix	[E x D*D]
dSdm	output covariance matrix wrt input mean	[E*E x D]
dSds	output cov wrt input cov	[E*E x D*D]
dMdp	output mean wrt policy parameters	[E x P]
dSdp	output covariance matrix wrt policy parameters	[E*E x P]

where P is the number of policy parameters.

6.5 Gradient Checks for Policy Evaluation

`valueT` checks the overall gradients $\partial J(\boldsymbol{\theta})/\partial \boldsymbol{\theta}$ of the expected long-term cost J with respect all policy parameters $\boldsymbol{\theta}$. Policy evaluation, i.e., computing J , and the corresponding gradients are implemented in `value.m` to be found in `<pilco_root>/base/`.

6.5.1 Interface

```
1 [d dy dh] = valueT(p, delta, m, s, dynmodel, policy, plant, cost, H)
```

Input arguments:

p	policy parameters (can be a structure)
.<>	fields that contain the policy parameters (nothing else)
m	mean of the input distribution
s	covariance of the input distribution
dynmodel	GP dynamics model (structure)
policy	policy structure
plant	plant structure
cost	cost structure
H	prediction horizon
delta	(optional) finite difference parameter. Default: 1e-4

Output arguments:

dd	relative error of analytical vs. finite difference gradient
dy	analytical gradient
dh	finite difference gradient

Chapter 7

Code and Auto-generated Documentation of the Main Functions

In the following, we include the auto-generated documentation of the most important functions of the PILCO learning framework. If you change the documentation in the *.m files, please run `<pilco_root>/doc/generate_docs.m` to update the rest of this chapter. HTML files will be generated as well and can be found in `<pilco_root>/doc/html/`.

7.1 Base Directory

7.1.1 `applyController.m`

Summary: Script to apply the learned controller to a (simulated) system

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2013-06-04

High-Level Steps

1. Generate a single trajectory rollout by applying the controller
2. Generate many rollouts for testing the performance of the controller
3. Save the data

Code

```
1 % 1. Generate trajectory rollout given the current policy
2 if isfield(plant, 'constraint'), HH = maxH; else HH = H; end
3 [xx, yy, realCost{j+J}, latent{j}] = ...
4   rollout(gaussian(mu0, S0), policy, HH, plant, cost);
5 disp(xx); % display states of observed trajectory
6 x = [x; xx]; y = [y; yy]; % augment training set
7 if plotting.verbosity > 0
8   if ~ishandle(3); figure(3); else set(0, 'CurrentFigure', 3); end
```

```

9   hold on; plot(1:length(realCost{J+j}),realCost{J+j}, 'r'); drawnow;
10  end
11
12  % 2. Make many rollouts to test the controller quality
13  if plotting.verbosity > 1
14      lat = cell(1,10);
15      for i=1:10
16          [~,~,~,lat{i}] = rollout(gaussian(mu0, S0), policy, HH, plant, cost);
17      end
18
19      if ~ishandle(4); figure(4); else set(0,'CurrentFigure',4); end; clf(4);
20
21      ldyno = length(dyno);
22      for i=1:ldyno % plot the rollouts on top of predicted error bars
23          subplot(ceil(ldyno/sqrt(ldyno)),ceil(sqrt(ldyno)),i); hold on;
24          errorbar( 0:length(M{j}(i,:))-1, M{j}(i,:), ...
25                  2*sqrt(squeeze(Sigma{j}(i,i,:))) );
26          for ii=1:10
27              plot( 0:size(lat{ii}(:,dyno(i)),1)-1, lat{ii}(:,dyno(i)), 'r' );
28          end
29          plot( 0:size(latent{j}(:,dyno(i)),1)-1, latent{j}(:,dyno(i)), 'g');
30          axis tight
31      end
32      drawnow;
33  end
34
35  % 3. Save data
36  filename = [basename num2str(j) '_H' num2str(H)]; save(filename);

```

7.1.2 propagate.m

Summary: Propagate the state distribution one time step forward.

`[Mnext, Snext] = propagate(m, s, plant, dynmodel, policy)`

Input arguments:

<code>m</code>	mean of the state distribution at time <code>t</code>	<code>[D x 1]</code>
<code>s</code>	covariance of the state distribution at time <code>t</code>	<code>[D x D]</code>
<code>plant</code>	plant structure	
<code>dynmodel</code>	dynamics model structure	
<code>policy</code>	policy structure	

Output arguments:

<code>Mnext</code>	mean of the successor state at time <code>t+1</code>	<code>[E x 1]</code>
<code>Snext</code>	covariance of the successor state at time <code>t+1</code>	<code>[E x E]</code>

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, Henrik Ohlsson, and Carl Edward Rasmussen.

Last modified: 2013-01-23

High-Level Steps

1. Augment state distribution with trigonometric functions
2. Compute distribution of the control signal
3. Compute dynamics-GP prediction
4. Compute distribution of the next state

```
1 function [Mnext, Snext] = propagate(m, s, plant, dynmodel, policy)
```

Code

```
1 % extract important indices from structures
2 ang_i = plant.ang_i; % angular indices
3 poli = plant.poli; % policy indices
4 dyn_i = plant.dyn_i; % dynamics-model indices
5 dif_i = plant.dif_i; % state indices where the model was trained on differences
6
7 D0 = length(m); % size of the input mean
8 D1 = D0 + 2*length(ang_i); % length after mapping all angles to sin/cos
9 D2 = D1 + length(policy.maxU); % length after computing control signal
10 D3 = D2 + D0; % length after predicting
11 M = zeros(D3,1); M(1:D0) = m; S = zeros(D3); S(1:D0,1:D0) = s; % init M and S
12
13 % 1) Augment state distribution with trigonometric functions -----
14 i = 1:D0; j = 1:D0; k = D0+1:D1;
15 [M(k), S(k,k), C] = gTrig(M(i), S(i,i), ang_i);
16 q = S(j,i)*C; S(j,k) = q; S(k,j) = q';
17
18 mm=zeros(D1,1); mm(i)=M(i); ss(i,i)=S(i,i)+diag(exp(2*dynmodel.hyp(end,:))/2);
19 [mm(k), ss(k,k), C] = gTrig(mm(i), ss(i,i), ang_i); % noisy state measurement
20 q = ss(j,i)*C; ss(j,k) = q; ss(k,j) = q';
21
22 % 2) Compute distribution of the control signal -----
23 i = poli; j = 1:D1; k = D1+1:D2;
24 [M(k), S(k,k), C] = policy.fcn(policy, mm(i), ss(i,i));
25 q = S(j,i)*C; S(j,k) = q; S(k,j) = q';
26
27 % 3) Compute dynamics-GP prediction -----
28 ii = [dyn_i D1+1:D2]; j = 1:D2;
29 if isfield(dynmodel, 'sub'), Nf = length(dynmodel.sub); else Nf = 1; end
30 for n=1:Nf % potentially multiple dynamics models
31 [dyn_i k] = sliceModel(dynmodel, n, ii, D1, D2, D3); j = setdiff(j, k);
32 [M(k), S(k,k), C] = dyn.fcn(dyn, M(i), S(i,i));
33 q = S(j,i)*C; S(j,k) = q; S(k,j) = q';
34
35 j = [j k]; % update 'previous' state vector
36 end
37
38 % 4) Compute distribution of the next state -----
39 P = [zeros(D0, D2) eye(D0)]; P(dif_i, dif_i) = eye(length(dif_i));
40 Mnext = P*M; Snext = P*S*P'; Snext = (Snext+Snext')/2;
```

```

1 function [dyn i k] = sliceModel(dynmodel,n,ii,D1,D2,D3) % separate sub-dynamics
2 % A1) Separate multiple dynamics models
3 if isfield(dynmodel,'sub')
4     dyn = dynmodel.sub{n}; do = dyn.dyno; D = length(ii)+D1-D2;
5     if isfield(dyn,'dyni'), di=dyn.dyni; else di=[]; end
6     if isfield(dyn,'dynu'), du=dyn.dynu; else du=[]; end
7     if isfield(dyn,'dynj'), dj=dyn.dynj; else dj=[]; end
8     i = [ii(di) D1+du D2+dj]; k = D2+do;
9     dyn.inputs = [dynmodel.inputs(:,[di D+du]) dynmodel.target(:,dj)]; % inputs
10    dyn.target = dynmodel.target(:,do); % targets
11 else
12     dyn = dynmodel; k = D2+1:D3; i = ii;
13 end

```

7.1.3 rollout.m

Summary: Generate a state trajectory using an ODE solver (and any additional dynamics) from a particular initial state by applying either a particular policy or random actions.

```
function [x y L latent] = rollout(start, policy, H, plant, cost)
```

Input arguments:

start	vector containing initial states (without controls)	[nX x 1]
policy	policy structure	
.fcn	policy function	
.p	parameter structure (if empty: use random actions)	
.maxU	vector of control input saturation values	[nU x 1]
H	rollout horizon in steps	
plant	the dynamical system structure	
.subplant	(opt) additional discrete-time dynamics	
.augment	(opt) augment state using a known mapping	
.constraint	(opt) stop rollout if violated	
.poli	indices for states passed to the policy	
.dyno	indices for states passed to cost	
.odei	indices for states passed to the ode solver	
.subi	(opt) indices for states passed to subplant function	
.augi	(opt) indices for states passed to augment function	
cost	cost structure	

Output arguments:

x	matrix of observed states	[H x nX+nU]
y	matrix of corresponding observed successor states	[H x nX]
L	cost incurred at each time step	[1 x H]
latent	matrix of latent states	[H+1 x nX]

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modification: 2013-05-21

High-Level Steps

1. Compute control signal u from state x : either apply policy or random actions
2. Simulate the true dynamics for one time step using the current pair (x, u)
3. Check whether any constraints are violated (stop if true)
4. Apply random noise to the successor state
5. Compute cost (optional)
6. Repeat until end of horizon

```
1 function [x y L latent] = rollout(start, policy, H, plant, cost)
```

Code

```
1 if isfield(plant, 'augment'), augi = plant.augi; % sort out indices!
2 else plant.augment = inline('[] '); augi = []; end
3 if isfield(plant, 'subplant'), subi = plant.subi;
4 else plant.subplant = inline('[] ',1); subi = []; end
5 odei = plant.odei; poli = plant.poli; dyno = plant.dyno; angi = plant.angi;
6 simi = sort([odei subi]);
7 nX = length(simi)+length(augi); nU = length(policy.maxU); nA = length(angi);
8
9 state(simi) = start; state(augi) = plant.augment(state); % initializations
10 x = zeros(H+1, nX+2*nA);
11 x(1,simi) = start' + randn(size(simi))*chol(plant.noise);
12 x(1,augi) = plant.augment(x(1,:));
13 u = zeros(H, nU); latent = zeros(H+1, size(state,2)+nU);
14 y = zeros(H, nX); L = zeros(1, H); next = zeros(1, length(simi));
15
16 for i = 1:H % ----- generate trajectory
17     s = x(i,dyno)'; sa = gTrig(s, zeros(length(s)), angi); s = [s; sa];
18     x(i,end-2*nA+1:end) = s(end-2*nA+1:end);
19
20     % 1. Apply policy ... or random actions -----
21     if isfield(policy, 'fcn')
22         u(i,:) = policy.fcn(policy, s(poli), zeros(length(poli)));
23     else
24         u(i,:) = policy.maxU.*(2*rand(1,nU)-1);
25     end
26     latent(i,:) = [state u(i,:)]; % latent state
27
28     % 2. Simulate dynamics -----
29     next(odei) = simulate(state(odei), u(i,:), plant);
30     next(subi) = plant.subplant(state, u(i,:));
31
32     % 3. Stop rollout if constraints violated -----
33     if isfield(plant, 'constraint') && plant.constraint(next(odei))
34         H = i-1;
35         fprintf('state constraints violated...\n');
36         break;
```



```

37 end
38
39 % 4. Augment state and randomize
40 state(simi) = next(simi); state(augi) = plant.augment(state);
41 x(i+1,simi) = state(simi) + randn(size(simi))*chol(plant.noise);
42 x(i+1,augi) = plant.augment(x(i+1,:));
43
44 % 5. Compute Cost
45 if nargout > 2
46     L(i) = cost.fcn(cost, state(dyno)', zeros(length(dyno)));
47 end
48 end
49
50 y = x(2:H+1,1:nX); x = [x(1:H,:) u(1:H,:)]';
51 latent(H+1, 1:nX) = state; latent = latent(1:H+1,:); L = L(1,1:H);

```

7.1.4 trainDynModel.m

Summary: Script to learn the dynamics model

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modification: 2013-05-20

High-Level Steps

1. Extract states and controls from x-matrix
2. Define the training inputs and targets of the GP
3. Train the GP

Code

```

1 % 1. Train GP dynamics model
2 Du = length(policy.maxU); Da = length(plant.angi); % no. of ctrl and angles
3 xaug = [x(:,dyno) x(:,end-Du-2*Da+1:end-Du)]; % x augmented with angles
4 dynmodel.inputs = [xaug(:,dyni) x(:,end-Du+1:end)]; % use dyni and ctrl
5 dynmodel.targets = y(:,dyno);
6 dynmodel.targets(:,difi) = dynmodel.targets(:,difi) - x(:,dyno(difi));
7
8 dynmodel = dynmodel.train(dynmodel, plant, trainOpt); % train dynamics GP
9
10 % display some hyperparameters
11 Xh = dynmodel.hyp;
12 % noise standard deviations
13 disp(['Learned noise std: ' num2str(exp(Xh(end,:)))]);
14 % signal-to-noise ratios (values > 500 can cause numerical problems)
15 disp(['SNRs : ' num2str(exp(Xh(end-1,:)-Xh(end,:)))]);

```

7.1.5 value.m

Summary: Compute expected (discounted) cumulative cost for a given (set of) initial state distributions

```
function [J, dJdp] = value(p, m0, S0, dynmodel, policy, plant, cost, H)
```

Input arguments:

p	policy parameters chosen by minimize
policy	policy structure
.fcn	function which implements the policy
.p	parameters passed to the policy
m0	matrix (D by k) of initial state means
S0	covariance matrix (D by D) for initial state
dynmodel	dynamics model structure
plant	plant structure
cost	cost function structure
.fcn	function handle to the cost
.gamma	discount factor
H	length of prediction horizon

Output arguments:

J	expected cumulative (discounted) cost
dJdp	(optional) derivative of J wrt the policy parameters

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modification: 2013-03-21

High-Level Steps

1. Compute distribution of next state
2. Compute corresponding expected immediate cost (discounted)
3. At end of prediction horizon: sum all immediate costs up

```
1 function [J, dJdp] = value(p, m0, S0, dynmodel, policy, plant, cost, H)
```

Code

```

1 policy.p = p; % overwrite policy.p with new parameters from minimize
2 p = unwrap(policy.p); dp = 0*p;
3 m = m0; S = S0; L = zeros(1,H);
4
5 if nargout <= 1 % no derivatives required
6
7     for t = 1:H % for all time steps in horizon
8         [m, S] = plant.prop(m, S, plant, dynmodel, policy); % get next state
9         L(t) = cost.gamma^t.*cost.fcn(cost, m, S); % expected discounted cost
10    end
11
12 else % otherwise, get derivatives
13
14     dm0dp = zeros([size(m0,1), length(p)]);
15     dS0dp = zeros([size(m0,1)*size(m0,1), length(p)]);
16
17     for t = 1:H % for all time steps in horizon
18         [m, S, dmdm0, dSdm0, dmdS0, dSdS0, dmdp, dSdp] = ...
19         plant.prop(m, S, plant, dynmodel, policy); % get next state
20
21         dmdp = dmdm0*dm0dp + dmdS0*dS0dp + dmdp;
22         dSdp = dSdm0*dm0dp + dSdS0*dS0dp + dSdp;
23
24         [L(t), dLdm, dLdS] = cost.fcn(cost, m, S); % predictive cost
25         L(t) = cost.gamma^t*L(t); % discount
26         dp = dp + cost.gamma^t*( dLdm(:)'*dmdp + dLdS(:)'*dSdp );
27
28         dm0dp = dmdp; dS0dp = dSdp; % bookkeeping
29     end
30
31 end
32
33 J = sum(L); dJdp = rewrap(policy.p, dp);

```

7.2 Control Directory

7.2.1 concat.m

Summary: Compute a control signal u from a state distribution $x \sim \mathcal{N}(x|m, s)$. Here, the predicted control distribution and its derivatives are computed by concatenating a controller "con" with a saturation function "sat", such as gSat.m.

```

function [M, S, C, dMdm, dSdm, dCdm, dMds, dSds, dCds, dMdp, dSdp, dCdp] ...
    = conCat(con, sat, policy, m, s)

```

Example call: conCat(@congp, @gSat, policy, m, s)

Input arguments:

```

con      function handle (controller)
sat      function handle (squashing function)
policy   policy structure

```

.maxU	maximum amplitude of control signal (after squashing)	
m	mean of input distribution	[D x 1]
s	covariance of input distribution	[D x D]

Output arguments:

M	control mean	[E x 1]
S	control covariance	[E x E]
C	inv(s)*cov(x,u)	[D x E]
dMdm	deriv. of expected control wrt input mean	[E x D]
dSdm	deriv. of control covariance wrt input mean	[E*E x D]
dCdm	deriv. of C wrt input mean	[D*E x D]
dMds	deriv. of expected control wrt input covariance	[E x D*D]
dSds	deriv. of control covariance wrt input covariance	[E*E x D*D]
dCds	deriv. of C wrt input covariance	[D*E x D*D]
dMdp	deriv. of expected control wrt policy parameters	[E x P]
dSdp	deriv. of control covariance wrt policy parameters	[E*E x P]
dCdp	deriv. of C wrt policy parameters	[D*E x P]

where P is the total number of policy parameters

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2012-07-03

High-Level Steps

1. Compute unsquashed control signal
2. Compute squashed control signal

```
1 function [M, S, C, dMdm, dSdm, dCdm, dMds, dSds, dCds, dMdp, dSdp, dCdp] ...
2 = conCat(con, sat, policy, m, s)
```

Code

```
1 maxU=policy.maxU; % amplitude limit of control signal
2 E=length(maxU); % dimension of control signal
3 D=length(m); % dimension of input
4
5 % pre-compute some indices
6 F=D+E; j=D+1:F; i=1:D;
7 % initialize M and S
8 M = zeros(F,1); M(i) = m; S = zeros(F); S(i,i) = s;
9
10 if nargin < 4 % without derivatives
11 [M(j), S(j,j), Q] = con(policy, m, s); % compute unsquashed control signal v
```

```

12  q = S(i,i)*Q; S(i,j) = q; S(j,i) = q'; % compute joint covariance S=cov(x,v)
13  [M, S, R] = sat(M, S, j, maxU); % compute squashed control signal u
14  C = [eye(D) Q]*R; % inv(s)*cov(x,u)
15  else % with derivatives
16  Mdm = zeros(F,D); Sdm = zeros(F*F,D); Mdm(1:D,1:D) = eye(D);
17  Mds = zeros(F,D*D); Sds = kron(Mdm,Mdm);
18
19  X = reshape(1:F*F,[F F]); XT = X'; % vectorized indices
20  I=0*X; I(j,j)=1; jj=X(I==1)'; I=0*X; I(i,j)=1; ij=X(I==1)'; ji=XT(I==1)';
21
22  % 1. Unsquashed controller
23  [M(j), S(j,j), Q, Mdm(j,:), Sdm(jj,:), dQdm, Mds(j,:), ...
24   Sds(jj,:), dQds, Mdp, Sdp, dQdp] = con(policy, m, s);
25  q = S(i,i)*Q; S(i,j) = q; S(j,i) = q'; % compute joint covariance S=cov(x,v)
26
27  % update the derivatives
28  SS = kron(eye(E),S(i,i)); QQ = kron(Q',eye(D));
29  Sdm(ij,:) = SS*dQdm; Sdm(ji,:) = Sdm(ij,:);
30  Sds(ij,:) = SS*dQds + QQ; Sds(ji,:) = Sds(ij,:);
31
32  % 2. Apply Saturation
33  [M, S, R, MdM, SdM, RdM, MdS, SdS, RdS] = sat(M, S, j, maxU);
34
35  % apply chain-rule to compute derivatives after concatenation
36  dMdm = MdM*Mdm + MdS*Sdm; dMds = MdM*Mds + MdS*Sds;
37  dSdm = SdM*Mdm + SdS*Sdm; dSds = SdM*Mds + SdS*Sds;
38  dRdm = RdM*Mdm + RdS*Sdm; dRds = RdM*Mds + RdS*Sds;
39
40  dMdp = MdM(:,j)*Mdp + MdS(:,jj)*Sdp;
41  dSdp = SdM(:,j)*Mdp + SdS(:,jj)*Sdp;
42  dRdp = RdM(:,j)*Mdp + RdS(:,jj)*Sdp;
43
44  C = [eye(D) Q]*R; % inv(s)*cov(x,u)
45  % update the derivatives
46  RR = kron(R(j,:)',eye(D)); QQ = kron(eye(E),[eye(D) Q]);
47  dCdm = QQ*dRdm + RR*dQdm;
48  dCds = QQ*dRds + RR*dQds;
49  dCdp = QQ*dRdp + RR*dQdp;
50  end

```

7.2.2 congpm

Summary: Implements the mean-of-GP policy (equivalent to a regularized RBF network. Compute mean, variance and input-output covariance of the control u using a mean-of-GP policy function, when the input x is Gaussian. The GP is parameterized using a pseudo training set size N . Optionally, compute partial derivatives wrt the input parameters.

This version sets the signal variance to 1, the noise to 0.01 and their respective lengthscales to zero. This results in only the lengthscales, inputs, and outputs being trained.

```

function [M, S, C, dMdm, dSdm, dCdm, dMds, dSds, dCds, dMdp, dSdp, dCdp] ...
    = congpm(policy, m, s)

```

Input arguments:

policy	policy (struct)	
.p	parameters that are modified during training	
.hyp	GP-log hyperparameters ($P_h = (d+2)*D$)	[P_h]
.inputs	policy pseudo inputs	[N x d]
.targets	policy pseudo targets	[N x D]
m	mean of state distribution	[d x 1]
s	covariance matrix of state distribution	[d x d]

Output arguments:

M	mean of the predicted control	[D x 1]
S	covariance of predicted control	[D x D]
C	$\text{inv}(s) \cdot \text{covariance between input and control}$	[d x D]
dMdm	deriv. of mean control wrt mean of state	[D x d]
dSdm	deriv. of control variance wrt mean of state	[$D \times D$ x d]
dCdm	deriv. of covariance wrt mean of state	[$d \times D$ x d]
dMds	deriv. of mean control wrt variance	[D x $d \times d$]
dSds	deriv. of control variance wrt variance	[$D \times D$ x $d \times d$]
dCds	deriv. of covariance wrt variance	[$d \times D$ x $d \times d$]
dMdp	deriv. of mean control wrt GP hyper-parameters	[D x P]
dSdp	deriv. of control variance wrt GP hyper-parameters	[$D \times D$ x P]
dCdp	deriv. of covariance wrt GP hyper-parameters	[$d \times D$ x P]

where $P = (d+2)*D + n*(d+D)$ is the total number of policy parameters.

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2013-01-24

High-Level Steps

1. Extract policy parameters from policy structure
2. Compute predicted control $u = \text{inv}(s) \cdot \text{covariance between input and control}$
3. Set derivatives of non-free parameters to zero
4. Merge derivatives

```

1 function [M, S, C, dMdm, dSdm, dCdm, dMds, dSds, dCds, dMdp, dSdp, dCdp] ...
2     = congpr(policy, m, s)

```

Code

```

1 % 1. Extract policy parameters
2 policy.hyp = policy.p.hyp;
3 policy.inputs = policy.p.inputs;

```

```

4 policy.targets = policy.p.targets;
5
6 % fix policy signal and the noise variance
7 % (avoids some potential numerical problems)
8 policy.hyp(end-1,:) = log(1); % set signal variance to 1
9 policy.hyp(end,:) = log(0.01); % set noise standard dev to 0.01
10
11 % 2. Compute predicted control u inv(s)*covariance between input and control
12 if nargout < 4 % if no derivatives are required
13     [M, S, C] = gp2(policy, m, s);
14 else % else compute derivatives too
15     [M, S, C, dMdm, dSdm, dCdm, dMds, dSds, dCds, dMdi, dSdi, dCdi, dMdt, ...
16      dSdt, dCdt, dMdh, dSdh, dCdh] = gp2d(policy, m, s);
17
18 % 3. Set derivatives of non-free parameters to zero: signal and noise variance
19 d = size(policy.inputs,2);
20 d2 = size(policy.hyp,1); dimU = size(policy.targets,2);
21 idx = bsxfun(@plus,(d+1:d2)',(0:dimU-1)*d2);
22 dMdh(:,idx(:)) = 0; dSdh(:,idx(:)) = 0; dCdh(:,idx(:)) = 0;
23
24 % 4. Merge derivatives
25 dMdp = [dMdh dMdi dMdt]; dSdp = [dSdh dSdi dSdt]; dCdp = [dCdh dCdi dCdt];
26 end

```

7.2.3 conlin.m

Summary: Affine controller $u = Wx + b$ with input dimension D and control dimension E . Compute mean and covariance of the control distribution $p(u)$ from a Gaussian distributed input $x \sim \mathcal{N}(x|m, s)$. Moreover, the $s^{-1}cov(x, u)$ is computed.

```

function [M, S, V, dMdm, dSdm, dVdm, dMds, dSds, dVds, dMdp, dSdp, dVdp] ...
    = conlin(policy, m, s)

```

Input arguments:

policy	policy structure	
.p	parameters that are modified during training	
.w	linear weights	[E x D]
.b	biases/offset	[E]
m	mean of state distribution	[D]
s	covariance matrix of state distribution	[D x D]

Output arguments:

M	mean of predicted control	[E]
S	variance of predicted control	[E x E]
C	inv(s) times input-output covariance	[D x E]
dMdm	deriv. of mean control wrt input mean	[E x D]
dSdm	deriv. of control covariance wrt input mean	[E*E x D]
dCdm	deriv. of C wrt input mean	[D*E x D]
dMds	deriv. of mean control wrt input covariance	[E x D*D]

dSds	deriv. of control covariance wrt input covariance	$[E \times E \times D \times D]$
dCds	deriv. of C wrt input covariance	$[D \times E \times D \times D]$
dMdp	deriv. of mean control wrt policy parameters	$[E \times P]$
dSdp	deriv. of control covariance wrt policy parameters	$[E \times E \times P]$
dCdp	deriv. of C wrt policy parameters	$[D \times E \times P]$

where $P = (D+1) \times E$ is the total number of policy parameters

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2012-07-03

High-Level Steps

1. Extract policy parameters from policy structure
2. Predict control signal
3. Compute derivatives if required

```
1 function [M, S, V, dMdm, dSdm, dVdm, dMds, dSds, dVds, dMdp, dSdp, dVdp] ...
2     = conlin(policy, m, s)
```

Code

```
1 % 1. Extract policy parameters from policy structure
2 w = policy.p.w; % weight matrix
3 b = policy.p.b; % bias/offset
4 [E D] = size(w); % dim of control and state
5
6 % 2. Predict control signal
7 M = w*m + b; % mean
8 S = w*s*w'; S = (S+S')/2; % covariance
9 V = w'; % inv(s)*input-output covariance
10
11 % 3. Compute derivatives if required
12 if nargin > 3
13     dMdm = w; dSdm = zeros(E,E,D); dVdm = zeros(D,E,D);
14     dMds = zeros(E,D*D); dSds = kron(w,w); dVds = zeros(D,E,D*D);
15
16     X=reshape(1:D*D,[D D]); XT=X'; dSds=(dSds+dSds(:,XT(:)))/2; % symmetrize
17     X=reshape(1:E*E,[E E]); XT=X'; dSds=(dSds+dSds(XT(:),:))/2;
18
19     wTdw = reshape(permute(reshape(eye(E*D),[E D E D]),[2 1 3 4]),[E*D E*D]);
20     dMdp = [eye(E) kron(m',eye(E))];
21     dSdp = [zeros(E*E,E) kron(eye(E),w*s)*wTdw + kron(w*s,eye(E))];
22     dSdp = (dSdp + dSdp(XT(:),:))/2; % symmetrize
23     dVdp = [zeros(D,E,E) wTdw];
24 end
```


7.3 GP Directory

7.3.1 train.m

Summary: Train a GP model with SE covariance function (ARD). First, the hyper-parameters are trained using a full GP. Then, if `gpmodel.induce` exists, indicating sparse approximation, if enough training exmples are present, train the inducing inputs (hyper-parameters are taken from the full GP). If no inducing inputs are present, then initialize thme to be a random subset of the training inputs.

```
function [gpmodel nlml] = train(gpmodel, dump, iter)
```

Input arguments:

<code>gpmodel</code>	GP structure	
<inputs></inputs>	GP training inputs	[N x D]
targets	GP training targets	[N x E]
hyp (optional)	GP log-hyper-parameters	[D+2 x E]
induce (optional)	pseudo inputs for sparse GP	
<code>dump</code>	not needed for this code, but required for compatibility reasons	
<code>iter</code>	optimization iterations for training	[1 x 2]
	[full GP, sparse GP]	

Output arguments:

<code>gpmodel</code>	updated GP structure
<code>nlml</code>	negative log-marginal likelihood

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2013-05-16

High-Level Steps

1. Initialization
2. Train full GP
3. If necessary: train sparse GP (FITC/SPGP)

```
1 function [gpmodel nlml] = train(gpmodel, dump, iter)
```

Code

```

1 % 1) Initialization
2 if nargin < 3, iter = [-500 -1000]; end % default training iterations
3
4 D = size(gpmodel.inputs,2); E = size(gpmodel.targets,2); % get variable sizes
5 covfunc = {'covSum', {'covSEard', 'covNoise'}}; % specify ARD covariance
6 curb.snr = 1000; curb.ls = 100; curb.std = std(gpmodel.inputs); % set hyp curb
7 if ~isfield(gpmodel, 'hyp') % if we don't pass in hyper-parameters, define them
8     gpmodel.hyp = zeros(D+2,E); nlml = zeros(1,E);
9     lh = repmat([log(std(gpmodel.inputs)) 0 -1]',1,E); % init hyp length scales
10    lh(D+1,:) = log(std(gpmodel.targets)); % signal std dev
11    lh(D+2,:) = log(std(gpmodel.targets)/10); % noise std dev
12 else
13     lh = gpmodel.hyp; % GP hyper-parameters
14 end
15
16 % 2a) Train full GP (always)
17 fprintf('Train hyper-parameters of full GP ...\n');
18 for i = 1:E % train each GP separately
19     fprintf('GP %i/%i\n', i, E);
20     try % BFGS training
21         [gpmodel.hyp(:,i) v] = minimize(lh(:,i), @hypCurb, iter(1), covfunc, ...
22             gpmodel.inputs, gpmodel.targets(:,i), curb);
23         catch % conjugate gradients (BFGS can be quite aggressive)
24             [gpmodel.hyp(:,i) v] = minimize(lh(:,i), @hypCurb, ...
25                 struct('length', iter(1), 'method', 'CG', 'verbosity', 1), covfunc, ...
26                 gpmodel.inputs, gpmodel.targets(:,i), curb);
27         end
28         nlml(i) = v(end);
29     end
30
31 % 2b) If necessary: sparse training using FITC/SPGP (Snelson & Ghahramani, 2006)
32 if isfield(gpmodel, 'induce') % are we using a sparse approximation?
33     [N D] = size(gpmodel.inputs);
34     [M uD uE] = size(gpmodel.induce);
35     if M >= N; return; end % if too few training examples, we don't need FITC
36     fprintf('Train pseudo-inputs of sparse GP ...\n');
37
38     if uD == 0 % we don't have inducing inputs yet?
39         gpmodel.induce = zeros(M,D,uE); % allocate space
40         iter = 3*iter; % train a lot for the first time (it's still cheap!)
41         [cidx, ctrs] = kmeans(gpmodel.inputs, M); % kmeans: initialize pseudo inputs
42         for i = 1:uD
43             % j = randperm(N);
44             % gpmodel.induce(:, :, i) = gpmodel.inputs(j(1:M), :); % random subset
45             gpmodel.induce(:, :, i) = ctrs;
46         end
47     end
48
49 % train sparse model
50 [gpmodel.induce nlml2] = minimize(gpmodel.induce, 'fite', iter(end), gpmodel);
51 fprintf('GP NLML, full: %e, sparse: %e, diff: %e\n', ...
52     sum(nlml), nlml2(end), nlml2(end)-sum(nlml));
53 end

```

7.3.2 hypCurb.m

Summary: Wrapper for GP training (via gpr.m), penalizing large SNR and extreme length-scales to

avoid numerical instabilities

```
function [f df] = hypCurb(lh, covfunc, x, y, curb)
```

Input arguments:

lh	log-hyper-parameters	[D+2 x E]
covfunc	covariance function, e.g.,	
	covfunc = \{'covSum', \{'covSEard', 'covNoise'\}\};	
x	training inputs	[n x D]
y	training targets	[n x E]
curb	(optional) parameters to penalize extreme hyper-parameters	
.ls	length-scales	
.snr	signal-to-noise ratio (try to keep it below 500)	
.std	additional parameter required for length-scale penalty	

Output arguments:

f	penalized negative log-marginal likelihood
df	derivative of penalized negative log-marginal likelihood wrt GP log-hyper-parameters

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2011-12-19

High-Level Steps

1. Compute the negative log-marginal likelihood (plus derivatives)
2. Add penalties and change derivatives accordingly

```
1 function [f df] = hypCurb(lh, covfunc, x, y, curb)
```

Code

```
1 if nargin < 5, curb.snr = 500; curb.ls = 100; curb.std = 1; end % set default
2
3 p = 30; % penalty power
4 D = size(x,2);
5 if size(lh,1) == 3*D+2; li = 1:2*D; sfi = 2*D+1:3*D+1; % 1D and DD terms
6 elseif size(lh,1) == 2*D+1; li = 1:D; sfi = D+1:2*D; % Just 1D terms
7 elseif size(lh,1) == D+2; li = 1:D; sfi = D+1; % Just DD terms
8 else error('Incorrect number of hyperparameters');
9 end
```

```

10
11 ll = lh(li); lsf = lh(sfi); lsn = lh(end);
12
13 % 1) compute the negative log-marginal likelihood (plus derivatives)
14 [f df] = gpr(lh, covfunc, x, y); % first, call gpr
15
16 % 2) add penalties and change derivatives accordingly
17 f = f + sum(((ll - log(curb.std'))./log(curb.ls)).^p); % length-scales
18 df(li) = df(li) + p*(ll - log(curb.std')).^(p-1)/log(curb.ls)^p;
19
20 f = f + sum(((lsf - lsn)/log(curb.snr)).^p); % signal to noise ratio
21 df(sfi) = df(sfi) + p*(lsf - lsn).^(p-1)/log(curb.snr)^p;
22 df(end) = df(end) - p*sum((lsf - lsn).^(p-1)/log(curb.snr)^p);

```

7.3.3 fitc.m

Summary: Compute the FITC negative log marginal likelihood and its derivatives with respect to the inducing inputs (we don't compute the derivatives with respect to the GP hyper-parameters)

```
function [nml dnml] = fitc(induce, gpmodel)
```

Input arguments:

induce	matrix of inducing inputs	[M x D x uE]
	M: number of inducing inputs	
	E: either 1 (inducing inputs are shared across target dim.)	
	or E (different inducing inputs for each target dim.)	
gpmodel	GP structure	
.hyp	log-hyper-parameters	[D+2 x E]
.inputs	training inputs	[N x D]
.targets	training targets	[N x E]
.noise (opt)	noise	

Output arguments:

nml	negative log-marginal likelihood
dnml	derivative of negative log-marginal likelihood wrt inducing inputs

Adapted from Ed Snelson's SPGP code.

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2013-05-21

High-Level Steps

1. Compute FITC marginal likelihood
2. Compute corresponding gradients wrt the pseudo inputs

```
1 function [nlml dnlml] = fitc(induce, gpmodel)
```

Code

```
1 ridge = 1e-06; % jitter to make matrix better conditioned
2
3 [N D] = size(gpmodel.inputs); E = size(gpmodel.targets,2);
4 [M uD uE] = size(induce);
5 if uD ~= D || (uE~=1 && uE ~= E); error('Wrong size of inducing inputs'); end
6
7 nlml = 0; dfxb = zeros(M, D); dnlml = zeros(M, D, E); % zero and allocate outputs
8
9 for j = 1:E
10     if uE > 1; u = induce(:,j); else u = induce; end
11     b = exp(gpmodel.hyp(1:D,j)); % length-scales
12     c = gpmodel.hyp(D+1,j); % log signal std dev
13     sig = exp(2.*gpmodel.hyp(D+2,j)); % noise variance
14
15     xb = bsxfun(@rdivide,u,b'); % divide inducing by lengthscales
16     x = bsxfun(@rdivide, gpmodel.inputs,b'); % divide inputs by length-scales
17     y = gpmodel.targets(:,j); % training targets
18
19     Kmm = exp(2*c-maha(xb,xb)/2) + ridge*eye(M);
20     Kmn = exp(2*c-maha(xb,x)/2);
21
22     % Check whether Kmm is no longer positive definite. If so, return
23     try
24         L = chol(Kmm)';
25     catch
26         nlml = Inf; dnlml = zeros(size(params));
27         return;
28     end
29     V = L\Kmn; % inv(sqrt(Kmm))*Kmn
30
31     if isfield(gpmodel,'noise')
32         Gamma = 1 + (exp(2*c)-sum(V.^2)'+gpmodel.noise(:,j))/sig;
33     else
34         Gamma = 1 + (exp(2*c)-sum(V.^2)')/sig; % Gamma = diag(Knn-Qnn)/sig + I
35     end
36
37     V = bsxfun(@rdivide,V,sqrt(Gamma)'); % inv(sqrt(Kmm))*Kmn * inv(sqrt(Gamma))
38     y = y./sqrt(Gamma);
39     Am = chol(sig*eye(M) + V*V'); % chol(inv(sqrt(Kmm))*A*inv(sqrt(Kmm)))
40     % V*V' = inv(chol(Kmm)')*K*inv(diag(Gamma))*K'*inv(chol(Kmm)')
41     Vy = V*y;
42     beta = Am\Vy;
43
44     nlml = nlml + sum(log(diag(Am))) + (N-M)/2*log(sig) + sum(log(Gamma))/2 ...
```

```

45         + (y'*y - beta'*beta)/2/sig + 0.5*N*log(2*pi);
46
47     if nargout == 2 % ... and if requested, its partial derivatives
48
49         At = L*Am; iAt = At\eye(M); % chol(sig*B) [Ed's thesis, p. 40]
50         iA = iAt'*iAt; % inv(sig*B)
51
52         iAmV = Am\V; % inv(Am)*V
53         B1 = At'\(iAmV);
54         b1 = At'\beta; % b1 = B1*y
55
56         iLV = L'\V; % inv(Kmm)*Kmm*inv(sqrt(Gamma))
57         iL = L\eye(M);
58         iKmm = iL'*iL;
59
60         mu = ((Am'\beta)*V)';
61         bs = y.*(beta'*iAmV)/sig - sum(iAmV.*iAmV)'/2 - (y.^2+mu.^2)/2/sig + 0.5;
62         TT = iLV*(bsxfun(@times,iLV',bs));
63         Kmm = bsxfun(@rdivide,Kmm,sqrt(Gamma)'); % overwrite Kmm
64
65         for i = 1:D % derivatives wrt inducing inputs
66             dsq_mn = bsxfun(@minus,xb(:,i),xb(:,i)').*Kmm;
67             dsq_mn = bsxfun(@minus,-xb(:,i),-x(:,i)').*Kmm;
68             dGamma = -2/sig*dsq_mn.*iLV;
69
70             dfxb(:,i) = -b1.*(dsq_mn*(y-mu)/sig + dsq_mn*b1) + dGamma*bs ...
71                 + sum((iKmm - iA*sig).*dsq_mn,2) - 2/sig*sum(dsq_mn.*TT,2);
72             dsq_mn = dsq_mn.*B1; % overwrite dsq_mn
73             dfxb(:,i) = dfxb(:,i) + sum(dsq_mn,2);
74             dfxb(:,i) = dfxb(:,i)/b(i);
75         end
76
77         dnllm(:,j) = dfxb;
78     end
79 end
80 if 1 == uE; dnllm = sum(dnllm,3); end % combine derivatives if sharing inducing

```

7.3.4 gp0.m

Summary: Compute joint predictions for multiple GPs with uncertain inputs. If `gpmodel.nigp` exists, individual noise contributions are added. Predictive variances contain uncertainty about the function, but no noise.

```
function [M, S, V] = gp0(gpmodel, m, s)
```

Input arguments:

<code>gpmodel</code>	GP model struct	
<code>hyp</code>	log-hyper-parameters	[D+2 x E]
<code>inputs</code>	training inputs	[n x D]
<code>targets</code>	training targets	[n x E]
<code>nigp</code>	(optional) individual noise variance terms	[n x E]
<code>m</code>	mean of the test distribution	[D x 1]
<code>s</code>	covariance matrix of the test distribution	[D x D]

Output arguments:

M	mean of pred. distribution	$[E \ x \ 1]$
S	covariance of the pred. distribution	$[E \ x \ E]$
V	inv(s) times covariance between input and output	$[D \ x \ E]$

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2013-05-24

High-Level Steps

1. If necessary, compute kernel matrix and cache it
2. Compute predicted mean and inv(s) times input-output covariance
3. Compute predictive covariance matrix, non-central moments
4. Centralize moments

```
1 function [M, S, V] = gp0(gpmodel, m, s)
```

Code

```
1 persistent K iK beta oldX oldn;
2 [n, D] = size(gpmodel.inputs); % number of examples and dimension of inputs
3 [n, E] = size(gpmodel.targets); % number of examples and number of outputs
4 X = gpmodel.hyp; % short hand for hyperparameters
5
6 % 1) if necessary: re-compute cached variables
7 if numel(X) ~= numel(oldX) || isempty(iK) || sum(any(X ~= oldX)) || n ~= oldn
8     oldX = X; oldn = n;
9     iK = zeros(n,n,E); K = zeros(n,n,E); beta = zeros(n,E);
10
11     for i=1:E % compute K and inv(K)
12         inp = bsxfun(@divide, gpmodel.inputs, exp(X(1:D,i)'));
13         K(:, :, i) = exp(2*X(D+1,i)-maha(inp,inp)/2);
14         if isfield(gpmodel, 'nigp')
15             L = chol(K(:, :, i) + exp(2*X(D+2,i))*eye(n) + diag(gpmodel.nigp(:,i)))';
16         else
17             L = chol(K(:, :, i) + exp(2*X(D+2,i))*eye(n))';
18         end
19         iK(:, :, i) = L' \ (L \ eye(n));
20         beta(:, i) = L' \ (L \ gpmodel.targets(:, i));
21     end
22 end
23
24 k = zeros(n,E); M = zeros(E,1); V = zeros(D,E); S = zeros(E);
25
26 inp = bsxfun(@minus, gpmodel.inputs, m'); % centralize inputs
27
```

```

28 % 2) compute predicted mean and inv(s) times input-output covariance
29 for i=1:E
30     iL = diag(exp(-X(1:D,i))); % inverse length-scales
31     in = inp*iL;
32     B = iL*s*iL+eye(D);
33
34     t = in/B;
35     l = exp(-sum(in.*t,2)/2); lb = l.*beta(:,i);
36     tiL = t*iL;
37     c = exp(2*X(D+1,i))/sqrt(det(B));
38
39     M(i) = sum(lb)*c; % predicted mean
40     V(:,i) = tiL'*lb*c; % inv(s) times input-output covariance
41     k(:,i) = 2*X(D+1,i)-sum(in.*in,2)/2;
42 end
43
44 % 3) compute predictive covariance, non-central moments
45 for i=1:E
46     ii = bsxfun(@rdivide,inp,exp(2*X(1:D,i)'));
47
48     for j=1:i
49         R = s*diag(exp(-2*X(1:D,i))+exp(-2*X(1:D,j)))+eye(D);
50         t = 1/sqrt(det(R));
51         ij = bsxfun(@rdivide,inp,exp(2*X(1:D,j)'));
52         L = exp(bsxfun(@plus,k(:,i),k(:,j)')+maha(ii,-ij,R\s/2));
53         if i==j
54             S(i,i) = t*(beta(:,i)'*L*beta(:,i) - sum(sum(iK(:, :, i).*L)));
55         else
56             S(i,j) = beta(:,i)'*L*beta(:,j)*t;
57             S(j,i) = S(i,j);
58         end
59     end
60
61     S(i,i) = S(i,i) + exp(2*X(D+1,i));
62 end
63
64 % 4) centralize moments
65 S = S - M*M';

```

7.3.5 gp1.m

Summary: Compute joint predictions for the FITC sparse approximation to multiple GPs with uncertain inputs. Predictive variances contain uncertainty about the function, but no noise. If gp-model.nigp exists, individual noise contributions are added.

```
function [M, S, V] = gp1d(gpmodel, m, s)
```

Input arguments:

gpmodel	GP model struct	
hyp	log-hyper-parameters	[D+2 x E]
inputs	training inputs	[n x D]
targets	training targets	[n x E]
nigp	(optional) individual noise variance terms	[n x E]

m	mean of the test distribution	[D x 1]
s	covariance matrix of the test distribution	[D x D]

Output arguments:

M	mean of pred. distribution	[E x 1]
S	covariance of the pred. distribution	[E x E]
V	inv(s) times covariance between input and output	[D x E]

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2013-03-05

High-Level Steps

1. If necessary, compute kernel matrix and cache it
2. Compute predicted mean and inv(s) times input-output covariance
3. Compute predictive covariance matrix, non-central moments
4. Centralize moments

```
1 function [M, S, V] = gp1(gpmodel, m, s)
```

Code

```
1 if ~isfield(gpmodel, 'induce') || numel(gpmodel.induce)==0,
2     [M, S, V] = gp0(gpmodel, m, s); return; end
3
4 persistent iK iK2 beta oldX;
5 ridge = 1e-6; % jitter to make matrix better conditioned
6 [n, D] = size(gpmodel.inputs); % number of examples and dimension of inputs
7 E = size(gpmodel.targets, 2); % number of examples and number of outputs
8 X = gpmodel.hyp; input = gpmodel.inputs; targets = gpmodel.targets;
9
10 [np pD pE] = size(gpmodel.induce); % number of pseudo inputs per dimension
11 pinput = gpmodel.induce; % all pseudo inputs
12
13 % 1) If necessary: re-compute cached variables
14 if numel(X) ~= numel(oldX) || isempty(iK) || isempty(iK2) || ... % if necessary
15     sum(any(X ~= oldX)) || numel(iK2) ~= E*np^2 || numel(iK) ~= n*np*E
16     oldX = X; % compute K, inv(K), inv(K2)
17     iK = zeros(np, n, E); iK2 = zeros(np, np, E); beta = zeros(np, E);
18
19     for i=1:E
20         pinp = bsxfun(@rdivide, pinput(:, :, min(i, pE)), exp(X(1:D, i)'));
21         inp = bsxfun(@rdivide, input, exp(X(1:D, i)'));
22         Kmm = exp(2*X(D+1, i)-maha(pinp, pinp)/2) + ridge*eye(np); % add small ridge
23         Kmn = exp(2*X(D+1, i)-maha(pinp, inp)/2);
24         L = chol(Kmm)';
```

```

25     V = L\Kmn; % inv(sqrt(Kmm))*Kmn
26     if isfield(gpmodel, 'nigp')
27         G = exp(2*X(D+1,i))-sum(V.^2)+gpmodel.nigp(:,i)';
28     else
29         G = exp(2*X(D+1,i))-sum(V.^2);
30     end
31     G = sqrt(1+G/exp(2*X(D+2,i)));
32     V = bsxfun(@rdivide,V,G);
33     Am = chol(exp(2*X(D+2,i))*eye(np) + V*V')';
34     At = L*Am; % chol(sig*B) [thesis, p. 40]
35     iAt = At\eye(np);
36 % The following is not an inverse matrix, but we'll treat it as such: multiply
37 % the targets from right and the cross-covariances left to get predictive mean.
38     iK(:, :, i) = ((Am\bsxfun(@rdivide,V,G))'*iAt)';
39     beta(:, i) = iK(:, :, i)*targets(:, i);
40     iB = iAt'*iAt.*exp(2*X(D+2,i)); % inv(B), [Ed's thesis, p. 40]
41     iK2(:, :, i) = Kmm\eye(np) - iB; % covariance matrix for predictive variances
42 end
43 end
44
45 k = zeros(np,E); M = zeros(E,1); V = zeros(D,E); S = zeros(E); % allocate
46 inp = zeros(np,D,E);
47
48 % 2) Compute predicted mean and inv(s) times input-output covariance
49 for i=1:E
50     inp(:, :, i) = bsxfun(@minus, pinput(:, :, min(i,pE)), m');
51
52     L = diag(exp(-X(1:D,i)));
53     in = inp(:, :, i)*L;
54     B = L*s*L+eye(D);
55
56     t = in/B;
57     l = exp(-sum(in.*t,2)/2); lb = 1.*beta(:,i);
58     tL = t*L;
59     c = exp(2*X(D+1,i))/sqrt(det(B));
60
61     M(i) = sum(lb)*c; % predicted mean
62     V(:, i) = tL'*lb*c; % inv(s) times input-output covariance
63     k(:, i) = 2*X(D+1,i)-sum(in.*in,2)/2;
64 end
65
66 % 3) Compute predictive covariance matrix, non-central moments
67 for i=1:E
68     ii = bsxfun(@rdivide, inp(:, :, i), exp(2*X(1:D,i)'));
69
70     for j=1:i
71         R = s*diag(exp(-2*X(1:D,i))+exp(-2*X(1:D,j)))+eye(D); t = 1./sqrt(det(R));
72         ij = bsxfun(@rdivide, inp(:, :, j), exp(2*X(1:D,j)'));
73         L = exp(bsxfun(@plus, k(:, i), k(:, j)')+maha(ii,-ij,R\s/2));
74         if i==j
75             S(i,i) = t*(beta(:,i)'*L*beta(:,i) - sum(sum(iK2(:, :, i).*L)));
76         else
77             S(i,j) = beta(:,i)'*L*beta(:,j)*t; S(j,i) = S(i,j);
78         end
79     end
80
81     S(i,i) = S(i,i) + exp(2*X(D+1,i));
82 end

```

```

83
84 % 4) Centralize moments
85 S = S - M*M';

```

7.3.6 gp2.m

Summary: Compute joint predictions and derivatives for multiple GPs with uncertain inputs. Does not consider the uncertainty about the underlying function (in prediction), hence, only the GP mean function is considered. Therefore, this representation is equivalent to a regularized RBF network. If `gpmodel.nigp` exists, individual noise contributions are added.

```
function [M, S, V] = gp2(gpmodel, m, s)
```

Input arguments:

<code>gpmodel</code>	GP model struct	
<code>hyp</code>	log-hyper-parameters	[D+2 x E]
<code>inputs</code>	training inputs	[n x D]
<code>targets</code>	training targets	[n x E]
<code>nigp</code>	(optional) individual noise variance terms	[n x E]
<code>m</code>	mean of the test distribution	[D x 1]
<code>s</code>	covariance matrix of the test distribution	[D x D]

Output arguments:

<code>M</code>	mean of pred. distribution	[E x 1]
<code>S</code>	covariance of the pred. distribution	[E x E]
<code>V</code>	inv(s) times covariance between input and output	[D x E]

Copyright (C) 2008-2013 by Marc Deisenroth, Andrew McHutchon, Joe Hall, and Carl Edward Rasmussen.

Last modified: 2013-03-05

High-Level Steps

1. If necessary, re-compute cached variables
2. Compute predicted mean and $\text{inv}(s)$ times input-output covariance
3. Compute predictive covariance matrix, non-central moments
4. Centralize moments

```
1 function [M, S, V] = gp2(gpmodel, m, s)
```

Code

```

1 persistent iK oldX oldIn oldOut beta oldn;
2 D = size(gpmodel.inputs,2); % number of examples and dimension of inputs
3 [n, E] = size(gpmodel.targets); % number of examples and number of outputs
4
5 input = gpmodel.inputs; target = gpmodel.targets; X = gpmodel.hyp;
6
7 % 1) if necessary: re-compute cached variables
8 if numel(X) ~= numel(oldX) || isempty(iK) || n ~= oldn || ...
9     sum(any(X ~= oldX)) || sum(any(oldIn ~= input)) || ...
10    sum(any(oldOut ~= target))
11    oldX = X; oldIn = input; oldOut = target; oldn = n;
12    K = zeros(n,n,E); iK = K; beta = zeros(n,E);
13
14    for i=1:E % compute K and inv(K)
15        inp = bsxfun(@rdivide, gpmodel.inputs, exp(X(1:D,i)'));
16        K(:, :, i) = exp(2*X(D+1,i)-maha(inp,inp)/2);
17        if isfield(gpmodel, 'nigp')
18            L = chol(K(:, :, i) + exp(2*X(D+2,i))*eye(n) + diag(gpmodel.nigp(:,i)))';
19        else
20            L = chol(K(:, :, i) + exp(2*X(D+2,i))*eye(n))';
21        end
22        iK(:, :, i) = L'\(L\eye(n));
23        beta(:, i) = L'\(L\gpmodel.targets(:,i));
24    end
25 end
26
27 k = zeros(n,E); M = zeros(E,1); V = zeros(D,E); S = zeros(E);
28
29 inp = bsxfun(@minus, gpmodel.inputs, m'); % centralize inputs
30
31 % 2) Compute predicted mean and inv(s) times input-output covariance
32 for i=1:E
33     iL = diag(exp(-X(1:D,i))); % inverse length-scales
34     in = inp*iL;
35     B = iL*s*iL+eye(D);
36
37     t = in/B;
38     l = exp(-sum(in.*t,2)/2); lb = l.*beta(:,i);
39     tL = t*iL;
40     c = exp(2*X(D+1,i))/sqrt(det(B));
41
42     M(i) = sum(lb)*c; % predicted mean
43     V(:,i) = tL'*lb*c; % inv(s) times input-output covariance
44     k(:,i) = 2*X(D+1,i)-sum(in.*in,2)/2;
45 end
46
47 % 3) Compute predictive covariance, non-central moments
48 for i=1:E
49     ii = bsxfun(@rdivide, inp, exp(2*X(1:D,i)'));
50
51     for j=1:i
52         R = s*diag(exp(-2*X(1:D,i))+exp(-2*X(1:D,j)))+eye(D);
53         t = 1/sqrt(det(R));
54         ij = bsxfun(@rdivide, inp, exp(2*X(1:D,j)'));
55         L = exp(bsxfun(@plus, k(:,i), k(:,j)')+maha(ii,-ij,R\s/2));

```

```

56     S(i,j) = t*beta(:,i)'*L*beta(:,j); S(j,i) = S(i,j);
57 end
58
59     S(i,i) = S(i,i) + 1e-6;           % add small jitter for numerical reasons
60
61 end
62
63 % 4) Centralize moments
64 S = S - M*M';

```

Bibliography

- [1] James A. Bagnell and Jeff G. Schneider. Autonomous Helicopter Control using Reinforcement Learning Policy Search Methods. In *Proceedings of the International Conference on Robotics and Automation*, 2001.
- [2] Marc 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.
- [3] Marc P. Deisenroth and Carl E. Rasmussen. PILCO: A Model-Based and Data-Efficient Approach to Policy Search. In *Proceedings of the International Conference on Machine Learning*, pages 465–472, New York, NY, USA, June 2011. ACM.
- [4] Marc P. Deisenroth, Carl E. Rasmussen, and Dieter Fox. Learning to Control a Low-Cost Manipulator using Data-Efficient Reinforcement Learning. In *Proceedings of the International Conference on Robotics: Science and Systems*, Los Angeles, CA, USA, June 2011.
- [5] David Forster. Robotic Unicycle. Report, Department of Engineering, University of Cambridge, UK, 2009.
- [6] Joaquin Quiñonero-Candela, Agathe Girard, Jan Larsen, and Carl E. Rasmussen. Propagation of Uncertainty in Bayesian Kernel Models—Application to Multiple-Step Ahead Forecasting. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, volume 2, pages 701–704, April 2003.
- [7] Carl E. Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA, 2006.
- [8] Edward Snelson and Zoubin 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.
- [9] Mark W. Spong and Daniel J. Block. The Pendubot: A Mechatronic System for Control Research and Education. In *Proceedings of the Conference on Decision and Control*, pages 555–557, 1995.
- [10] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, 1998.