

# Reinforcement Learning and Trajectory Planning based on Model Approximation with Neural Networks applied to Transition Problems

Max Pritzkoleit\* Carsten Knoll\* Klaus Röbenack\*

\* TU Dresden, Faculty of Electrical and Computer Engineering,  
Institute of Control Theory, Dresden, Germany (e-mail:  
{max.pritzkoleit, carsten.knoll, klaus.roebenack}@tu-dresden.de).

---

**Abstract:** In this paper we use a multilayer neural network to approximate the dynamics of nonlinear (mechanical) control systems. Furthermore, these neural network models are combined with offline trajectory planning, to form a model-based reinforcement learning (RL) algorithm, suitable for transition problems of nonlinear dynamical systems. We evaluate the algorithm on the swing-up of the cart-pole benchmark system and observe a significant performance gain in terms of data efficiency compared to a state-of-the-art model-free RL method (Deep Deterministic Policy Gradient (DDPG)). Additionally, we present first experimental results on a cart-triple-pole system test bench. For a simple transition problem, the proposed algorithm shows a good controller performance.

*Keywords:* Trajectory planning, Reinforcement learning, Learning control, Neural-network models, Model approximation, Tracking control,

---

## 1. INTRODUCTION

In recent years, reinforcement learning gained a lot of attention in- and outside the scientific community through advances made possible by incorporating deep neural networks for function approximation. These model-free<sup>1</sup> algorithms (i.e. DDPG, DQN) try to approximate the Hamiltonian (often called “Q-function”) of the underlying optimal control problem based on data, which is generated by interacting with the system. The main drawback of these algorithms is the poor data efficiency, which means a lot of data has to be gathered for convergence. This makes it hard to apply them to real world systems due to wear-out and time consumption.

In contrast, model-based algorithms, which learn a representation of the system dynamics, need less data to achieve a certain task (i.e., equilibrium transition) when applicable. These algorithms try to automate the classical process of obtaining a model of the system dynamics by first principles, identification of physical parameters, and controller design. In the context of reinforcement learning, a priori knowledge about the system is assumed to be not available or very limited. This makes the control problem unnecessarily hard from a control theoretic point of view, but reduces model bias and empowers generalization of the learning algorithm.

In this work we combine neural network dynamics models with trajectory planning to form a model-based reinforcement learning algorithm. This approach is similar to Yam-

aguchi and Atkeson (2016), where a stochastic trajectory planning is performed. Other works use Locally Weighted Projection Regression (LWPR) (Mitrovic et al., 2008) or Gaussian Process (GP) (Bechtel et al., 2019; Lee et al., 2017) for model learning.

This paper is structured as follows: In Sec. 2 the basic concepts of model-based reinforcement learning are presented. In Sec. 3 it is shown, how a neural network can be used to learn a dynamics model. In Sec. 4 the iterative Linear-Quadratic Regulator (iLQR) is described, which is used for trajectory planning. Finally we combine the two concepts to form a model-based RL algorithm in Sec. 5 and present results in simulation and on a test bench in Sec. 6.

## 2. PRELIMINARIES

Reinforcement learning (RL) deals with general decision making problems similar to optimal control. In reinforcement learning, an agent interacts with a system by applying a control  $\mathbf{u}_k \in \mathbb{U}_k$  to it. The system (1) transits from a state  $\mathbf{x}_k \in \mathbb{X}_k$  at time step  $k$  to a successor state  $\mathbf{x}_{k+1} \in \mathbb{X}_{k+1}$  at time step  $k+1$  and emits a cost signal  $c_k \in \mathbb{R}$ , which holds information about the fulfillment of the given task.

In contrast to the classical RL setting, the system dynamics

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \quad (1)$$

consists of a deterministic sampled-data system with underlying continuous dynamics and the incremental cost function

$$c_k : \mathbb{X}_k \times \mathbb{U}_k \mapsto \mathbb{R} \quad (2)$$

---

<sup>1</sup> Model-free in the sense, that the system dynamics are not explicitly represented.

is known.

The agent's goal is to minimize the total cost

$$J = \sum_{k=0}^{N-1} c_k(\mathbf{x}_k, \mathbf{u}_k) + c_N(\mathbf{x}_N) \quad (3)$$

over a given time-horizon  $N$ , starting from an initial state  $\mathbf{x}_0$  by choosing appropriate controls. If the system transits to a terminal state  $\mathbf{x}_{k+1} \in \mathbb{X}^-$  or  $k = N$ , the episode terminates, and the system is reset to  $\mathbf{x}_0$ . At every time step the transition tuple  $(\mathbf{x}, \mathbf{u}, \mathbf{x}')$ , where  $\mathbf{x}' = \mathbf{f}(\mathbf{x}, \mathbf{u})$ , is added to a data set  $\mathcal{D} = \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}_{i=0}^T$ , where  $T$  is the total experience of the agent.

### 2.1 Problem statement

This leads to a classical formulation of the optimal control problem:

$$\min_{\mathbf{u}_0, \dots, \mathbf{u}_{N-1}} J(\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{u}_{N-1}) \quad (4a)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \forall k \in [0, N-1], \quad (4b)$$

$$\mathbf{u}_k \in \mathbb{U}_k, \mathbf{x}_k \in \mathbb{X}_k, \quad \text{giv. } \mathbf{x}_0. \quad (4c)$$

### 2.2 Transition problems

It is assumed, that  $\mathbf{x}_0$  is a stable equilibrium of (1) and the goal is to drive the system into a goal state  $\mathbf{x}^*$  in  $N$  time steps. An equilibrium  $\mathbf{x}_e$  of system (1) is defined as:

$$\mathbf{x}_e = \mathbf{f}(\mathbf{x}_e, \mathbf{u}_e), \quad (5)$$

where  $\mathbf{u}_e \in \mathbb{U}$  is an admissible control, such that (5) is satisfied. To transcribe the transition problem to an optimal control problem (4), the cost function has to be designed, such that the cost in the equilibrium vanishes:

$$c_k(\mathbf{x}_e, \mathbf{u}_e) \stackrel{\text{!}}{=} 0, \quad c_N(\mathbf{x}_e) \stackrel{\text{!}}{=} 0. \quad (6)$$

A quadratic cost function is thus a suitable choice.

### 2.3 Model-based RL

In the context of RL, the dynamics  $\mathbf{f}$  are initially unknown and problem (4) can't be solved. Therefore, in model-based RL the dynamics are approximated by a parameterized model  $\mathbf{f}_\theta$  using  $\mathcal{D}$ . With the learned dynamics model an approximate optimal control problem can be solved:

$$\min_{\mathbf{u}_0, \dots, \mathbf{u}_{N-1}} J(\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{u}_{N-1}) \quad (7a)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_\theta(\mathbf{x}_k, \mathbf{u}_k), \quad \forall k \in [0, N-1] \quad (7b)$$

$$\mathbf{u}_k \in \mathbb{U}_k, \mathbf{x}_k \in \mathbb{X}_k, \quad \text{giv. } \mathbf{x}_0 \quad (7c)$$

There are numerous ways to solve problem (7). In *policy search* the global solution is approximated by optimizing a parameterized control law (policy)  $\pi(\mathbf{x}; \phi)$ , see e.g. Deisenroth and Rasmussen (2011). Nagabandi et al. (2018) and Lee et al. (2017); Mitrovic et al. (2008) use a model predictive control scheme, where (7) is solved over a short horizon  $\leq N$  at every time step.

The optimized control law is applied to system (1) and the resulting trajectory  $\tau = \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}_{i=0}^N$  is added to  $\mathcal{D}$ . With the aggregated data set, the model  $\mathbf{f}_\theta$  can be improved and the process repeats until a sufficiently good control law is found, as summarized in Alg. 1.

---

### Algorithm 1 Model-based RL

---

**Require:** initial model  $\mathbf{f}_\theta$ , data set  $\mathcal{D}$   
 initialize  $\mathcal{D}$  by random sampling  
**while** not converged **do**  
     1) train  $\mathbf{f}_\theta$  using  $\mathcal{D}$   
     2) solve the optimal control problem (7)  
     3) rollout the resulting control law on system (1)  
     4) add the resulting trajectory  $\tau$  to the data set  $\mathcal{D}$   
**end while**

---

## 3. LEARNING A NEURAL NETWORK DYNAMICS MODEL

To approximate the dynamics, the neural network proposed by Nagabandi et al. (2018) is used.

The dynamics  $\mathbf{f}$  are split into two terms:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{f}_d(\mathbf{x}_k, \mathbf{u}_k). \quad (8)$$

To approximate the difference model  $\mathbf{f}_d$ , a feedforward neural network  $\mathbf{f}_{d,\theta}$  is used:

$$\mathbf{f}_{d,\theta}(\mathbf{x}_k, \mathbf{u}_k) \stackrel{\text{!}}{\approx} \mathbf{f}_d(\mathbf{x}_k, \mathbf{u}_k). \quad (9)$$

The network is composed of two hidden layers and uses ReLU activation functions:

$$\mathbf{h}^{[1]} = \text{ReLU} \left( \mathbf{W}^{[0]} (\mathbf{x}_k^T, \mathbf{u}_k^T)^T + \mathbf{b}^{[0]} \right), \quad (10a)$$

$$\mathbf{h}^{[2]} = \text{ReLU} \left( \mathbf{W}^{[1]} \mathbf{h}^{[1]} + \mathbf{b}^{[1]} \right), \quad (10b)$$

$$\mathbf{f}_{d,\theta} = \mathbf{W}^{[2]} \mathbf{h}^{[2]} + \mathbf{b}^{[2]}. \quad (10c)$$

Hereby  $\mathbf{h}^{[1]}$  and  $\mathbf{h}^{[2]}$  denote the intermediate results of the internal network layers.

The networks weight matrices and bias vectors together form the parameter vector  $\theta$ :

$$\theta := \{\mathbf{W}^{[0]}, \mathbf{b}^{[0]}; \mathbf{W}^{[1]}, \mathbf{b}^{[1]}; \mathbf{W}^{[2]}, \mathbf{b}^{[2]}\}.$$

To learn an approximate dynamics model, the mean-squared error (MSE) of the difference model derived from (8) and (9) is minimized:

$$\min_{\theta} \frac{1}{T} \sum_{i=0}^T \left\| \left( \mathbf{x}'^{(i)} - \mathbf{x}^{(i)} \right) - \mathbf{f}_{d,\theta} \left( \mathbf{x}^{(i)}, \mathbf{u}^{(i)} \right) \right\|_2^2, \quad (11)$$

where the superscript  $(i)$  labels data points in the  $i$ -th tuple of  $\mathcal{D}$ . Due to a possibly different range of values in the state dimensions the error residuals are not weighted equally. As proposed in Nagabandi et al. (2018), we use data standardization to preprocess the training data.

### 3.1 Data Standardization

The data standardization is performed by computing the standard deviation  $\sigma$  and mean  $\mu$  vectors of  $\Delta \mathbf{x}_k := \mathbf{x}_{k+1} - \mathbf{x}_k$ ,  $\mathbf{x}_k$  and  $\mathbf{u}_k$  over the whole data set  $\mathcal{D}$ :

$$\begin{aligned} \mu_{\Delta \mathbf{x}} &:= \mathbb{E}_{\mathcal{D}}[\Delta \mathbf{x}] & \sigma_{\Delta \mathbf{x}} &:= \sqrt{\text{Var}_{\mathcal{D}}[\Delta \mathbf{x}]} \\ \mu_{\mathbf{x}} &:= \mathbb{E}_{\mathcal{D}}[\mathbf{x}] & \sigma_{\mathbf{x}} &:= \sqrt{\text{Var}_{\mathcal{D}}[\mathbf{x}]} \\ \mu_{\mathbf{u}} &:= \mathbb{E}_{\mathcal{D}}[\mathbf{u}] & \sigma_{\mathbf{u}} &:= \sqrt{\text{Var}_{\mathcal{D}}[\mathbf{u}]} \end{aligned}$$

and applying the following transformation to the data:

$$\begin{aligned} \Delta \hat{\mathbf{x}}_k &:= \text{diag}(\sigma_{\Delta \mathbf{x}}^{-1})(\Delta \mathbf{x} - \mu_{\Delta \mathbf{x}}), \\ \hat{\mathbf{x}}_k &:= \text{diag}(\sigma_{\mathbf{x}}^{-1})(\mathbf{x} - \mu_{\mathbf{x}}), \\ \hat{\mathbf{u}}_k &:= \text{diag}(\sigma_{\mathbf{u}}^{-1})(\mathbf{u} - \mu_{\mathbf{u}}). \end{aligned}$$

The neural network  $\hat{f}_{d,\theta}$  is trained in the standardized coordinates:

$$\min_{\theta} \frac{1}{T} \sum_{i=0}^T \left\| \Delta \hat{x}^{(i)} - \hat{f}_{d,\theta}(\hat{x}^{(i)}, \hat{u}^{(i)}) \right\|_2^2. \quad (13)$$

A retransformation to the original coordinates leads to the final model:

$$f_{d,\theta} := \text{diag}(\sigma_{\Delta x}) \hat{f}_{d,\theta}(\hat{x}_k, \hat{u}_k) + \mu_{\Delta x}. \quad (14)$$

### 3.2 Data acquisition heuristic

When the agent interacts with the system and collects data points  $(\mathbf{x}, \mathbf{u}, \mathbf{x}')$ , these data points are not sampled equally distributed in the  $\mathbb{X} \times \mathbb{U} \times \mathbb{X}$  space. At the beginning of the learning process they are obviously dense near the initial state  $(\mathbf{x}_0, \star, \star)$ . When a trajectory is rolled out, that does not lie in the distribution of the data set  $\mathcal{D}$ , these new data points are sparse and therefore treated as outliers in the training process. To make sure the neural network model also captures the information of these data points the density of the data in  $\mathcal{D}$  has to be bounded. Therefore, a simple heuristic is used, that makes sure a data point is only added to  $\mathcal{D}$  if it is not already captured by the model. Namely, if the prediction error of the latest sample  $(\mathbf{x}, \mathbf{u}, \mathbf{x}')$  is higher than a specified bound  $\epsilon$ , the sample is added to  $\mathcal{D}$ :

$$\left\| \Delta \mathbf{x}^{(i)} - f_{d,\theta}(\mathbf{x}^{(i)}, \mathbf{u}^{(i)}) \right\|_2^2 > \epsilon. \quad (15)$$

### 3.3 Exploiting the Model Structure of Mechanical Systems

When dealing with mechanical systems, some properties of the system dynamics can be exploited to ease the learning process.

When adopting a Lagrangian perspective, the state of a mechanical system  $\mathbf{x}$  is composed of its generalized coordinates and the corresponding velocities  $\mathbf{x} := (\mathbf{q}^T \dot{\mathbf{q}}^T)^T$ . The continuous time state space model is then given by:

$$\dot{\mathbf{x}} = \begin{pmatrix} \dot{\mathbf{q}} \\ \mathbf{g}(\dot{\mathbf{q}}, \mathbf{q}, \mathbf{u}) \end{pmatrix}, \quad (16)$$

where  $\dot{\mathbf{q}} = \mathbf{g}$  are the relations stemming from the Lagrangian equations of motion (which are considered as a first principle). For a sufficiently small sample-time  $\Delta t$  an approximation of the discrete time dynamics is given by:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t \begin{pmatrix} \dot{\mathbf{q}}_k \\ \mathbf{g}(\dot{\mathbf{q}}_k, \mathbf{q}_k, \mathbf{u}_k) \end{pmatrix}. \quad (17)$$

The upper hyper-row of the system dynamics consists of known definitional equations which can be exactly represented (i. e. hardcoded) easily by a neural network. Thus only the mapping  $\mathbf{g}$ , which models the acceleration dynamics has to be learned.

If the system has angles as generalized coordinates, the system dynamics are typically not dependent on the absolute value of the angle, but depend on its sine and/or cosine components. Thus, instead of the state component  $x_i$  the two pseudo-state-components  $(\sin x_i, \cos x_i)$  can be considered. If this is done for all angle-coordinates, one obtains the pseudo-state-vector  $\mathbf{o}$ . In the learning process,  $\mathbf{o}$  is used instead of  $\mathbf{x}$  as the input of the neural network, but its output remains  $\Delta \mathbf{x}$ .

## 4. TRAJECTORY PLANNING

For solving the approximate optimal control problem (7) with (14) as  $f_{\theta}$ , trajectory planning techniques can be used. Instead of only solving for the open loop controls  $\mathbf{u}_0^*, \dots, \mathbf{u}_{N-1}^*$ , the optimization performed in this work leads to a time-varying feedback-controller with feedforward:

$$\mathbf{u}_k = \mathbf{K}_k(\mathbf{x}_k - \mathbf{x}_k^*) + \mathbf{u}_k^*, \forall k \in [0, N-1], \quad (18)$$

that stabilizes the system around the optimized state trajectory  $\mathbf{x}_0^*, \dots, \mathbf{x}_{N-1}^*$ .

### 4.1 Iterative Linear-Quadratic Regulator (iLQR)

The iLQR is a second-order gradient method based on Differential Dynamic Programming (DDP) (Mayne, 1966), that solves problem (4) in an iterative fashion. In this work, the control-limited variant of iLQR from Tassa et al. (2014) is used. In this section, a brief overview of the algorithm is presented, for details see Li and Todorov (2004); Tassa et al. (2014). Neglecting a regularization heuristic and the performed line-search, the basic algorithm can be boiled down to three steps as shown in Alg. 2.

---

#### Algorithm 2 iLQR (simplified)

---

**Require:**  $\mathbf{f}$ ,  $c_k$ ,  $c_N$

choose initial controls  $\mathbf{U}^{[0]} := \{\mathbf{u}_0^{[0]}, \dots, \mathbf{u}_{N-1}^{[0]}\}$

$\mathbf{X}^{[0]} := \{\mathbf{x}_0^{[0]}, \dots, \mathbf{x}_N^{[0]}\}$  rollout by applying  $\mathbf{U}^{[0]}$  to  $\mathbf{f}$   
**for**  $j = [0, \dots, j_{\max}]$  **do**

1) approximate (4) around  $\tau^{[j]} := \{\mathbf{X}^{[j]}, \mathbf{U}^{[j]}\}$

2) solve the resulting time-varying LQR problem

3)  $\tau^{[j+1]}$ : rollout the resulting control law on (1)

**if**  $J(\tau^{[j]}) - J(\tau^{[j+1]}) < \epsilon_J$  **then**

**break**

**end if**

**end for**

**return**  $\tau^* = \{\mathbf{X}^*, \mathbf{U}^*\} := \tau^{[j+1]}$

---

#### 1) Local approximation of the optimal control problem

Given the current trajectory  $\{\tilde{\mathbf{X}}, \tilde{\mathbf{U}}\} =: \tau^{[j]}$  a Taylor expansion of the system dynamics and cost is performed around it. In contrast to DDP, iLQR only uses a first order expansion of the dynamics:

$$\tilde{\mathbf{x}}_{k+1} \approx \mathbf{A}_k \tilde{\mathbf{x}}_k + \mathbf{B}_k \tilde{\mathbf{u}}_k, \quad (19)$$

with

$$\tilde{\mathbf{x}}_k = \mathbf{x}_k - \bar{\mathbf{x}}_k, \quad \tilde{\mathbf{u}}_k = \mathbf{u}_k - \bar{\mathbf{u}}_k, \forall k. \quad (20)$$

The cost approximation is of second order:

$$\tilde{c}_k \approx (\tilde{\mathbf{x}}_k^T \tilde{\mathbf{u}}_k^T) \left[ \begin{pmatrix} \mathbf{c}_{\mathbf{x},k} \\ \mathbf{c}_{\mathbf{u},k} \end{pmatrix} + \frac{1}{2} \begin{pmatrix} \mathbf{C}_{\mathbf{x}\mathbf{x},k} & \mathbf{C}_{\mathbf{u}\mathbf{x},k} \\ \mathbf{C}_{\mathbf{x}\mathbf{u},k} & \mathbf{C}_{\mathbf{u}\mathbf{u},k} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{x}}_k \\ \tilde{\mathbf{u}}_k \end{pmatrix} \right], \quad (21a)$$

$$\tilde{c}_N \approx \tilde{\mathbf{x}}_N^T \mathbf{c}_{\mathbf{x},N} + \frac{1}{2} \tilde{\mathbf{x}}_N^T \mathbf{C}_{\mathbf{x}\mathbf{x},N} \tilde{\mathbf{x}}_N. \quad (21b)$$

This leads to the optimization objective:

$$\tilde{J} = \sum_{k=0}^{N-1} \tilde{c}_k(\tilde{\mathbf{x}}_k, \tilde{\mathbf{u}}_k) + \tilde{c}_N(\tilde{\mathbf{x}}_N). \quad (22)$$

2) *Solving the local approximation of the optimal control problem* From the linear-quadratic cost function (22) and the linear time-varying dynamics (19) an optimal control problem can be formulated:

$$\min_{\tilde{\mathbf{u}}_0, \dots, \tilde{\mathbf{u}}_{N-1}} \tilde{J}(\tilde{\mathbf{x}}_0, \tilde{\mathbf{u}}_0, \dots, \tilde{\mathbf{u}}_{N-1}), \quad (23a)$$

$$\text{s.t. } \tilde{\mathbf{x}}_{k+1} = \mathbf{A}_k \tilde{\mathbf{x}}_k + \mathbf{B}_k \tilde{\mathbf{u}}_k, \forall k \in [0, N-1], \quad (23b)$$

$$\mathbf{u}_k \in \mathbb{U}_k, \mathbf{x}_k \in \mathbb{X}_k, \text{ giv. } \mathbf{x}_0. \quad (23c)$$

Problem (23) can be solved by dynamic programming (Bertsekas, 2005):

$$\tilde{V}_N(\tilde{\mathbf{x}}_N) = \tilde{c}_N(\tilde{\mathbf{x}}_N), \quad (24a)$$

$$\tilde{V}_k(\tilde{\mathbf{x}}_k) = \min_{\tilde{\mathbf{u}}_k} \underbrace{\left[ \tilde{c}_k(\tilde{\mathbf{x}}_k, \tilde{\mathbf{u}}_k) + \tilde{V}_{k+1}(\tilde{\mathbf{x}}_{k+1}) \right]}_{=: \tilde{Q}_k(\tilde{\mathbf{x}}_k, \tilde{\mathbf{u}}_k)}. \quad (24b)$$

The optimal controller can be derived in analogy to the discrete-time LQR backwards in time starting from  $k = N-1$  with

$$\mathbf{V}_{\mathbf{x}\mathbf{x},N} = \mathbf{C}_{\mathbf{x}\mathbf{x},N} \quad \mathbf{v}_{\mathbf{x},N} = \mathbf{c}_{\mathbf{x},N}. \quad (25)$$

The Hamiltonian  $\tilde{Q}_k$  is linear-quadratic:

$$\tilde{Q}_k = \begin{pmatrix} \tilde{\mathbf{x}}_k^T & \tilde{\mathbf{u}}_k^T \end{pmatrix} \left[ \begin{pmatrix} \mathbf{q}_{\mathbf{x},k} \\ \mathbf{q}_{\mathbf{u},k} \end{pmatrix} + \frac{1}{2} \begin{pmatrix} \mathbf{Q}_{\mathbf{x}\mathbf{x},k} & \mathbf{Q}_{\mathbf{x}\mathbf{u},k} \\ \mathbf{Q}_{\mathbf{x}\mathbf{u},k} & \mathbf{Q}_{\mathbf{u}\mathbf{u},k} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{x}}_k \\ \tilde{\mathbf{u}}_k \end{pmatrix} \right], \quad (26)$$

with

$$\mathbf{Q}_{\mathbf{x}\mathbf{x},k} = \mathbf{C}_{\mathbf{x}\mathbf{x},k} + \mathbf{A}_k^T \mathbf{V}_{\mathbf{x}\mathbf{x},k+1} \mathbf{A}_k, \quad (27a)$$

$$\mathbf{Q}_{\mathbf{x}\mathbf{u},k} = \mathbf{Q}_{\mathbf{u}\mathbf{x},k}^T = \mathbf{C}_{\mathbf{x}\mathbf{u},k} + \mathbf{A}_k^T \mathbf{V}_{\mathbf{x}\mathbf{x},k+1} \mathbf{B}_k, \quad (27b)$$

$$\mathbf{Q}_{\mathbf{u}\mathbf{u},k} = \mathbf{C}_{\mathbf{u}\mathbf{u},k} + \mathbf{B}_k^T \mathbf{V}_{\mathbf{x}\mathbf{x},k+1} \mathbf{B}_k, \quad (27c)$$

$$\mathbf{q}_{\mathbf{x},k} = \mathbf{c}_{\mathbf{x},k} + \mathbf{A}_k^T \mathbf{v}_{\mathbf{x},k+1}, \quad (27d)$$

$$\mathbf{q}_{\mathbf{u},k} = \mathbf{c}_{\mathbf{u},k} + \mathbf{B}_k^T \mathbf{v}_{\mathbf{x},k+1}. \quad (27e)$$

Solving for the optimal control law leads to the following feedback law with feedforward:

$$\tilde{\mathbf{u}}_k^* = \mathbf{K}_k \tilde{\mathbf{x}}_k + \mathbf{k}_k, \forall k \in [0, N-1]. \quad (28)$$

The feedback matrix and feedforward are given by:

$$\mathbf{K}_k = -\mathbf{Q}_{\mathbf{u}\mathbf{u},k}^{-1} \mathbf{Q}_{\mathbf{u}\mathbf{x},k}, \quad \mathbf{k}_k = -\mathbf{Q}_{\mathbf{u}\mathbf{u},k}^{-1} \mathbf{q}_{\mathbf{u},k}. \quad (29)$$

Substituting the optimal control law (28) into the Hamiltonian (26) results in the cost-to-go:

$$\tilde{V}_k(\tilde{\mathbf{x}}_k) = \tilde{\mathbf{x}}_k^T \mathbf{v}_{\mathbf{x},k} + \frac{1}{2} \tilde{\mathbf{x}}_k^T \mathbf{V}_{\mathbf{x}\mathbf{x},k} \tilde{\mathbf{x}}_k + \text{const.}, \quad (30)$$

with

$$\mathbf{V}_{\mathbf{x}\mathbf{x},k} = \mathbf{Q}_{\mathbf{x}\mathbf{x},k} + \mathbf{K}_k^T (\mathbf{Q}_{\mathbf{u}\mathbf{x},k} + \mathbf{Q}_{\mathbf{u}\mathbf{u},k} \mathbf{K}_k) + \mathbf{Q}_{\mathbf{x}\mathbf{u},k} \mathbf{K}_k, \quad (31a)$$

$$\mathbf{v}_{\mathbf{x},k} = \mathbf{q}_{\mathbf{x},k} + \mathbf{K}_k^T (\mathbf{q}_{\mathbf{u},k} + \mathbf{Q}_{\mathbf{u}\mathbf{u},k} \mathbf{k}_k) + \mathbf{Q}_{\mathbf{x}\mathbf{u},k} \mathbf{k}_k. \quad (31b)$$

3) *Rollout of the optimal control law* The optimal control law (28) for problem (23) can be solved for  $\mathbf{u}_k$  resulting in (18):

$$\mathbf{u}_k = \underbrace{\mathbf{K}_k (\mathbf{x}_k - \bar{\mathbf{x}}_k)}_{=: \mathbf{x}_k^*} + \underbrace{\mathbf{k}_k + \bar{\mathbf{u}}_k}_{=: \mathbf{u}_k^*}, \forall k. \quad (32)$$

Applying (18) to (1) leads to a new trajectory starting from  $\mathbf{x}_0 := \bar{\mathbf{x}}_0 \forall k \in [0, N-1]$ :

$$\mathbf{u}_k = \mathbf{K}_k (\mathbf{x}_k - \bar{\mathbf{x}}_k) + \mathbf{k}_k + \bar{\mathbf{u}}_k, \quad (33a)$$

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k). \quad (33b)$$

The new trajectory  $\tau^{[j+1]} := \{\mathbf{X}, \mathbf{U}\}$  is used for the next iteration.

## 4.2 Computing the final control law

Dealing with transition problems, it is assumed that  $\mathbf{x}_k = \mathbf{x}_e$  and  $\mathbf{u}_k = \mathbf{u}_e$  for  $k > N$ . To stabilize the system in  $\mathbf{x}_e$ , a standard method is to design a LQ controller

$$\mathbf{u}_k = \mathbf{K}_{LQR} (\mathbf{x}_k - \mathbf{x}_e) + \mathbf{u}_e, \quad k > N. \quad (34)$$

Under the given assumptions  $\mathbf{K}_N = \mathbf{K}_{LQR}$  can be achieved, i.e. the time-varying policy (18) transits smoothly to the constant LQR policy. To accomplish this, the DARE<sup>2</sup> of the LQR problem for  $\mathbf{A}_e = \partial_{\mathbf{x}} \mathbf{f}_{\theta}(\mathbf{x}_e, \mathbf{u}_e)$  and  $\mathbf{B}_e = \partial_{\mathbf{u}} \mathbf{f}_{\theta}(\mathbf{x}_e, \mathbf{u}_e)$  is solved, with  $\mathbf{R} := \mathbf{C}_{\mathbf{u}\mathbf{u},N}$ ,  $\mathbf{Q} := \mathbf{C}_{\mathbf{x}\mathbf{x},N}$ ,  $\mathbf{N} := \mathbf{C}_{\mathbf{x}\mathbf{u},N}$ . The resulting matrix  $\mathbf{P}$ , which is the solution of the DARE, is used to initialize the backward pass in step 2):

$$\mathbf{V}_{\mathbf{x}\mathbf{x},N} := \mathbf{P}, \quad \mathbf{v}_{\mathbf{x},N} := \mathbf{0}. \quad (35)$$

However, while the trajectory optimization yields  $\mathbf{x}_N^* \approx \mathbf{x}_e$  it does not guarantee  $\mathbf{x}_N^* \equiv \mathbf{x}_e$ . To stabilize the system in the desired state  $\mathbf{x}_e$ , for  $k > N$ , the following control law is applied instead of (18):

$$\mathbf{u}_k = \mathbf{K}_N (\mathbf{x}_k - \mathbf{x}_e) + \mathbf{u}_e. \quad (36)$$

## 5. COMBINING MODEL APPROXIMATION AND TRAJECTORY OPTIMIZATION

The proposed neural network model in Sec. 3 for approximating the dynamics (1) and the trajectory planning approach in Sec. 4 can be combined to form the model-based reinforcement learning algorithm Alg. 3, which is a special case of Alg. 1.

---

**Algorithm 3** Model-based RL with offline trajectory planning

---

**Require:** model  $\mathbf{f}_{d,\theta}$ , data set  $\mathcal{D}$

initialize  $\mathcal{D}$  by random sampling

**while** not converged **do**

1) train  $\mathbf{f}_{d,\theta}$  using (13) on the data set  $\mathcal{D}$

2) solve the optimal control problem (7) with iLQR

3) apply the resulting policy (18) on system (1)

4) add the resulting trajectory  $\tau$  to the data set  $\mathcal{D}$

**end while**

---

## 6. RESULTS

### 6.1 Simulation Results<sup>3</sup>

As a first simulation experiment, Alg. 3 is applied to the cart-pole system (Fig. 1).

The state is composed of the position  $q_0$  and velocity  $\dot{q}_0$  of the cart, as well as the angle  $q_1$  and angular velocity  $\dot{q}_1$  of the pole:

$$\mathbf{x} := (q_0 \quad q_1 \quad \dot{q}_0 \quad \dot{q}_1). \quad (37)$$

The system is controlled by the force  $F_0$ , that acts on the cart. The continuous dynamics in an input-output linearized form, where the control input  $\mathbf{u} = (u_1)$  is the acceleration of the cart  $\ddot{q}_0$ , are given by:

<sup>2</sup> Discrete Algebraic Riccati Equation

<sup>3</sup> Code available at: <https://github.com/TUD-RST/pygent>

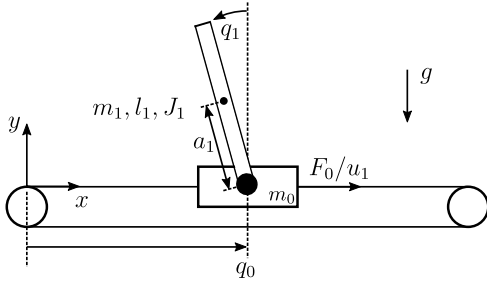


Fig. 1. Cart-pole system

$$\dot{x}_1 = x_3, \quad \dot{x}_2 = x_4, \quad \dot{x}_3 = u_1, \quad (38a)$$

$$\dot{x}_4 = a_1 m_1 (g \sin(x_2) + u_1 \cos(x_2)) - \frac{d_1}{J_1 + a_1^2 m_1} x_4, \quad (38b)$$

where  $a_1$  is the position of the center of gravity of the pole,  $m_1$  is the mass of the pole,  $J_1$  is the angular momentum of the pole,  $d_1$  is a damping coefficient, and  $g$  is gravity. The control objective is to swing-up and balance the pole in the upward unstable equilibrium  $\mathbf{x}^* = (0 \ 0 \ 0 \ 0)^T$  starting from  $\mathbf{x}_0 = (0 \ \pi \ 0 \ 0)^T$ .

*Controller performance* The learning curve of Alg.3 applied to the cart-pole system is shown in Fig.2. For the first five episodes, the agent samples controls from a uniform distribution to initialize the data set  $\mathcal{D}$ . At the 14-th episode (about 2500 time steps), the swing-up and stabilization of the pole in the upward position is accomplished for the first time. After 19 episodes (about 3600 time steps), the learning based controller asymptotically matches the performance of a controller, that has access to the true system dynamics. Compared to DDPG (Lillicrap et al., 2016), a popular model-free RL method, our approach is about 70 times more data-efficient on the swing-up task and has a better final performance, cf. (Pritzkoleit, 2019).

*Prediction performance of the approximated model* To examine the prediction quality of the neural network model, we consider the error along the (optimal) trajectory<sup>4</sup>, i. e.,  $\mathbf{e}(k) := \hat{\mathbf{f}}_{d,\theta}(\mathbf{x}_k^*, \mathbf{u}_k^*) - \mathbf{f}_d(\mathbf{x}_k^*, \mathbf{u}_k^*)$  at different training stages (episodes). Because the first two components of  $\mathbf{f}$  are definitional equations (cf. (17) and (38a)), only the last two components – which correspond to the acceleration of the cart and the pendulum angle – have to be considered. From Fig. 3 we see good convergence with increasing episode number. Furthermore, we notice that, even after about 49 episodes the approximated model still differs perceptibly but that does not impair the algorithm from successfully performing the swing-up task.

## 6.2 Experimental results

To further investigate the performance of Alg.3, it is applied on a real cart-triple-pole test bench system. The goal of the experiment is to side-step the downward hanging poles from  $q_0 = 0$  to  $q_0 = 0.8$ , while the poles must rest at the end of the transition. To initialize the data set  $\mathcal{D}$  we first planned a twice differentiable

<sup>4</sup> The optimal trajectory was calculated with the true dynamics using iLQR.

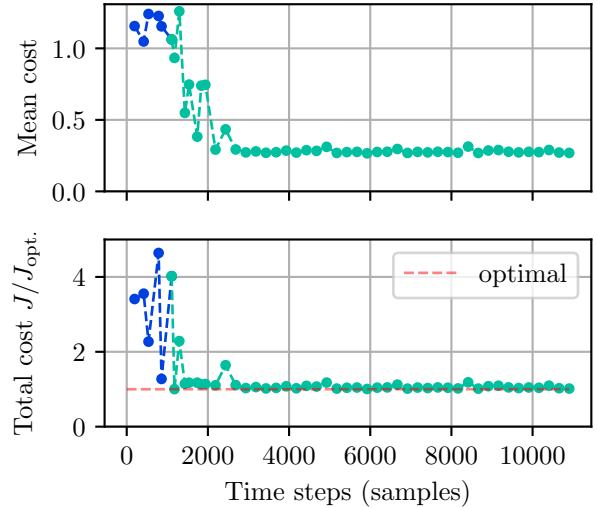


Fig. 2. Learning Process: After around 3600 time steps (72 s interaction time), the learning based algorithm converges to the performance of the controller with access to the true dynamics (red dashed). The total cost is normalized. The initialization episodes (blue) have a high mean cost, but a varying total cost, due to early termination of these episodes.

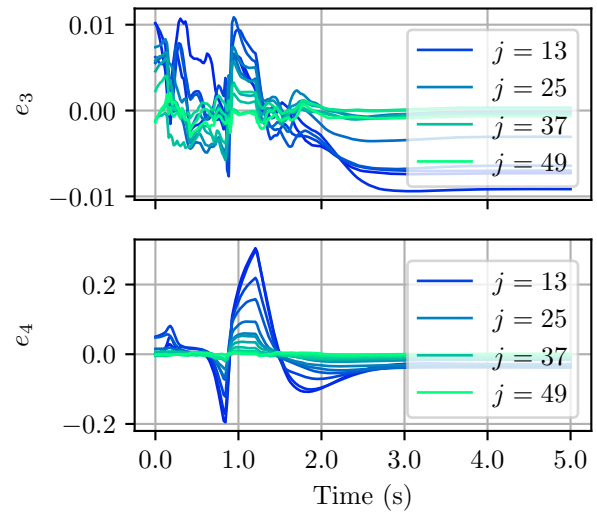


Fig. 3. Approximation error: The values of the neural network model  $\hat{\mathbf{f}}_{d,\theta}$  converge to the real model  $\mathbf{f}_d$  with increasing episode number  $j$ .

trajectory for  $q_0$ , which transitions smoothly between several randomly chosen positions. Based on the planned trajectory, a feedforward control for the input  $u_1 = \ddot{q}_0$  is derived. After three<sup>5</sup> iterations of Alg.3 the transition can be accomplished (see Fig. 4 and the respective video at <https://github.com/TUD-RST/pygent>).

<sup>5</sup> After only one iteration, the algorithm drives the system to 0.8, but the poles are not perfectly at rest.

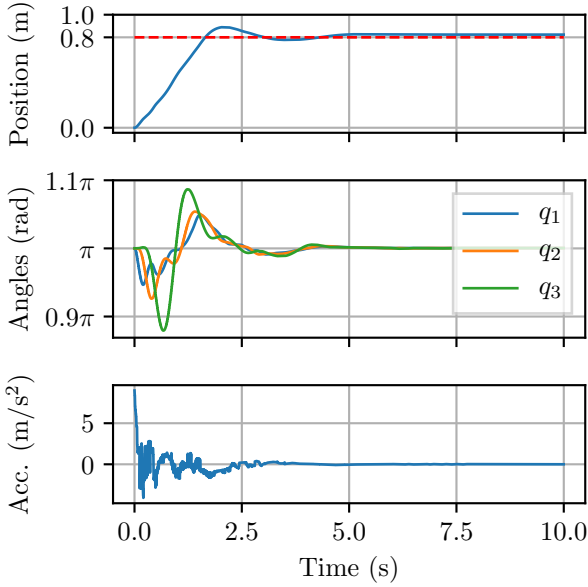


Fig. 4. Cart-triple-pole system transition from  $q_0 = 0.0$  to  $q_0 = 0.8$  with downward hanging poles.

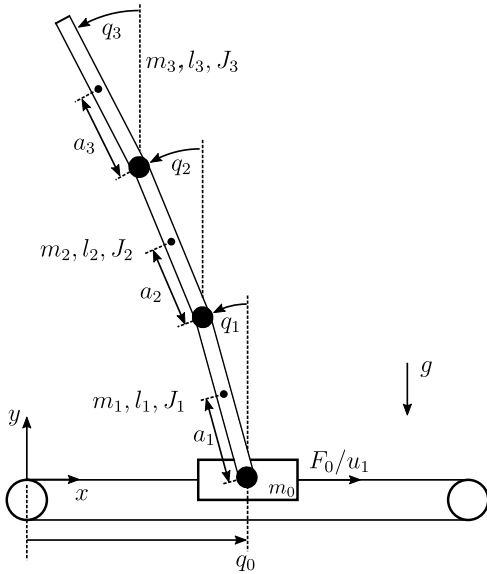


Fig. 5. Cart-triple-pole system.

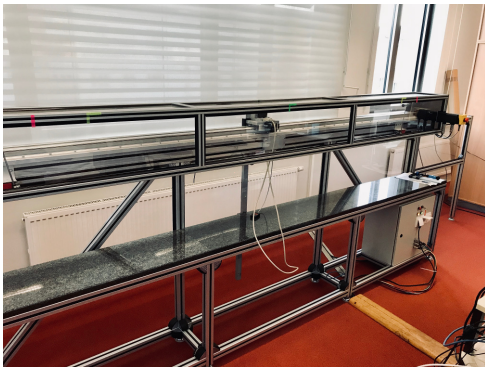


Fig. 6. Cart-triple-pole test bench at our institute.

## 7. CONCLUSION AND OUTLOOK

In this work we demonstrated, that the combination of neural network dynamics models and offline trajectory planning is suitable for solving certain transition problems on the single and cart-triple-pole systems without prior knowledge of the dynamics through a reinforcement learning scheme. Thereby, the *data efficiency* of the swing-up task of the cart-pole system is significantly higher, compared to model-free methods and matches that of other model-based algorithms (Chua et al., 2018, Fig. 3). One advantage of our approach solving the swing-up task of the cart-pole system over PILCO (Deisenroth and Rasmussen, 2011) and Deep PILCO (Gal et al., 2016) is the lower *computational effort*. The authors in Gal et al. (2016) report 31.05 (PILCO)/8.78 (Deep PILCO) minutes of average training time per episode, where the presented approach only requires 4.13 minutes<sup>6</sup> per episode and scales linear with the data. In contrast to using a GP dynamics model, the applicability of the neural network model is thus not as limited by the size of the data set  $\mathcal{D}$ .

An open research question is, if Alg. 3 is also suited for learning the swing-up of the cart-triple-pole system. Here, the remaining model error, as pointed out in Sec. 6.1 could be more problematic, due to the strong nonlinearity of the system and its chaotic behaviour. The small sample-time further increases the problem of error propagation in multi-step trajectory planning. To tackle this problem, one approach could be to use a multi-step prediction error instead of (11).

One general drawback with the presented deterministic modeling approach is that *the model does not know what it does not know*. To overcome this problem, the prediction uncertainty could be incorporated by using so called deep ensemble models (Lakshminarayanan et al., 2017) which were used in the RL setting in Chua et al. (2018).

In this work we used iLQR as a planning algorithm, which sometimes got stuck in local minima in the early stages of trajectory planning, when using the neural network model of dynamical systems that were not included in this paper. Therefore, other trajectory planning algorithms like collocation could be compared to iLQR in future works.

From a control theoretic point of view it may seem unnecessary effortful to learn the full dynamics, when a reasonable good model can be found by first-principles and system identification. The cart-triple-pole system however is rather difficult to control and an approach that combines model learning with a physical modeling based prior could be beneficial. Kaheman et al. (2019) and Lee et al. (2017), where only the modeling error is learned, are promising approaches.

## ACKNOWLEDGEMENTS

We thank the Center for Information Services and High Performance Computing (ZIH) at TU Dresden for generous allocations of compute resources. We also thank the anonymous reviewers for their valuable comments.

<sup>6</sup> Hardware: 1.3 GHz Dual-Core Intel Core i5 (2013).

## REFERENCES

- Bechtle, S., Rai, A., Lin, Y., Righetti, L., and Meier, F. (2019). Curious iLQR: Resolving uncertainty in model-based RL. *arXiv:1904.06786*.
- Bertsekas, D.P. (2005). *Dynamic Programming and Optimal Control*, volume 1. Athena Scientific, Belmont, MA, 3rd edition.
- Chua, K., Calandra, R., McAllister, R., and Levine, S. (2018). Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In *Advances in Neural Information Processing Systems*, 4754–4765.
- Deisenroth, M. and Rasmussen, C.E. (2011). PILCO: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, 465–472.
- Gal, Y., McAllister, R., and Rasmussen, C.E. (2016). Improving pilco with bayesian neural network dynamics models.
- Kaheman, K., Kaiser, E., Strom, B., Kutz, J.N., and Brunton, S.L. (2019). Learning discrepancy models from experimental data. *arXiv:1909.08574*.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. (2017). Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, 6402–6413.
- Lee, G., Srinivasa, S.S., and Mason, M.T. (2017). GP-iLQG: Data-driven robust optimal control for uncertain nonlinear dynamical systems. *arXiv:1705.05344*.
- Li, W. and Todorov, E. (2004). Iterative linear quadratic regulator design for nonlinear biological movement systems. In *Proceedings of the 1st International Conference on Informatics in Control, Automation and Robotics*, 222–229.
- Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. (2016). Continuous control with deep reinforcement learning. *arXiv: 1509.02971*.
- Mayne, D. (1966). A second-order gradient method for determining optimal trajectories of non-linear discrete-time systems. *International Journal of Control*, 3(1), 85–95.
- Mitrovic, D., Klanke, S., and Vijayakumar, S. (2008). Optimal control with adaptive internal dynamics models. In *Proceedings of the 5th International Conference on Informatics in Control, Automation and Robotics*.
- Nagabandi, A., Kahn, G., Fearing, R.S., and Levine, S. (2018). Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 7559–7566.
- Pritzkolet, M. (2019). *Bestärkendes Lernen zur Steuerung und Regelung nichtlinearer dynamischer Systeme*. Diploma thesis, TU Dresden, Germany. URL <https://github.com/TUD-RST/pygent>.
- Tassa, Y., Mansard, N., and Todorov, E. (2014). Control-limited differential dynamic programming. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, 1168–1175. IEEE.
- Yamaguchi, A. and Atkeson, C.G. (2016). Neural networks and differential dynamic programming for reinforcement learning problems. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, 5434–5441.