

Exercise sheet 4 - Planning and Learning

Please prepare the following exercises for the upcoming tutorial.

Task 1: Path Planning

- a) The original RRT algorithm has been proposed in [LaValle, 1998]. The main idea is to build a tree of feasible trajectories, rooted at the initial state \mathbf{x}_{init} . The algorithm starts with a graph which includes a single vertex. With each iteration a new vertex is created. Therefore, a random point $\mathbf{x}_{\text{rand}} \in \mathcal{X}_{\text{free}}$ from the search space is sampled. The nearest vertex $\mathbf{x}_{\text{nearest}}$ to \mathbf{x}_{rand} from the graph \mathcal{V} is used in order to generate a new point \mathbf{x}_{new} using a steering algorithm. If the connection between $\mathbf{x}_{\text{nearest}}$ and \mathbf{x}_{new} is free from obstacles, the new point \mathbf{x}_{new} is added to the graph \mathcal{V} and an edge $(\mathbf{x}_{\text{new}}, \mathbf{v})$ to the set of edges \mathcal{E} . The algorithm is stopped as soon as a predefined number of iteration steps n is reached. The shortest path from the initial state to the goal region is generated by starting at a node in the goal region and go backwards along the tree until the initial state is reached. In algorithm 1 an outline of the RRT algorithm is given. The search space \mathcal{X} is divided into an obstacle-free space $\mathcal{X}_{\text{free}} \subseteq \mathcal{X}$ and an obstacle space $\mathcal{X}_{\text{obstacle}} \subseteq \mathcal{X}$. A connection between two points is feasible if a straight line between the two do not intersect with space of obstacles $\mathcal{X}_{\text{obstacle}}$. In Figure 1 an example is given, where the algorithm is used to find a trajectory from the initial position $\mathbf{x}_{\text{init}} = [0 \ 0]^T$ to the goal position $\mathbf{x}_{\text{goal}} = [1.2 \ 1.2]^T$. The blue crosses are the randomly generated points \mathbf{x}_{rand} , the black lines represent the tree structure and the red path is the path with lowest cost from \mathbf{x}_{init} to \mathbf{x}_{goal} . The initial and goal states are marked with red circles, where the initial state is located in the bottom left corner.
- b) In [Karaman and Frazzoli, 2011] the RRT* algorithm has been proposed. Both algorithms, RRT and RRT*, are probabilistically complete. Probabilistically completeness means that, if a solution exists, the probability that a solution is produced approaches 1 if sufficient time is spent. The main difference between RRT and RRT* is, that RRT* is not only probabilistically complete but also asymptotically optimal. Hence, the cost of the calculated solution converges to the optimum. In order to achieve this goal, RRT* uses rewiring of the tree in every iteration step. Therefore, a set $\mathcal{X}_{\text{near}}$ from \mathcal{V} with vertices near \mathbf{x}_{new} is generated. The radius defining if a vertex is near another can be dependent of the cardinality of \mathcal{V} . In order to generate an optimal path, an additive cost function is used and all points in $\mathcal{X}_{\text{near}}$ and $\mathbf{x}_{\text{nearest}}$ are checked in order to find the connection with the lowest cost as long as such a connection is obstacle free. The lowest cost connection $\{(\mathbf{x}_{\text{min}}, \mathbf{x}_{\text{new}})\}$ is added to the edge set \mathcal{E} . If there is a feasible connection from

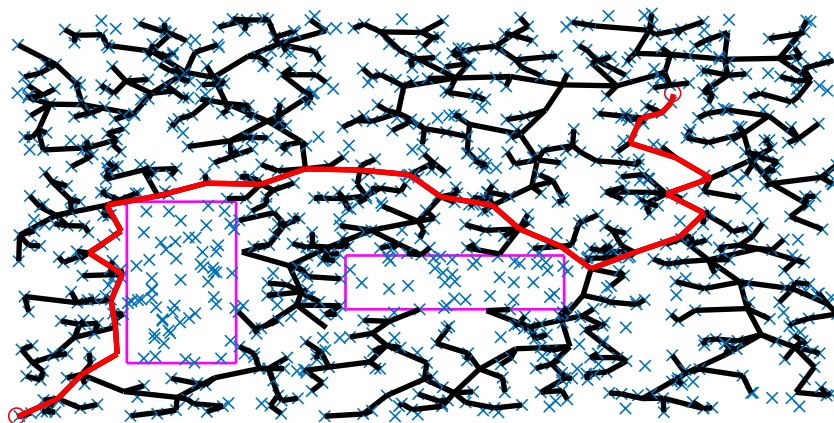


Figure 1 Rapidly-exploring Random Tree example. The branches of the tree are depicted in black and the final trajectory in red. Obstacles are represented in magenta.

Algorithm 1 RRT, adjusted from [Karaman and Frazzoli, 2011]

```

1:  $\mathcal{V} \leftarrow \{\mathbf{x}_{\text{init}}\}; \mathcal{E} \leftarrow \emptyset$ 
2: for  $i = 1, \dots, n$  do
3:    $\mathbf{x}_{\text{rand}} \leftarrow \text{SampleRandom}$ 
4:    $\mathbf{x}_{\text{nearest}} \leftarrow \text{Nearest}(\mathcal{V}, \mathbf{x}_{\text{rand}})$ 
5:    $\mathbf{x}_{\text{new}} \leftarrow \text{Steer}(\mathbf{x}_{\text{nearest}}, \mathbf{x}_{\text{rand}})$ 
6:   if  $\text{ObstacleFree}(\mathbf{x}_{\text{nearest}}, \mathbf{x}_{\text{new}})$  then
7:      $\mathcal{V} \leftarrow \mathcal{V} \cup \{\mathbf{x}_{\text{new}}\}$ 
8:      $\mathcal{E} \leftarrow \mathcal{E} \cup \{(\mathbf{x}_{\text{nearest}}, \mathbf{x}_{\text{new}})\}$ 
9:   end if
10: end for
11: return  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ 

```

\mathbf{x}_{new} to $\mathbf{x}_{\text{near}} \in \mathcal{X}_{\text{near}}$ which produces a lower cost path as the cost to reach \mathbf{x}_{near} at this time, the edge $\{(\mathbf{x}_{\text{parent}}, \mathbf{x}_{\text{near}})\}$, which connects the point with the tree, is deleted and a new edge $\{(\mathbf{x}_{\text{new}}, \mathbf{x}_{\text{near}})\}$ added. In algorithm 2 the RRT* algorithm is presented and in Figure 2 an example is shown. The algorithm is used to find a trajectory from the initial position $\mathbf{x}_{\text{init}} = [0 \ 0]^\top$ to the goal position $\mathbf{x}_{\text{goal}} = [1.2 \ 1.2]^\top$. The green lines are representing the rewiring.

- c) In this section the STOMP algorithm [Kalakrishnan et al., 2011] is reviewed. Here, the algorithm for a 1-dimensional trajectory is derived. For higher dimensions the derivation is similar. Considering a 1-dimensional discretized trajectory given by the vector $\boldsymbol{\theta} \in \mathbb{R}^N$, where N is the number of waypoints. The initial waypoint as well as the goal waypoint are given and kept fixed. The optimization problem is given by

$$\min_{\boldsymbol{\theta}} \mathbb{E} \left[\sum_{i=1}^N q(\hat{\boldsymbol{\theta}}_i) + \frac{1}{2} \hat{\boldsymbol{\theta}}^\top \mathbf{R} \hat{\boldsymbol{\theta}} \right], \quad (1)$$

where $\hat{\boldsymbol{\theta}} = \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma})$ is a noisy parameter vector and $q(\hat{\boldsymbol{\theta}}_i)$ the state dependent cost. The matrix \mathbf{R} is a positive semi-definite matrix which represents the control costs. Hence, let \mathbf{A} be a finite difference matrix then $\boldsymbol{\theta}^\top \mathbf{R} \boldsymbol{\theta}$ with $\mathbf{R} = \mathbf{A}^\top \mathbf{A}$ leads to the sum of squared accelerations

$$\ddot{\boldsymbol{\theta}}^\top \ddot{\boldsymbol{\theta}} = \boldsymbol{\theta}^\top (\mathbf{A}^\top \mathbf{A}) \boldsymbol{\theta}. \quad (2)$$

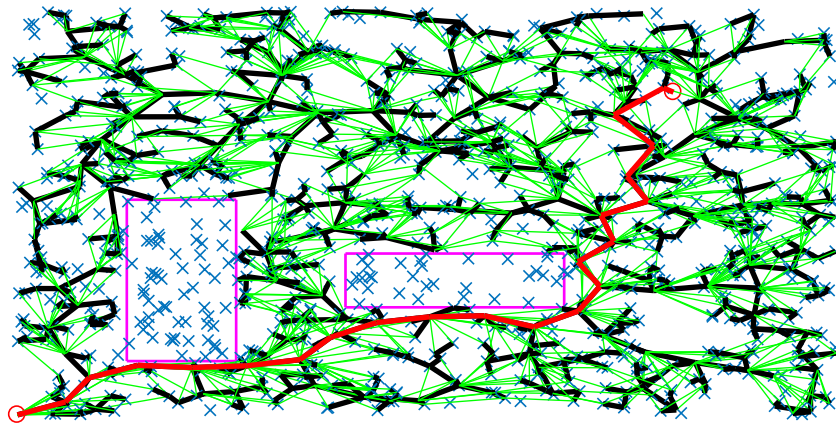


Figure 2 RRT* example. The rewiring is depicted using green lines.

Algorithm 2 RRT*, adjusted from [Karaman and Frazzoli, 2011]

```

1:  $\mathcal{V} \leftarrow \{\mathbf{x}_{\text{init}}\}; \mathcal{E} \leftarrow \emptyset$ 
2: for  $i = 1, \dots, n$  do
3:    $\mathbf{x}_{\text{rand}} \leftarrow \text{SampleRandom}$ 
4:    $\mathbf{x}_{\text{nearest}} \leftarrow \text{Nearest}(\mathcal{V}, \mathbf{x}_{\text{rand}})$ 
5:    $\mathbf{x}_{\text{new}} \leftarrow \text{Steer}(\mathbf{x}_{\text{nearest}}, \mathbf{x}_{\text{rand}})$ 
6:   if  $\text{ObstacleFree}(\mathbf{x}_{\text{nearest}}, \mathbf{x}_{\text{new}})$  then
7:      $\mathcal{X}_{\text{near}} \leftarrow \text{Near}(\mathcal{V}, \mathbf{x}_{\text{new}}, r(\text{card}(\mathcal{V})))$ 
8:      $\mathcal{V} \leftarrow \mathcal{V} \cup \{\mathbf{x}_{\text{new}}\}$ 
9:      $\mathbf{x}_{\text{min}} \leftarrow \mathbf{x}_{\text{nearest}}$ 
10:     $c_{\text{min}} \leftarrow \text{Cost}(\mathbf{x}_{\text{nearest}}) + c(\text{Line}(\mathbf{x}_{\text{nearest}}, \mathbf{x}_{\text{new}}))$ 
11:    for all  $\mathbf{x}_{\text{near}} \in \mathcal{X}_{\text{near}}$  do
12:      if  $\left( \begin{array}{l} \text{ObstacleFree}(\mathbf{x}_{\text{near}}, \mathbf{x}_{\text{new}}) \wedge \\ (\text{Cost}(\mathbf{x}_{\text{near}}) + c(\text{Line}(\mathbf{x}_{\text{near}}, \mathbf{x}_{\text{new}}))) < c_{\text{min}} \end{array} \right)$  then
13:         $\mathbf{x}_{\text{min}} \leftarrow \mathbf{x}_{\text{near}}$ 
14:         $c_{\text{min}} \leftarrow \text{Cost}(\mathbf{x}_{\text{near}}) + c(\text{Line}(\mathbf{x}_{\text{near}}, \mathbf{x}_{\text{new}}))$ 
15:      end if
16:    end for
17:     $\mathcal{E} \leftarrow \mathcal{E} \cup \{(\mathbf{x}_{\text{min}}, \mathbf{x}_{\text{new}})\}$ 
18:    for all  $\mathbf{x}_{\text{near}} \in \mathcal{X}_{\text{near}}$  do
19:      if  $\left( \begin{array}{l} \text{ObstacleFree}(\mathbf{x}_{\text{new}}, \mathbf{x}_{\text{near}}) \wedge \\ (\text{Cost}(\mathbf{x}_{\text{new}}) + c(\text{Line}(\mathbf{x}_{\text{new}}, \mathbf{x}_{\text{near}}))) < \text{Cost}(\mathbf{x}_{\text{near}}) \end{array} \right)$  then
20:         $\mathbf{x}_{\text{parent}} \leftarrow \text{Parent}(\mathbf{x}_{\text{near}})$ 
21:      end if
22:       $\mathcal{E} \leftarrow (\mathcal{E} \setminus \{(\mathbf{x}_{\text{parent}}, \mathbf{x}_{\text{near}})\}) \cup \{(\mathbf{x}_{\text{new}}, \mathbf{x}_{\text{near}})\}$ 
23:    end for
24:  end if
25: end for
26: return  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ 

```

In order to find the optimum of the optimization problem, the derivative of (1) is set to zero:

$$\nabla_{\hat{\theta}} \left(\mathbb{E} \left[\sum_{i=1}^N q(\hat{\theta}_i) + \frac{1}{2} \hat{\theta}^T \mathbf{R} \hat{\theta} \right] \right) = 0. \quad (3)$$

This leads to the expectation of the waypoints

$$\mathbb{E}(\hat{\theta}) = -\mathbf{R}^{-1} \mathbb{E} \left(\nabla_{\hat{\theta}} \left[\sum_{i=1}^N q(\hat{\theta}_i) \right] \right) = -\mathbf{R}^{-1} \delta \hat{\theta}_{\text{G}}. \quad (4)$$

In order to solve (4), STOMP proposes an estimated gradient

$$\delta \hat{\theta}_{\text{G}} = \int \delta \theta \, d\mathbf{P}. \quad (5)$$

In comparison to other approaches [Ratliff et al., 2009], which use the analytical gradient of the cost function to calculate an iterative gradient update, the STOMP approach allows the utilization of non-differentiable or non-smooth cost functions. With (5) the expectation of $\delta \theta$, which is the noise in $\hat{\theta}$, under the probability metric \mathbf{P} , can be calculated. The probability metric is defined as $\mathbf{P} = \exp \left(-\frac{1}{\lambda} \mathbf{S}(\hat{\theta}) \right)$ with

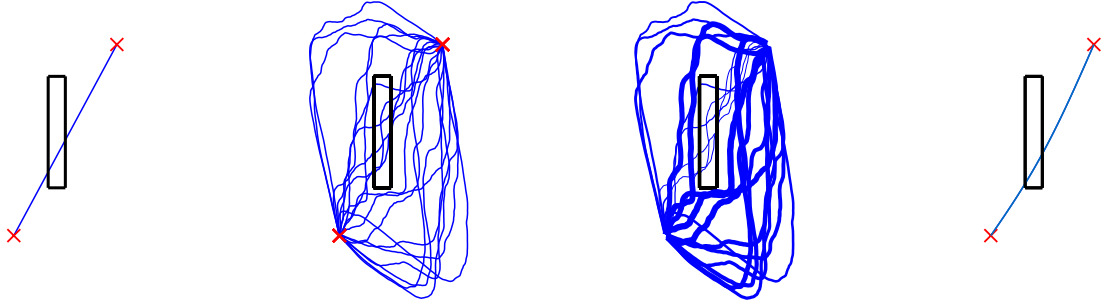


Figure 3 STOMP Procedure: Figure (a) shows the initial trajectory. In (b), the randomly sampled trajectories are depicted. Figure (c) shows the probabilities of the sampled trajectories. The higher the probability of a trajectory, the thicker the trajectory is drawn. In (d) the updated trajectory is presented.

$\mathbf{S}(\hat{\theta}) = \left[\sum_{i=1}^N q(\hat{\theta}_i) \right]$ being the state dependent cost defined on the trajectory. Considering the statements above, the stochastic gradient can be written as

$$\delta \hat{\theta}_G = \int \exp\left(-\frac{1}{\lambda} \mathbf{S}(\theta)\right) \delta \theta \, d\delta(\theta). \quad (6)$$

Using the given probability metric with $\mathbf{S}(\hat{\theta})$ being the cost of a path, paths with higher costs get a lower probability. Hence, they have a lower contribution to the update of the initial trajectory as paths with lower costs. In Figure 3 the STOMP procedure is exemplarily shown. The first figure shows the initial trajectory, a linear path from the start to the goal point, where the rectangular represents an obstacle, which should not be passed by the trajectory. The second picture shows the randomly sampled trajectories and the third depicts the different probabilities of the sampled trajectories. The higher the probability of a trajectory, the thicker the trajectory is drawn. The last picture shows the new generated trajectory. The algorithm starts again using this trajectory as new initial one. As exploration strategy for the algorithm, samples with variance $\Sigma = \mathbf{R}^{-1}$ are used. This ensures low control costs and smoothness. In Figure 4 thirty random samples drawn from a zero mean normal distribution $\epsilon = \mathcal{N}(\mathbf{0}, \mathbf{R}^{-1})$ are shown. In algorithm 3 pseudo-code for the reviewed STOMP method is presented.

Task 2: Reinforcement Learning

- The general Idea of Reinforcement Learning is to find a policy $\pi(\mathbf{u}|\mathbf{x}, \theta)$ which minimizes a cost function $\mathbf{R}(\tau)$ for a trajectory $\tau = (\mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \dots)$. Here \mathbf{u} are the control inputs into the system, \mathbf{x} are the states of the system and θ are the policy parameters. In Figure 5 a general overview over the subareas in Reinforcement Learning is shown.
- From exercise 0 we know that the differential equation for the pendulum can be written as

$$\ddot{\varphi} + \omega^2 \sin(\varphi) = u; \quad \omega = \sqrt{\frac{g}{L}} \quad (7)$$

for $m = 1$. With the given restriction of $u_{\max} = 5$ we can show that can not swing up the pendulum at once. Therefore we simulate the system with $u = u_{\max}, \forall t$ which leads to a result as depicted in Figure 6.

Algorithm 3 STOMP, based on [Kalakrishnan et al., 2011]

– **Require**

- * Start and goal positions x_0 and x_N
- * An initial 1-D discretized trajectory vector θ
- * A state-dependent cost function $q(\theta_i)$

– **Precompute**

- * \mathbf{A} : A central finite difference matrix
- * $\mathbf{R}^{-1} = (\mathbf{A}^\top \mathbf{A})^{-1}$
- * $\mathbf{M} = \mathbf{R}^{-1}$, each column scaled such that the maximum element is $1/N$

1: **procedure** STOMP(Repeat until convergence of trajectory cost $Q(\theta)$)

2: **for** $k = 1, \dots, K$ **do**

3: $\epsilon_k = \mathcal{N}(0, \mathbf{R}^{-1})$

4: $\hat{\theta}_k = \theta + \epsilon_k$

5: **for** $i = 2, \dots, N - 1$ **do**

6: $S(\hat{\theta}_{k,i}) = q(\hat{\theta}_{k,i})$

7: **end for**

8: **end for**

9: **for** $k = 1, \dots, K$ **do**

10: **for** $i = 2, \dots, N - 1$ **do**

11: $P(\hat{\theta}_{k,i}) = e^{-\frac{1}{\lambda} S(\hat{\theta}_{k,i})}$

12: **end for**

13: $P(\hat{\theta}_k) = \sum_{i=1}^{N-1} (P(\hat{\theta}_{k,i}))$

14: **end for**

15: $P(\hat{\theta}_k) \leftarrow \frac{P(\hat{\theta}_k)}{\sum_{k=1}^K P(\hat{\theta}_k)}$

16: $\delta \hat{\theta} = \sum_{k=1}^K P(\hat{\theta}_k) \epsilon_k$

17: $\delta \theta = \mathbf{M} \delta \hat{\theta}$

18: $\theta \leftarrow \theta + \delta \theta$

19: $Q(\theta) = \sum_{i=1}^N q(\theta_i) + \frac{1}{2} \theta^\top \mathbf{R} \theta$

20: **end procedure**

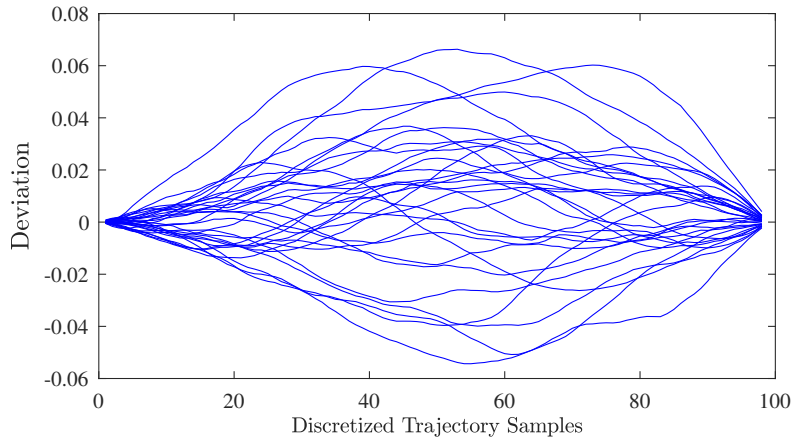


Figure 4 Thirty random samples of the noise drawn from a zero mean normal distribution $\epsilon = \mathcal{N}(0, \mathbf{R}^{-1})$.

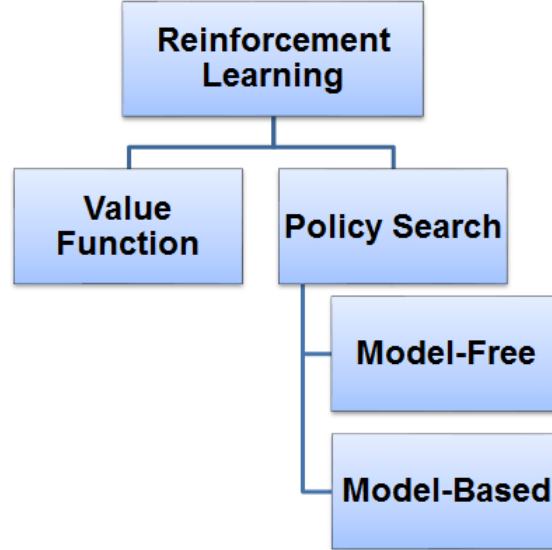


Figure 5 Reinforcement Learning

Thus we have to swing up the pendulum first to one and then to the other side in order to get enough momentum. Therefore we will use the following parametrized policy for the input signals \mathbf{u}

$$u(t) = \begin{cases} K_{p1}(g_1 - \varphi) - K_{d1}\dot{\varphi} & t \leq d_1 \\ K_{p2}(g_{\text{end}} - \varphi) - K_{d2}\dot{\varphi} & t > d_1 \end{cases} \quad (8)$$

Here $g_{\text{end}} = \pi$ is the end position at time T . The parameters which we have to determine are g_1 , the via-point, K_{p1} , K_{d1} , K_{p2} and K_{d2} are the gains for the PD-controller and d_1 is the time at when we change from steering to the via-point into steering to the end position. In the following we will write

$$\boldsymbol{\theta} = [K_{p1}, K_{d1}, K_{p2}, K_{d2}, d_1, g_1]^T. \quad (9)$$

In order to find convenient parameters we will define a cost function which has to be optimized. The intrinsic costs are defined as

$$L(\tau; \boldsymbol{\theta}) = \sum_{t=0}^T l(\varphi_t, u_t, t; \boldsymbol{\theta}) + c_p(\varphi_t, u_t) \quad (10)$$

with

$$l(\varphi_t, u_t, t; \boldsymbol{\theta}) = \frac{1}{N}(\varphi_t - g_i)^2. \quad (11)$$

Here N is the number of evaluations. We have decomposed our movement into two phases as described above such that g_1 is active as long $t \leq d_1$ and otherwise g_{end} is active. We consider the second term in 10 to be zero. Finally we require extrinsic costs which evaluate the overall task performance. We set them to

$$C = 100(x_T - g_{\text{end}}). \quad (12)$$

Thus, we have to reach the desired end state, otherwise our costs will explode. We can now start to find convenient parameters for our policy. Therefore we have to find the minimum for our costs. We will use Covariance Matrix Adaptation Evolution Strategy in order to search for a minimum. Matlab code is

available at https://www.lri.fr/~hansen/cmaes_inmatlab.html#matlab. As a optimization result we get

$$\theta = [K_{p1}, K_{d1}, K_{p2}, K_{d2}, d_1, g_1]^T = [6.13, -5.95, 33.83, 6.62, 1.74, -0.22]^T \quad (13)$$

as parameters. In Figure 7 the simulation result is shown.

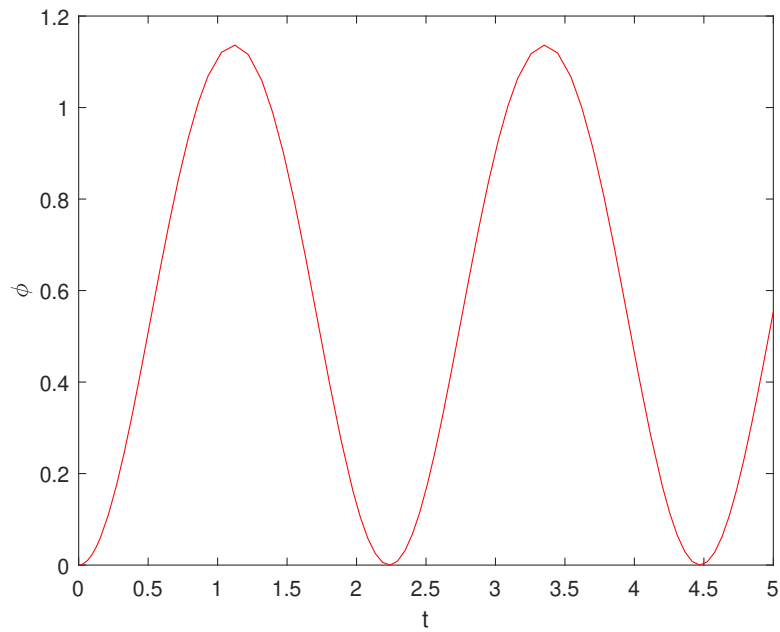


Figure 6 Pendulum with $u = u_{\max}$

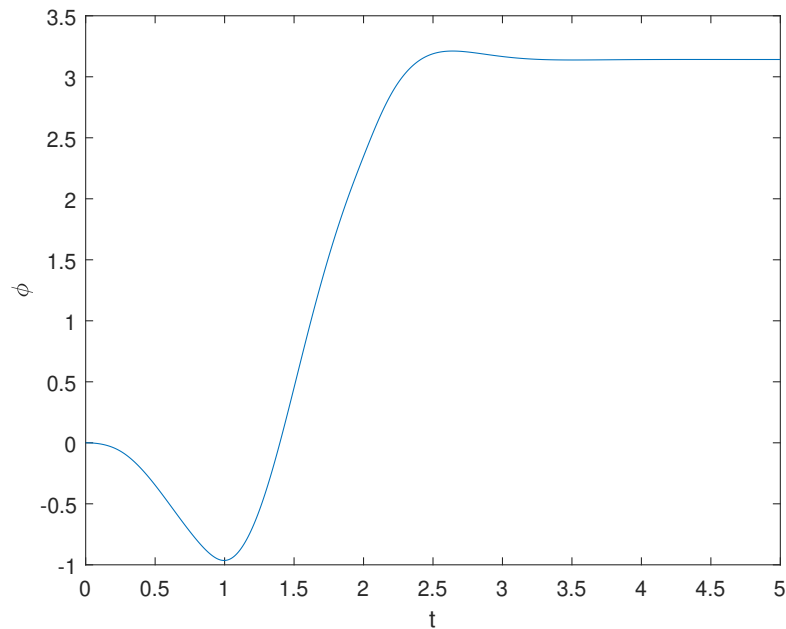


Figure 7 Up-Swing with via-point

Literatur

- [Kalakrishnan et al., 2011] Kalakrishnan, M., Chitta, S., Theodorou, E., Pastor, P., and Schaal, S. (2011). STOMP: Stochastic Trajectory Optimization for Motion Planning. In *International Conference on Robotics and Automation (ICRA), 2011*, pages 4569–4574. IEEE.
- [Karaman and Frazzoli, 2011] Karaman, S. and Frazzoli, E. (2011). Sampling-based Algorithms for Optimal Motion Planning. *The International Journal of Robotics Research*, 30(7):846–894.
- [LaValle, 1998] LaValle, S. M. (1998). Rapidly-exploring Random Trees: A new Tool for Path Planning.
- [Ratliff et al., 2009] Ratliff, N., Zucker, M., Bagnell, J. A., and Srinivasa, S. (2009). CHOMP: Gradient Optimization Techniques for Efficient Motion Planning. In *International Conference on Robotics and Automation (ICRA), 2009*, pages 489–494. IEEE.