

Robot Learning



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Winter Semester 2016, Homework 3 (54 points + 8 bonus)

Prof. Dr. J. Peters, F. Veiga, S. Parisi

Due Date: Tuesday, 31 January 2017 (3 pm)

Problem 3.1 Optimal Control [20 Points]

In this exercise, we consider a finite-horizon discrete time-varying Stochastic Linear Quadratic Regulator with Gaussian noise and time-varying quadratic reward function. Such system is defined as

$$\mathbf{s}_{t+1} = \mathbf{A}_t \mathbf{s}_t + \mathbf{B}_t \mathbf{a}_t + \mathbf{w}_t, \quad (1)$$

where \mathbf{s}_t is the state, \mathbf{a}_t is the control signal, $\mathbf{w}_t \sim \mathcal{N}(\mathbf{b}_t, \Sigma_t)$ is Gaussian additive noise with mean \mathbf{b}_t and covariance Σ_t and $t = 0, 1, \dots, T$ is the time horizon. The control signal \mathbf{a}_t is computed as

$$\mathbf{a}_t = -\mathbf{K}_t \mathbf{s}_t + \mathbf{k}_t \quad (2)$$

and the reward function r_t is

$$r_t = \begin{cases} -(\mathbf{s}_t - \mathbf{r}_t)^\top \mathbf{R}_t (\mathbf{s}_t - \mathbf{r}_t) - \mathbf{a}_t^\top \mathbf{H}_t \mathbf{a}_t & \text{when } t = 0, 1, \dots, T-1 \\ -(\mathbf{s}_t - \mathbf{r}_t)^\top \mathbf{R}_t (\mathbf{s}_t - \mathbf{r}_t) & \text{when } t = T \end{cases} \quad (3)$$

Note: r_t and \mathbf{r}_t are different!

Note 2: the notation used in Marc Toussaint's notes "*(Stochastic) Optimal Control*" is different from the one used in the lecture's slides.

a) Implementation [8 Points]

Implement the LQR with the following properties

$$\begin{aligned} \mathbf{s}_0 &\sim \mathcal{N}(0, 1) & T &= 50 \\ \mathbf{A}_t &= \begin{bmatrix} 1 & 0.1 \\ 0 & 1 \end{bmatrix} & \mathbf{B}_t &= \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} \\ \mathbf{b}_t &= \begin{bmatrix} 5 \\ 0 \end{bmatrix} & \Sigma_t &= 0.01 \\ \mathbf{K}_t &= \begin{bmatrix} 5 & 0.3 \end{bmatrix} & \mathbf{k}_t &= 0.3 \\ H_t &= 1 & \mathbf{R}_t &= \begin{cases} \begin{bmatrix} 100000 & 0 \\ 0 & 0.1 \end{bmatrix} & \text{if } t = 14, 40 \\ \begin{bmatrix} 0.01 & 0 \\ 0 & 0.1 \end{bmatrix} & \text{otherwise} \end{cases} & \mathbf{r}_t &= \begin{cases} \begin{bmatrix} 10 \\ 0 \end{bmatrix} & \text{if } t = 0, 1, \dots, 14 \\ \begin{bmatrix} 20 \\ 0 \end{bmatrix} & \text{if } t = 15, 16, \dots, T \end{cases} \end{aligned}$$

Execute the system 20 times. Plot the mean and 95% confidence over the different experiments of the state \mathbf{s}_t and of the control signal \mathbf{a}_t over time. How does the system behave? Compute and write down the mean and the standard deviation of the cumulative reward over the experiments. Attach a snippet of your code.

Robot Learning - Homework 3

Name, Vorname: _____ Matrikelnummer: □□□□□□□□

b) LQR as a P controller [4 Points]

The LQR can also be seen as a simple P controller of the form

$$\mathbf{a}_t = \mathbf{K}_t (\mathbf{s}_t^{\text{des}} - \mathbf{s}_t) + \mathbf{k}_t, \quad (4)$$

which corresponds to the controller used in the canonical LQR system with the introduction of the target $\mathbf{s}_t^{\text{des}}$.

Assume as target

$$\mathbf{s}_t^{\text{des}} = \mathbf{r}_t = \begin{cases} \begin{bmatrix} 10 \\ 0 \end{bmatrix} & \text{if } t = 0, 1, \dots, 14 \\ \begin{bmatrix} 20 \\ 0 \end{bmatrix} & \text{if } t = 15, 16, \dots, T \end{cases} \quad (5)$$

Use the same LQR system as in the previous exercise and run 20 experiments. Plot in one figure the mean and 95% confidence of the first state, for both $\mathbf{s}_t^{\text{des}} = \mathbf{r}_t$ and $\mathbf{s}_t^{\text{des}} = \mathbf{0}$.

c) Optimal LQR [8 Points]

To compute the optimal gains \mathbf{K}_t and \mathbf{k}_t , which maximize the cumulative reward, we can use an analytic optimal solution. This controller recursively computes the optimal action by

$$\mathbf{a}^* = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T (\mathbf{V}_{t+1} (\mathbf{A}_t \mathbf{s}_t + \mathbf{b}_t) - v_{t+1}), \quad (6)$$

which can be decomposed into

$$\mathbf{K}_t = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t, \quad (7)$$

$$\mathbf{k}_t = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T (\mathbf{V}_{t+1} \mathbf{b}_t - v_{t+1}). \quad (8)$$

where

$$\mathbf{M}_t = \mathbf{B}_t (\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t \quad (9)$$

$$\mathbf{V}_t = \begin{cases} \mathbf{R}_t + (\mathbf{A}_t - \mathbf{M}_t)^T \mathbf{V}_{t+1} \mathbf{A}_t & \text{when } t = 1 \dots T-1 \\ \mathbf{R}_t & \text{when } t = T \end{cases} \quad (10)$$

$$v_t = \begin{cases} \mathbf{R}_t \mathbf{r}_t + (\mathbf{A}_t - \mathbf{M}_t)^T (\mathbf{V}_{t+1} \mathbf{b}_t - v_{t+1}) & \text{when } t = 1 \dots T-1 \\ \mathbf{R}_t \mathbf{r}_t & \text{when } t = T \end{cases} \quad (11)$$

Run 20 experiments with $\mathbf{s}_t^{\text{des}} = \mathbf{0}$ computing the optimal gains \mathbf{K}_t and \mathbf{k}_t . Plot the mean and 95% confidence of both states for all three different controllers used so far. Use one figure per state. Report the mean and std of the cumulative reward for each controller and comment the results. Attach a snippet of your code.

Robot Learning - Homework 3

Name, Vorname: _____ Matrikelnummer:

Problem 3.2 Reinforcement Learning [34 Points + 8 Bonus]

You recently acquired a robot for cleaning you apartment but you are not happy with its performance and you decide to reprogram it using the latest AI algorithms. As a consequence the robot became self-aware and, whenever you are away, it prefers to play with toys rather than cleaning the apartment. Only the cat has noticed the strange behavior and attacks the robot. The robot is about to start its day and its current perception of the environment is as following

Your graphic could be here. It just wasn't included.

The black squares denote extremely dangerous states that the robot must avoid to protect its valuable sensors. The reward of such states is set to $r_{\text{danger}} = -10^5$ (NB: the robot can still go through these states!). Moreover, despite being waterproof, the robot developed a phobia of water (W), imitating the cat. The reward of states with water is $r_{\text{water}} = -100$. The robot is also afraid of the cat (C) and tries to avoid it at any cost. The reward when encountering the cat is $r_{\text{cat}} = -3000$. The state containing the toy (T) has a reward of $r_{\text{toy}} = 1000$, as the robot enjoys playing with them. Some of the initial specification still remain, therefore the robot receives $r_{\text{dirt}} = 35$ in states with dirt (D).

State rewards can be collected at every time the robot is at that state. The robot can perform the following actions: *down*, *right*, *up*, *left* and *stay*.

In our system we represent the actions with the an ID (0, 1, 2, 3, 4), while the grid is indexed as {row, column}. The robot can't leave the grid as it is surrounded with walls. A skeleton of the gridworld code and some plotting functions are available at the webpage. For all the following questions, always attach a snippet of your code.

a) Finite Horizon Problem [14 Points]

In the first exercise we consider the finite horizon problem, with horizon $T = 15$ steps. The goal of the robot is to maximize the expected return

$$J_{\pi} = \mathbb{E}_{\pi} \left[\sum_{t=1}^{T-1} r_t(s_t, a_t) + r_T(s_T) \right], \quad (12)$$

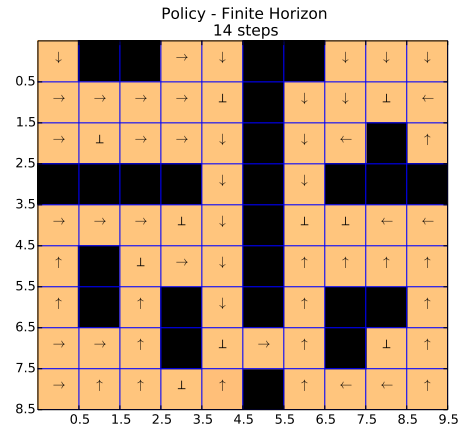
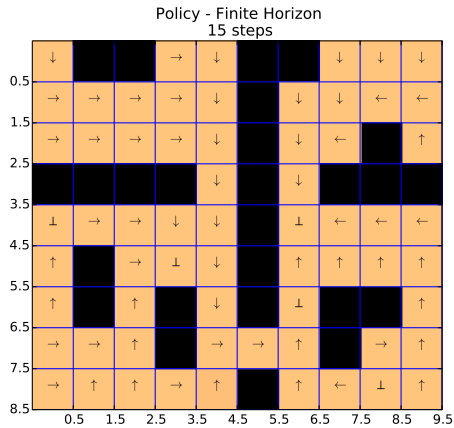
according to policy π , state s , action a , reward r , and horizon T . Since rewards in our case are independent of the action and the actions are deterministic, Equation (12) becomes

$$J_{\pi} = \sum_{t=1}^T r_t(s_t). \quad (13)$$

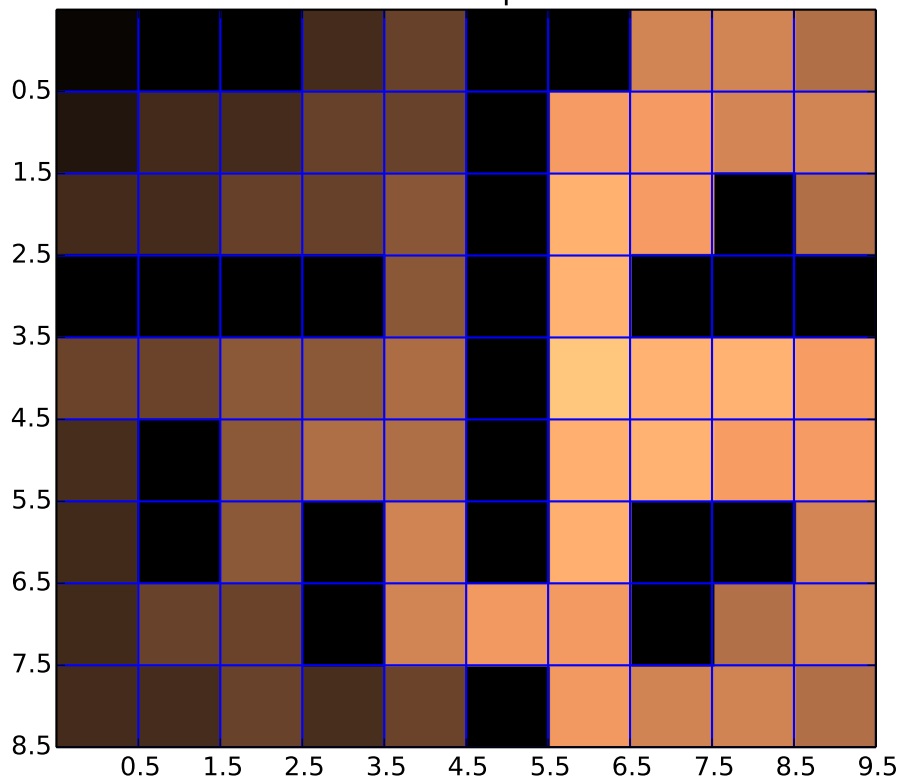
Using the Value Iteration algorithm, determine the optimal action for each state when the robot has 15 steps left. Attach the plot of the policy to your answer and a mesh plot for the value function. Describe and comment the policy: is the robot avoiding the cat and the water? Is it collecting dirt and playing with the toy? Which would be the time horizon that makes the robot acts differently in state (9, 4)?

Robot Learning - Homework 3

Name, Vorname: _____ Matrikelnummer:



Value Function - Finite Horizon
15 steps



Our policy seems appropriate given our explanation. The robot does collect dirt if it is far enough away from the toy (5,1) but otherwise will prioritize getting towards the toy. It seems so reckless that it even walks over the cat space to get to it when starting in (9,4). This however changes once the robot is limited to 14 steps forcing our robot to simply collect dirt beside the cat in (9,4).

Robot Learning - Homework 3

Name, Vorname: _____ Matrikelnummer:

b) Infinite Horizon Problem - Part 1 [4 Points]

We now consider the infinite horizon problem, where $T = \infty$. Rewrite Equation (12) for the infinite horizon case adding a discount factor γ . Explain briefly why the discount factor is needed.

c) Infinite Horizon Problem - Part 2 [6 Points]

Calculate the optimal actions with the infinite horizon formulation. Use a discount factor of $\gamma = 0.8$ and attach the new policy and value function plots. What can we say about the new policy? Is it different from the finite horizon scenario? Why?

d) Finite Horizon Problem with Probabilistic Transition Function [10 Points]

After a fight with the cat, the robot experiences control problems. For each of the actions *up*, *left*, *down*, *right*, the robot has now a probability 0.7 of correctly performing it and a probability of 0.1 of performing another action according to the following rule: if the action is *left* or *right*, the robot could perform *up* or *down*. If the action is *up* or *down*, the robot could perform *left* or *right*. Additionally, the action can fail causing the robot to remain on the same state with probability 0.1. Using the finite horizon formulation, calculate the optimal policy and the value function. Use a time horizon of $T = 15$ steps as before. Attach your plots and comment them: what is the most common action and why does the learned policy select it?

e) Reinforcement Learning - Other Approaches [8 Bonus Points]

What are the two assumptions that let us use the Value Iteration algorithm? What if they would have been not satisfied? Which other algorithm would you have used? Explain it with your own words and write down its fundamental equation.