

**Deep Learning and Artificial Intelligence**  
 WS 2024/25

**Exercise 11: Modelfree Reinforcement Learning**

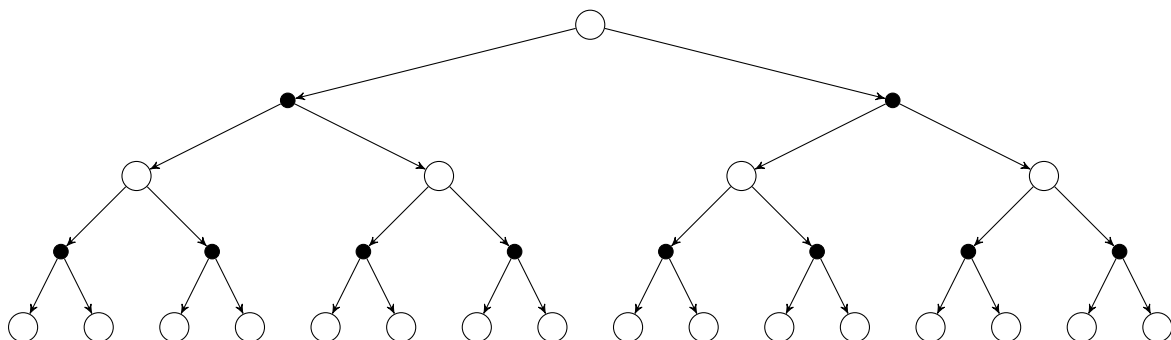
**Exercise 11-1 RL Theory Questions**

- In Reinforcement Learning, what do the terms ‘prediction’ and ‘control’ refer to?
- What are the differences between model-free and model-based methods? What do they have in common? Which category do Dynamic Programming (DP), Temporal-Difference-Learning (TD) and Monte Carlo (MC) -methods belong to?
- Write down the Bellman equations for optimal state values (utilities), denoted  $u_*(s)$ , and optimal state-action values, denoted  $q_*(s, a)$ . How are  $u_*(s)$  and  $q_*(s, a)$  related to each other?
- What does on- and off-policy learning mean? Name an example of each.
- Why is Q-learning an off-policy method? What is the difference to Sarsa?
- When do we need exploration and why? Why is there a trade-off between exploitation and exploration?

**Exercise 11-2 Backup Strategies**

In the lecture you learned three different backup strategies for the Bellman equation (Dynamic Programming (DP), Monte Carlo (MC) and Temporal Difference Learning (TD)).

- Assume we have an MDP with exactly two steps (two times an action is performed). The figure below visualizes this MDP. An unfilled circle represents a state and a filled circle represents an action respectively.  
 For each of the different backup strategies, mark the paths that are used. For DP, assume a fixed policy  $\pi(S_t)$  in the first step (instead of max). For MC and TD, just choose arbitrary actions. Also write down the formulas for these updates.



- (b) How do the three backup strategies compare regarding variance, efficiency, necessity of a model and bias?

### **Exercise 11-3      Model-free RL in Python**

On moodle course page you can find a Jupyter notebook file with a programming exercise for model-free reinforcement learning. Please follow the instructions in the notebook.