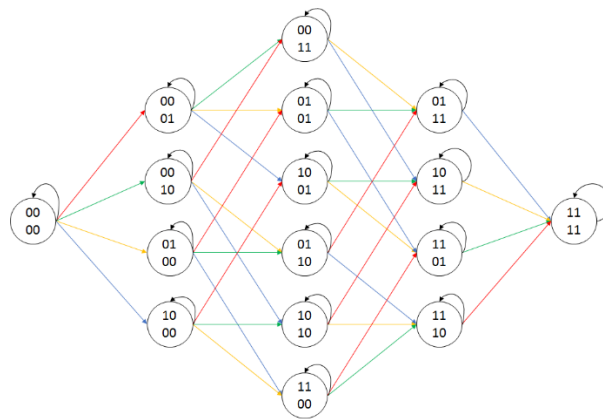# 1.  Exercise 1

## 1.1    Markov Decision Processes

a) Formalize the above described problem as a Markov Decision Process. You do not have to write out the individual elements completely (i.e. name all states explicitly or provide a full transition matrix), but define the contents in unambiguous expressions.



**State set:**

Each state is represented as an encoding of the solved tasks. E.g.:

01 10 means, that the tasks 2 and 3 have been solved, but the tasks 1 and 4 have not been solved yet (0110 = 6)

00 11 means, that the tasks 1 and 2 have been solved, but the tasks 3 and 4 have not been solved yet (0011 = 3)

**Action set:**

The actions are referenced by colors: $\mathcal{A} = \{a_{Task1}, a_{Task2}, a_{Task3}, a_{Task4}\}$

**Transition function:**

The state transition is defined by an arrow, colored accordingly to the action, which must be taken to perform the state transition. The color encoding contains the information of the probability to propagate from one state to another: 0.1 for action 1, 0.8 for action 2, 0.3 for action 3 and 0.5 for action 4. For better readability, all self-loops have been colored black which means, that all the outgoing edges to other states can also keep the agent in its current state.

**Reward function:**

Like the actions the rewards are encoded on the colored edges. The agent is rewarded 4 points for action 1, 1 point for action 2, 3 points for action 3 and 2 points for action 4. Black edges have a reward of 0 points.

b) How would you model the risk of failing the exam in this scenario?

*Hint: You can introduce terminal rewards.*

A terminal reward of -4 is introduced for all states were the sum of the points is lower than 4.

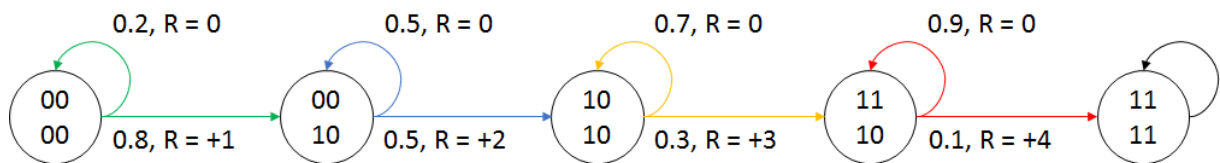If the final reward is 0 or higher the student passed the exam. If it is lower, the student failed.

All probability of taking a branch where the student ends up in a failed task are then summed up.

c) Student $S$ considers two possible policies, $\pi_A^S$ and $\pi_B^S$, that determine in which order the tasks will be solved. The first policy follows the tasks in increasing difficulty, i.e. with decreasing possibility of success $p_i^S$ . Using policy $\pi_B^S$, the tas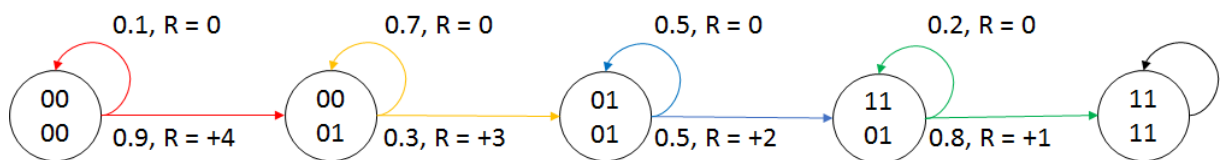ks are tackled in reversed order (reversed with respect to $\pi_A^S$). Compare both policies by determining the path costs for both policies $\pi_A^S$ and $\pi_B^S$.

*Remember: The student knows about the success of solving a task after dealing with it and may choose to solve a task again in the case of failing.*

**Markov Chain for $\pi_A^S$:**



**Markov Chain for $\pi_B^S$:**



$$v_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1})|S_t = s]$$

*Remark: We put a lot of effort into writing a general solving algorithm for this type of MDP problems (see code). However, after the meeting on Friday we knew, we made things too complicated and we had understanding issues what really is expected from us to program for the finite horizon problem.*

d) Derive a stationary policy $\pi_C^S$ that is more promising than $\pi_A^S$ and $\pi_B^S$. In what order does this strategy choose the tasks?

The stationary policy $\pi_C^S$ chooses the action with the highest expected direct reward:

$$\pi_C^S(a|s) = \begin{cases} 1 & if \ a = arg\max_{a\in\mathcal{A}} \mathbb{E}_\pi[R_t|S_t = s, A_t = a] \\ 0 & otherwise \end{cases}$$

The order is 4-3-2-1.

e) Suggest a method that reduces the probability of failing for each of the above policies. You can describe in non-formal sentences. *Hint: Use non-stationary policies.*

For policy $\pi_A^S$: When the agent solves task 2 on the third try, it will try to solve task 3 directly because solving task 4 will not help the agent preventing to fail.

For policy $\pi_B^S$: When the agent did not solve task 1 at least on the third try, it will try to solve task 2.

-   In case of success it will try to solve task 3 to get the needed 4 total points.
-   In case of failure it will try to solve task 1 again.

For policy $\pi_C^S$: When the agent did not solve task 4 at least on the third try, it will try to solve task 2.

-   In case of success it will try to solve task 3 to get the needed 4 total points.
-   In case of failure it will try to solve task 1 to get the for points at once.

f) Student $\mathcal{T}$ has learned selectively and did not prepare the topics of tasks 2-4 ($p_i^{\mathcal{T}} = 0.0$ for $i \in \{2,3,4\}$). How good must student $\mathcal{T}$ be prepared for the topic of task 1 (i.e. how high has $p_1^{\mathcal{T}}$ to be at least) in order to be at least as successful as student $S$, if student $S$ follows policy $\pi_A^S$?

Student $\mathcal{T}$ must prepare for task 1 so that the $p_1^{\mathcal{T}} \approx 22.97\%$ to perform as good as student $S$ with $\pi_A^S$. This result has been calculated with Excel.

## 1.2    Bellman Equation

a) Show exemplary for state $s_{3,3}$ in the middle of the grid with $v(s_{3,3} = 0.7)$ that the Bellman equation is satisfied for all neighboring states.

$$v(s_{3,3}) = \mathcal{R}_{s_{3,3}} + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{s_{3,3}s'} v(s')$$

$$\mathcal{R}_{s_{3,3}} = 0, \qquad \gamma = 0.9$$

$$\mathcal{P}_{s_{3,3}s_{2,3}} = \mathcal{P}_{s_{3,3}s_{4,3}} = \mathcal{P}_{s_{3,3}s_{3,2}} = \mathcal{P}_{s_{3,3}s_{3,4}} = 0.25$$

$$v(s_{2,3}) = 2.3, \qquad v(s_{4,3}) = -0.4, \qquad v(s_{3,2}) = 0.7, \qquad v(s_{3,4}) = 0.4,$$

$$v(s) = 0 + 0.9 * [0.25 * (0.4 + 0.7 + 2.3 - 0.4)] = 0.675 \approx 0.7$$

b) Explain why the value of state $B$ is higher than the direct reward. Why does this not hold for state $A$

A state valuation function tries to judge the goodness of the actual state. For that it takes the expectation of the direct reward **and** the discounted expected state valuation of the next possible states into account. This can directly be seen by examining the Bellman equation:

$$v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s')$$

**Expected direct reward (first part of the Bellman equation $\mathcal{R}_s$):**

For both states A and B the actual reward is determined because every action will give the reward for jumping into A' and respectively B'.

$$\mathcal{R}_A = 10, \qquad \mathcal{R}_B = 5$$

**Expected state valuation (second part of the Bellman equation $\gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s')$):**

When the agent will be in state A, every action will take it to state A', respectively from B to B'.

$$\mathcal{P}_{AA'} = \mathcal{P}_{BB'} = 1$$

$$v(A) = \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{As'} v(s') = 0.9 * [1 * (-1.3)] = -1.17 \approx -1.2$$

$$v(B) = \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{Bs'} v(s') = 5 + 0.9 * [1 * (0.4)] = 0.36 \approx 0.4$$

The expectation of the next states valuation is positive for B and negative for A.

**Both parts combined:**

$$v(A) = \mathcal{R}_A + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{As'} v(s') = 10 - 1.2 \approx 8.8$$

$$v(B) = \mathcal{R}_B + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{Bs'} v(s') = 5 + 0.4 = 5.4$$

The second part of the equation makes the difference between state valuation and direct reward.
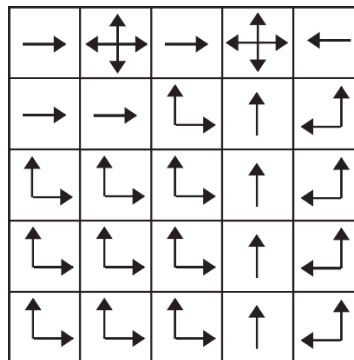
The optimal policy $\pi_*$ is shown in Figure 2. Now assume that we reduce the reward for jumping from A to $A'$ to 4.

(c) How does the optimal policy change? What are the new values of states $A$ and $B$?

Before the reward change the optimal policy for the agent to do was to loop the behavior of striving towards A and being teleported to A'. (Only exception is B, because it doesn't matter what the agent is doing, it will be teleported to B').

After the reward change it becomes more rewarding to loop between B and B' because the distance between B and B' is shorter than the distance between A and A'. Plus, the reward is higher for looping between B and B' than looping between A and A'.

In a graphical way, the optimal policy looks now like this:



The policy defines the action so that the agent should move towards B with every step. An interesting behavior to point out is that the agent chooses to move from state $s_{1,1}$ to A. The agent might lose some reward for getting teleported to A' and being therefore teleported away from B but the direct reward outweighs this 'distance loss'.

The following state valuations of A and B have been calculated with the program 'GridWorld.py'

$$v_{\pi_*}(A) \approx 12.8$$

$$v_{\pi_*}(B) \approx 18.45$$

## 1.3    Bonus: Experiences

Submit an *experiences.txt*, where you provide a brief summary of your experience with this exercise, the corresponding lecture and the last meeting. As a minimum, say how much time you invested and if you had major problems – and if yes, where.

Nico

|  | Hours | Experience | Problems |
|---|---|---|---|
| **Exercise 1** | 13+3 | Bellman part quite easy. Just use the formula. MDP/MRP task was good for understanding. In general, the exercise sheet conception is excellent. | Problem formulation for the MDP/MRP part was a confusing. Especially Task 1c). |
| **Lecture 2** | 4 | Really had to look up the mathematical definitions sometimes to be absolute clear what David meant. | Posted in the Google Doc |
| **Meeting** | 1.5 | More time should be used on the exercises. | |
| **General** | | Flipped classroom concept works out nicely. Splitting tasks between teachers benefits a lot. Kahoot is cool, but takes a lot of time. Somehow you are forced to answer with time pressure which is a good exercise, however, | |

Hendrik

|  | Hours | Experience | Problems |
|---|---|---|---|
| **Exercise 1** | 12+8 | **Task1:** subtasks a) and b) were good to understand and pretty much straightforward but I did not manage to solve c) and the following subtasks on my own. **Task 2:** good and well-structured task | 1c) was very confusing for me with the path cost expression and I did not know, what value is wanted here |
| **Lecture 2** | 3 | | |
| **Meeting** | 1.5 | | |
| **General** | | | |