

Modul - Unsupervised and Reinforcement Learning (URL)

Bachelor Programme AAI

09 - Markov Decision Processes (MDP)

Prof. Dr. Marcel Tilly

Faculty of Computer Science, Cloud Computing

Agenda



Code can be found in

- [RL_Policy_Optimization.ipynb](#)
- or on GitHub or [hosted on myBinder](#)

On the menu for today:

- Reinforcement Learning
 - Agent, Rewards, States and Actions
 - Markov Decision Processes
 - Bellmann equation/function
 - Policy Evaluation, Improvement and Iteration



Agent and Environment

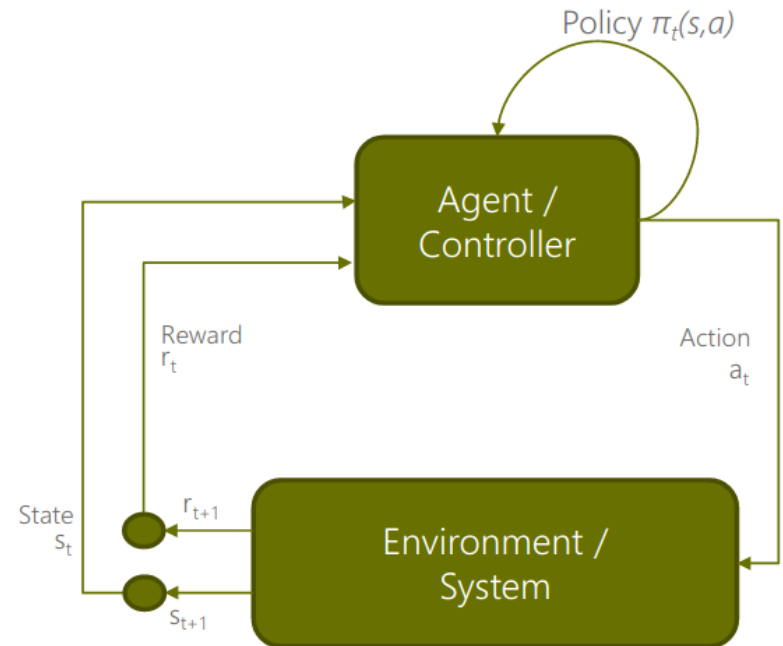


At each step the agent

- receives state observations S_t
- receives scalar reward R_t
- executes action A_t

The environment

- receives action A_t
- emits state observation S_{t+1}
- emits scalar reward R_{t+1}



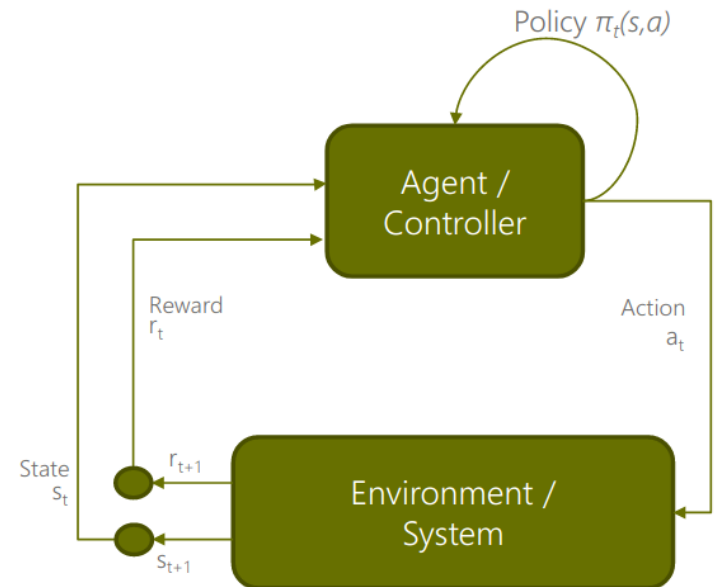
RL Agent: Policy



Telling the agent what to do:

- Deterministic $a = \pi_t(s)$
- Stochastic
 $\pi_t(s, a) = Pr(a_t = a | s_t = s)$

Given the situation at time t is state s , the policy gives the probability the agent's action will be a .



Goal is to find the policy(π) that maximizes rewards!

Deterministic Policy

Example: Frozenlake

```
# Make the environment based on deterministic policy (is_slippery=False)
env = gym.make('FrozenLake-v1', is_slippery=False, render_mode="ansi")
# Go right once (action = 2)
env.reset()
action = 2
(observation, reward, done, prob, a) = env.step(action)
print(env.render())
# Observation = 0: move to the right once from grid 0 to grid 1
# Prob = 1: deterministic policy, if we choose to go right, we will go right
print(observation, reward, done, prob, a)
```

		(Right)
[S]FFF		S[F]FF
FHFH		FHFH
FFFH	==>	FFFH
HFFG		HFFG

Example: Frozenlake

```
# Make the environment based on non-deterministic policy (is_slippery=True)
env = gym.make('FrozenLake-v1', is_slippery=True, render_mode="ansi")
# Go right once (action = 2)
env.reset()
action = 2
(observation, reward, done, prob, a) = env.step(action)
print(env.render())
# Observation = 0: move to the right once from grid 0 to grid 1
# Prob = 1/3: deterministic policy, if we choose to go right, we might not go right
print(observation, reward, done, prob, a)
```

[S]FFF		(Right)
FHFH		SFFF
FFFH	==>	[F]HFH
HFFG		FFFH
		HFFG

- Typically we can frame all RL tasks as MDPs
 - The key in MDPs is the **Markov Property**
 - Essentially the future depends on the present and not the past
 - More specifically, the future is independent of the past given the present
 - There's an assumption the present state encapsulates past information.
- Putting into the context of what we have covered so far: our agent can (1) **control its action based** on its current (2) **completely known state**

Two main characteristics for MDPs

- Control over state transitions
- States completely observable

Types of Markov Models

Permutations of whether there is presence of the two main characteristics would lead to different Markov models:

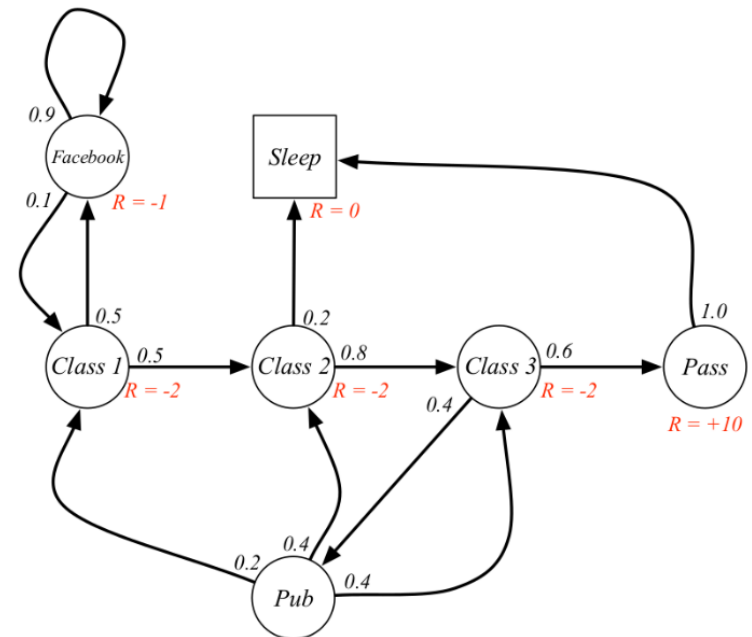
- Control over state transitions and completely observable states: MDPs
- Control over state transitions and partially observable states: Partially Observable MDPs (POMDPs)
- No control over state transitions and completely observable states: Markov Chain
- No control over state transitions and partially observable states: Hidden Markov Model

Markov Decision Process (MDP)



MDP is represented by five important elements:

- A set of states S the agent can actually be in
- A set of actions A that can be performed by an agent, for moving from one state to another
- A transition probability $P_{ss'}^a$, which is the probability of moving from one state S to another S' by performing an action a
- A reward probability $R_{ss'}^a$, which is the probability of a reward acquired by the agent for moving from one state S to another S' by performing some action a
- A discount factor γ , which controls the importance of immediate and future rewards



We may or may not know our model

- **Model-based RL:** this is where we can **clearly define** our (1) transition probabilities and/or (2) reward function
 - A global minima can be attained via Dynamic Programming (DP)
- **Model-free RL:** this is where we cannot clearly define our (1) transition probabilities and/or (2) reward function
 - Most real-world problems are under this category so we will mostly place our attention on this category

Dynamic Programming

Dynamic programming is breaking down a problem into smaller sub-problems, solving each sub-problem and storing the solutions to each of these sub-problems in an array (or similar data structure) so each sub-problem is only calculated once.

- It is both a mathematical optimization method and a computer programming method.
- Richard Bellman invented DP in the 1950s.

Example:

```
def fibonacciVal(n):  
    memo[0], memo[1] = 0, 1  
    for i in range(2, n+1):  
        memo[i] = memo[i-1] + memo[i-2]  
    return memo[n]
```

- When the agent acts given its state under the policy ($\pi(a|s)$), the transition probability function $P_{ss'}^a$, determines the subsequent state (s')

$$P = P(s' | s, a) = P[S_{t+1} = s' | S_t = s, A_t = a]$$

- When the agent act based on its policy ($\pi(a|s)$) and transited to the new state determined by the transition probability function $P_{ss'}^a$, it gets a reward based on the reward function R_s^a as a feedback

$$R = E[R_{t+1} | S_t = s, A_t = a]$$

- Rewards are short-term, given as feedback after the agent takes an action and transits to a new state. Summing all future rewards and discounting them would lead to our return G_t

$$G = \sum \gamma^i R_{t+1+i}$$

Goal: Maximize total rewards!

- Immediate reward r at any time t_x (each step).
- If a reward r at any time t_x is defined as r_x then the cumulative future reward R_t at time t is defined as where T is the terminal time:

$$R_t = r_{t+1} + r_{t+2} + \dots + r_T$$

- with $T = \infty$, we use discounted future rewards :

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

where γ is the discount factor given as $0 \leq \gamma \leq 1$

- If $\gamma > 0$ then we say the algorithm is myopic
- If $\gamma > 1$ then the algorithm is far-sighted (it maximizes future rewards!)

State-Value Function

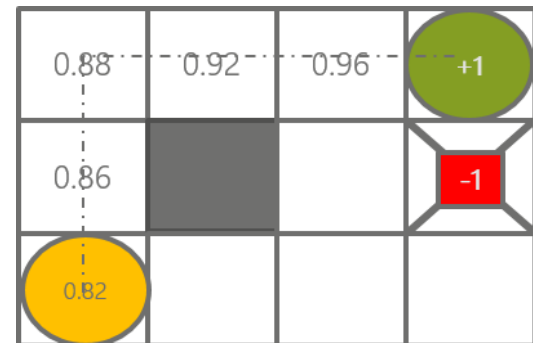


- We need to calculate the accumulative reward from being in that state and follow the policy until the goal state, where the episode ends.
- We could calculate the value of 0.82 for the first given state.
- Also, we could calculate the same way we did, for each of the states that make up the path to the goal state.
- We are looking for the utility of a state under a given policy

The utility value is the summary of the reward of itself and all the future states following that policy.

$$V^{\pi}(s) = \mathbb{E}[\sum_{t \geq 0} \gamma^t r_t | s_t = s, \pi]$$

For each **state s**
It yields the **expected return**
If the agent **starts in state s**
And then uses **the policy**
To choose its action



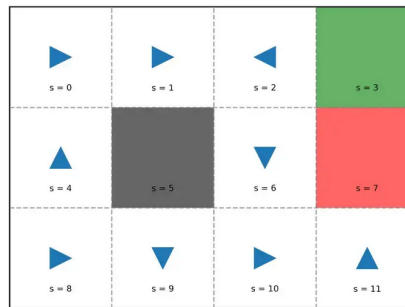
Utility or Value of States

$s = 0$ $r = -0.04$	$s = 1$ $r = -0.04$	$s = 2$ $r = -0.04$	$s = 3$ $r = 1.0$
$s = 4$ $r = -0.04$	$s = 5$ $r = \text{nan}$	$s = 6$ $r = -0.04$	$s = 7$ $r = -1.0$
$s = 8$ $r = -0.04$	$s = 9$ $r = -0.04$	$s = 10$ $r = -0.04$	$s = 11$ $r = -0.04$

Let's consider a *non-deterministic* MDP:

- We are standing in $s = 8$. If we go to $s = 9$, we receive a reward of -0.04 ; if we go to $s = 4$, we also receive a reward of -0.04 .
 - If we choose $s = 9$, it may lead us to the path of $[9, 10, 11, 7]$ corresponding to rewards of $[-0.04, -0.04, -0.04, -1]$. Since MDP has the property of additive rewards, the total reward along this path is thus -1.12 .
 - If we choose $s = 4$, it may lead us to the path of $[4, 0, 1, 2, 3]$ corresponding to rewards of $[-0.04, -0.04, -0.04, -0.04, +1]$ with a total reward of 0.84 .
- Here -1.12 is the utility of state $s = 9$ while 0.84 is the utility of state $s = 4$. We usually use V to denote the *utility of states* or the *state values*.

Deterministic case



For instance, if we look at $s = 6$ and we follow the policy π , we have the path $[6, 10, 11, 7]$ in front of us. Therefore, the utility of $s = 6$ is the sum of $r(6)$, $r(10)$, $r(11)$, and $r(7)$.

$$V(6)_{\pi=[6,10,11,7]} = r(6) + r(10) + r(11) + r(7)$$

or

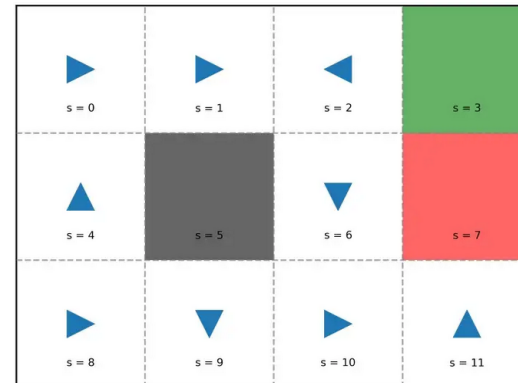
$$V(6)_{\pi=[6,10,11,7]} = r(6) + V(10)$$



Stochastic case

In a stochastic environment, even if we follow the policy, we are not guaranteed to move along the path [6, 10, 11, 7]

$$\begin{aligned} v(6) &= 0.8 \times \text{discounted rewards following policy}[6, 10\dots] \\ &\quad + 0.1 \times \text{discounted rewards following policy}[6, 7] \\ &\quad + 0.1 \times \text{discounted rewards following policy}[6, 6\dots] \\ &= 0.8 \times 0.8 \times \text{discounted rewards following policy}[6, 10, 11\dots] \\ &\quad + 0.8 \times 0.1 \times \text{discounted rewards following policy}[6, 10, 10\dots] \\ &\quad + 0.1 \times \text{discounted rewards following policy}[6, 7] \\ &\quad + 0.1 \times 0.8 \times \text{discounted rewards following policy}[6, 6, 10\dots] \\ &\quad + 0.1 \times 0.1 \times \text{discounted rewards following policy}[6, 6, 7] \\ &\quad + 0.1 \times 0.1 \times \text{discounted rewards following policy}[6, 6, 6\dots] \end{aligned}$$



We can generalize to:

$$V^{\pi}(s) = E\left[\sum_{t \geq 0} \gamma^t r_t \mid s_t = s, \pi\right]$$

Stochastic case

With stochastic, the utility of a state under a certain policy can also be represented as the sum of its immediate reward and the utility of its successor state following a probability distribution:

$$V(s) = R_s + \gamma \sum_{s'} P_{ss'}^a V(s')$$

Example:

$$V(6) = r(6) + \gamma[0.8V(10) + 0.1V(6) + 0.1V(7)]$$

Finally, we will have 11 equations with 11 unknowns.

$$\begin{aligned} v(0) &= r(0) + \gamma [0.8v(1) + 0.1v(0) + 0.1v(4)] \\ v(1) &= r(1) + \gamma [0.8v(2) + 0.1v(1) + 0.1v(1)] \\ v(2) &= r(2) + \gamma [0.8v(1) + 0.1v(6) + 0.1v(2)] \\ v(3) &= r(3) + \gamma [v(3)] \\ v(4) &= r(4) + \gamma [0.8v(0) + 0.1v(4) + 0.1v(4)] \\ v(6) &= r(6) + \gamma [0.8v(10) + 0.1v(6) + 0.1v(7)] \\ v(7) &= r(7) + \gamma [v(7)] \\ v(8) &= r(8) + \gamma [0.8v(9) + 0.1v(4) + 0.1v(8)] \\ v(9) &= r(9) + \gamma [0.8v(9) + 0.1v(8) + 0.1v(10)] \\ v(10) &= r(10) + \gamma [0.8v(11) + 0.1v(6) + 0.1v(10)] \\ v(11) &= r(11) + \gamma [0.8v(7) + 0.1v(10) + 0.1v(11)] \end{aligned}$$

Policy evaluation: for a given policy, we evaluate that policy by determining the utility of each state.

Instead of directly solving the linear equations, we can use a *dynamic programming approach* called *iterative policy evaluation*.

1. We first initialize the utility of each state as zero, then we loop through the states.
2. We repeat until the changes of utility values between consecutive sweeps are marginal.



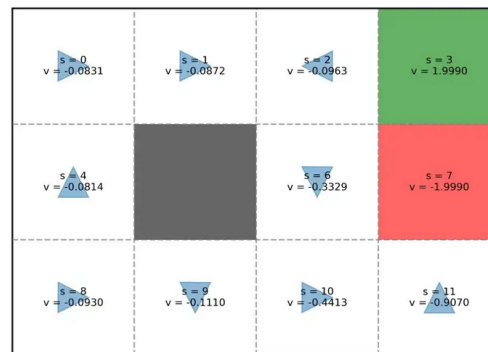
$$\begin{aligned}
 v(0) &= r(0) + \gamma [0.8v(1) + 0.1v(0) + 0.1v(4)] \\
 &= -0.04 + 0.5 \times [0.8 \times 0 + 0.1 \times 0 + 0.1 \times 0] = -0.04 \\
 v(1) &= r(1) + \gamma [0.8v(2) + 0.1v(1) + 0.1v(1)] \\
 &= -0.04 + 0.5 \times [0.8 \times 0 + 0.1 \times 0 + 0.1 \times 0] = -0.04 \\
 v(2) &= r(2) + \gamma [0.8v(1) + 0.1v(6) + 0.1v(2)] \\
 &= -0.04 + 0.5 \times [0.8 \times (-0.04) + 0.1 \times 0 + 0.1 \times 0] = -0.056 \\
 v(3) &= r(3) + \gamma [v(3)] \\
 &= 1 + 0.5 \times 0 = 1 \\
 v(4) &= r(4) + \gamma [0.8v(0) + 0.1v(4) + 0.1v(4)] \\
 &= -0.04 + 0.5 \times [0.8 \times (-0.04) + 0.1 \times 0 + 0.1 \times 0] = -0.056 \\
 v(6) &= r(6) + \gamma [0.8v(10) + 0.1v(6) + 0.1v(7)] \\
 &= -0.04 + 0.5 \times [0.8 \times 0 + 0.1 \times 0 + 0.1 \times 0] = -0.04 \\
 v(7) &= r(7) + \gamma [v(7)] \\
 &= -1 + 0.5 \times 0 = -1 \\
 v(8) &= r(8) + \gamma [0.8v(9) + 0.1v(4) + 0.1v(8)] \\
 &= -0.04 + 0.5 \times [0.8 \times 0 + 0.1 \times (-0.056) + 0.1 \times 0] = -0.0428 \\
 v(9) &= r(9) + \gamma [0.8v(9) + 0.1v(8) + 0.1v(10)] \\
 &= -0.04 + 0.5 \times [0.8 \times 0 + 0.1 \times (-0.0428) + 0.1 \times 0] = -0.04214 \\
 v(10) &= r(10) + \gamma [0.8v(11) + 0.1v(6) + 0.1v(10)] \\
 &= -0.04 + 0.5 \times [0.8 \times 0 + 0.1 \times (-0.04) + 0.1 \times 0] = -0.042 \\
 v(11) &= r(11) + \gamma [0.8v(7) + 0.1v(10) + 0.1v(11)] \\
 &= -0.04 + 0.5 \times [0.8 \times (-1) + 0.1 \times (-0.042) + 0.1 \times 0] = -0.4421
 \end{aligned}$$

Pseudo-Code

```
input : reward function  $r(s)$ , transitional model  $p(s'|s, a)$ ,  
        discounted factor  $\gamma$ , convergence threshold  $\theta$   
        policy  $\pi(s)$ , value  $v(s)$   
output: converged value  $v(s)$   
1 converge  $\leftarrow$  false  
2 while converge = false do  
3    $\Delta \leftarrow 0$   
4   for  $s \in S$  do  
5     temp  $\leftarrow v(s)$   
6      $v(s) \leftarrow r(s) + \gamma \sum_{s'} p(s'|s, a = \pi(s))v(s')$   
7      $\Delta \leftarrow \max(\Delta, |\text{temp} - v(s)|)$   
8   end  
9   if  $\Delta < \theta$  then  
10    converge  $\leftarrow$  true  
11  end  
12 end  
13 return  $v(s)$ 
```

After we get the utility for this policy, it is time to improve the policy!

- The utility represents how good a state is.
 - When choosing actions for a state, we should prefer successor states with higher utility.
 - For instance, if we look at state $s = 2$, potential successor states include 1, 3, and 6. Clearly $s = 3$ is the best choice as it has the highest utility of 1.999. Another example, if we look at state $s = 8$, potential successor states include 4 and 9 with utility of -0.0814 and -0.1110 respectively. As such, $s=4$ is preferable than $s=9$.



- Since our MDP is stochastic, selecting a preferable successor state does not guarantee we will reach it.
- Rather than successor states, what we should compare is actions. For state $s = 6$, we have four possible actions:
 1. UP : $0.8V(2) + 0.1V(6) + 0.1V(7) = -0.3093$
 2. LEFT : $0.8V(6) + 0.1V(10) + 0.1V(2) = -0.3201$
 3. DOWN : $0.8V(10) + 0.1V(6) + 0.1V(7) = -0.5862$
 4. RIGHT: $0.8V(7) + 0.1V(2) + 0.1V(10) = -1.6530$

Comparing the outcomes of these four possible actions, clearly UP is the best choice. Therefore, we should update the policy of the state $s = 6$ to UP.

We perform this process for each state:

$$\pi(s) = \underset{s'}{\operatorname{argmax}} \left[\sum P_{ss'}^a V(s') \right]$$

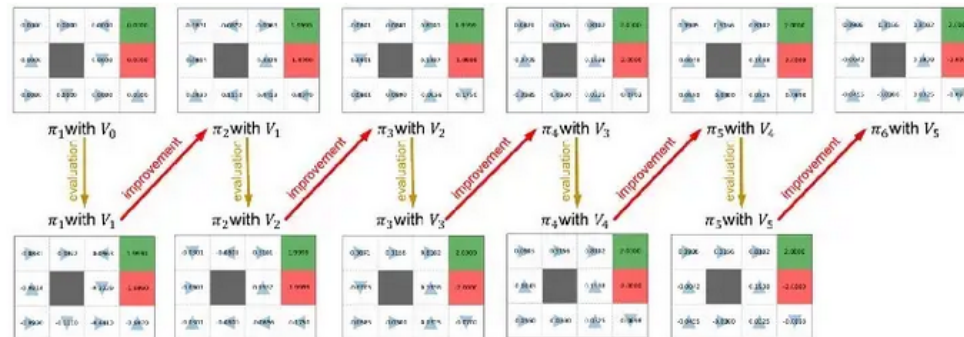
Policy Improvement

Pseudo-Code

```
input : transitional model  $p(s'|s, a)$ ,  
        policy  $\pi(s)$ , value  $v(s)$   
output: updated policy  $\pi(s)$ ,  
        binary indicating whether any change occurs  
1 change  $\leftarrow$  false  
2 for  $s \in S$  do  
3   | temp  $\leftarrow \pi(s)$   
4   |  $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s'} p(s'|s, a)v(s')$   
5   | if  $\text{temp} \neq \pi(s)$  then  
6   |   | change  $\leftarrow$  true  
7   | end  
8 end  
9 return  $\pi(s)$ , change
```

We are not done yet!

- We do get a policy better than our initial one, but it is not necessarily the best one?
- What we need is to repeat this whole process of **Policy Evaluation** and then **Policy Improvement** again and again.
- When we perform policy evaluation again, we use for utility the results of policy iteration.



Pseudo-Code

```
input : reward function  $r(s)$ , transitional model  $p(s'|s, a)$ ,  
         discounted factor  $\gamma$ , convergence threshold  $\theta$   
output: optimal policy  $\pi^*(s)$   
1 initialize  $\pi(s)$  randomly  
2 initialize  $v(s)$  with zeros  
3  $\text{stable} \leftarrow \text{false}$   
4 while  $\text{stable} = \text{false}$  do  
5    $v(s) \leftarrow \text{policy evaluation}(r(s), p(s'|s, a), \gamma, \theta, \pi(s), v(s))$   
6    $\pi(s), \text{change} \leftarrow \text{policy improvement}(p(s'|s, a), \pi(s), v(s))$   
7   if  $\text{change} = \text{false}$  then  
8      $\text{stable} \leftarrow \text{true}$   
9   end  
10 end  
11  $\pi^*(s) \leftarrow \pi(s)$   
12 return  $\pi^*(s)$ 
```

State-Value Function



- State value (How good is it to be in state s ?)

$$V^{\pi}(s) = E\left[\sum_{t \geq 0} \gamma^t r_t \mid s_t = s, \pi\right]$$

$$= \sum_a \pi(s, a) [R_{ss'}^a + \gamma \sum_{s'} P_{ss'}^a V^{\pi}(s')]$$

Expected cumulative reward from following the policies from state s

State	Value
State0	0.705
State5	0.762
State9	0.812
State2	0.655

0.812	0.868	0.918	+1
0.762		0.660	-1
0.705	0.655	0.611	0.388

Bellman Expectations Equations

- Basic: State-value function $V_{\pi}(s)$
 - Current state S
 - Multiple possible actions determined by stochastic policy $\pi(a|s)$
 - Each possible action is associated with a action-value function $Q_{\pi}(s, a)$ returning a value of that particular action
 - Multiplying the possible actions with the action-value function and summing them gives us an indication of how good it is to be in that state

$$V_{\pi}(s) = \sum_a \pi(a|s)Q(s, a)$$

Loose intuitive interpretation: **state-value = sum(policy determining actions * respective action-values)**

Action-Value Function

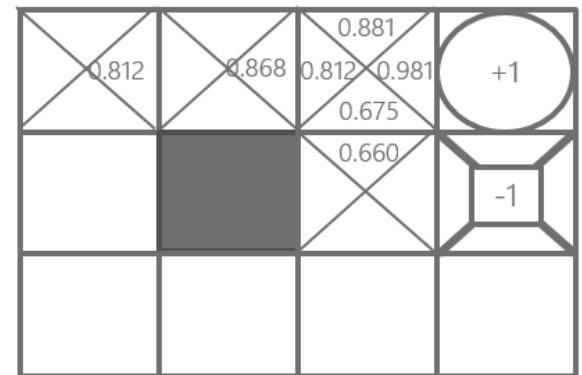
- Action value (How good is a (state, action) pair)

$$Q^{\pi}(s, a) = E\left[\sum_{t \geq 0} \gamma^t r_t \mid s_t = s, a_t = a, \pi\right]$$

$$= R_s^a + \gamma \sum_{s'} P_{ss'}^a \sum_{a'} \pi(a' | s') Q(s', a')$$

Expected cumulative reward from taking action a in state s and then following the policy.

State	Action	Value
State8	Action0 (←)	0.812
State8	Action1 (↑)	0.881
State8	Action2 (→)	0.981
State8	Action3 (↓)	0.675



Bellman Expectations Equations

- Basic: Action-value function $Q_{\pi}(s, a)$
- With a list of possible multiple actions, there is a list of possible subsequent states s' associated with:
 - state value function $V_{\pi}(s')$
 - transition probability function $P_{ss'}^a$, determining where the agent could land in based on the action
 - reward R_s^a for taking the action

$$Q_{\pi}(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V_{\pi}(s')$$

Loose intuitive interpretation: **state-value = sum(action-value = reward + sum(transition outcomes determining states * respective state-values))**

Value Iteration

$$Q^{\pi}(s, a) = \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma \sum_{a'} Q^{\pi}(s', a')]$$

1. Initialize the random value function (random value for each state)
2. Compute Q function for all state action pairs of Q(s,a)
3. Update value function with the max value from Q(s,a)
4. Repeat steps until the **change in the value function is very small**

1

State	Value
A	0 -> 0.3
B	0
C	0

3 max Q(s,a)

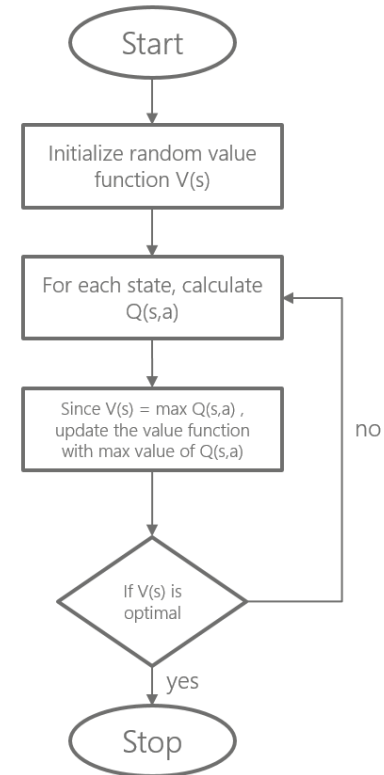
2

State	Action	Value
A	0	-0.1
A	1	0.3
B	0	

State	Action	Next State	TP	RP
A	0	A	0.1	0
A	0	B	0.4	-1.0
A	0	C	0.3	1.0
A	1	A	0.3	0
A	1	B	0.1	-2.0
A	1	C	0.5	1.0

$$Q(A,0) = (0.1 \cdot (0+0)) + (0.4 \cdot (-1.0+0)) + (0.3 \cdot (1.0+0)) = -0.1$$

$$Q(A,1) = (0.3 \cdot (0+0)) + (0.1 \cdot (-2.0+0)) + (0.5 \cdot (1.0+0)) = 0.3$$



31 / 34

- Both *value iteration* and *policy iteration* compute the same thing (all optimal values)
- In **policy iteration**:
 - We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
 - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
 - The new policy will be better (or we're done)
- In **value iteration**:
 - Every iteration updates both the values and (implicitly) the policy
 - We don't track the policy, but taking the max over actions implicitly recomputes it
- Both are dynamic programs for solving MDPs

Lessons learned today:

- **Model based:** Attempts to model the environment to find the best policy
- **Policy:** Evaluation, Improvement, and Optimization
- **State-value based:** Search for the optimal state-value function (goodness of action in the state)
- **Action-value based:** Search for the optimal action-value function (goodness of policy)

Code can be found in

- [RL_Policy_Optimization.ipynb](#)
- or on GitHub or [hosted on myBinder](#)



Exercise

https://inf-git.fh-rosenheim.de/aai-url/09_uebung

- Value and Policy Iteration by hand
- Questions?