

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

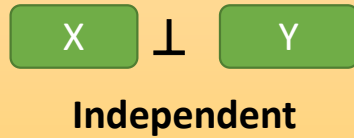
TP Reinforcement Learning

Adriana TAPUS and Juan Jose Garcia Cardinas and Adnan Saood

adriana.tapus@ensta-paris.fr & juan-jose.garcia@ensta.fr &
adnan.saood@ensta-paris.fr

1. Basic knowledge

■ Random Variable



■ Stochastic Process

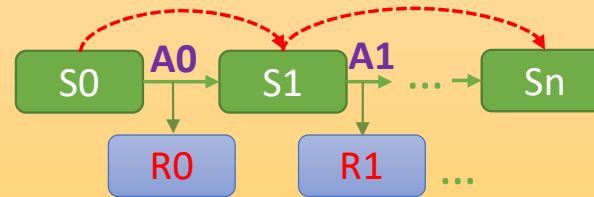


■ Markov Process/Chain

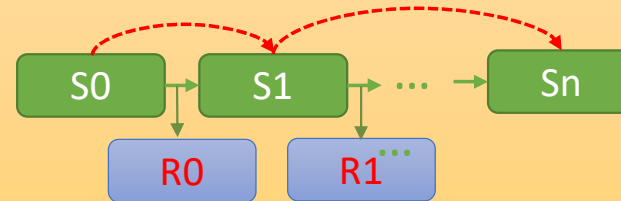


$$P(X_n | X_{n-1}) = P(X_n | X_0, X_1, \dots, X_{n-1}).$$

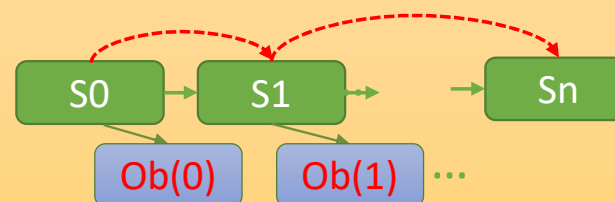
■ Markov Decision Process (Markov Chain+Reward+Action)



■ Markov Reward Process (Markov Chain+Reward)



■ State Space Model (Markov Chain+Observation)



e.g. HMM
e.g. Kalman Filter
e.g. Particle Filter

2. Terminologies

Return : cumulative future reward

$$U_t = R_t + R_{t+1} + R_{t+2} + \dots$$



R_t and R_{t+1} are important equally?

Discounted Return : discounted cumulative future reward

$$U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \dots$$

Action-value function $Q_\pi(s_t, a_t)$

$$Q_\pi(s_t, a_t) = \mathbb{E}[U_t | S_t = s_t, A_t = a_t]$$

State-value function (value used in task 2)

$$V_\pi(s_t) = \mathbb{E}_A[Q_\pi(s_t, A)] = \sum_a \pi(a|s_t) \cdot Q_\pi(s_t, a)$$

Reinforcement Learning (RL)

- Learning through Interaction with Environment
- Agent is in State s
- Agent executes Action a
- Agent receives a *Reward* $r(s,a)$ from the environment
- Goal: Maximize *long-term discounted Reward*

Value-Based RL

- Policy Iteration:
 - Start with random policy π_0
 - Estimate Value-Function of π_i
 - Improve $\pi_i \rightarrow \pi_{i+1}$ by making it greedy w.r.t. to the learned value function
 - Exploration: Try out random actions to explore the state-space
 - Repeat until Convergence
- Learning Algorithms:
 - Q-Learning (off-policy), SARSA (on-policy)
 - Actor-Critic Methods, etc.

3. Value iteration

Value iteration is a method of computing an optimal policy for an MDP (Markov Decision Process) and its value.

$$\begin{aligned} Q_{k+1}(s, a) &= R(s, a) + \gamma * \sum_{s'} P(s' | s, a) * V_k(s') \\ V_k(s) &= \max_a Q_k(s, a) \end{aligned}$$

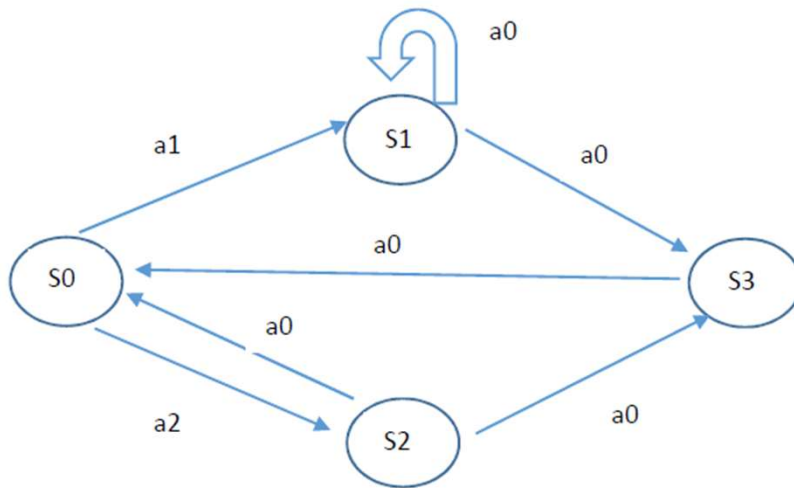
$$\pi[s] = \operatorname{argmax}_a R(s, a) + \gamma * \sum_{s'} P(s' | s, a) \cdot V_k[s']$$

[Interactive Example](https://perso.ensta-paris.fr/~saood/external/RL)

<https://perso.ensta-paris.fr/~saood/external/RL>

<https://artint.info/2e/html/ArtInt2e.Ch9.S5.SS2.html>

4. Task 1



In the figure above, the states are depicted by circles (S0, S1, S2, and S3) and the associated actions are indicated on the arrows: a0, a1, and a2. The transition functions for all the actions are shown below.

	S0	S1	S2	S3	s' (future)
S0	0	0	0	0	
S1	0	1-x	0	x	
S2	1-y	0	0	y	
S3	1	0	0	0	

S(Current) ↓

$$T(S, a0, S') = \begin{pmatrix} 0 & 0 & 1-y & 1 \\ 0 & 1-x & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$T(S, a1, S') = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$T(S, a2, S') = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

Each of the parameters x and y are in the interval $[0, 1]$, and the discounted factor $\gamma \in [0, 1]$

The reward is:

$$R(s) = \begin{cases} 10, & \text{for state } S3 \\ 1, & \text{for state } S2 \\ 0, & \text{otherwise} \end{cases}$$

4. Task 1

Question 1:

Enumerate all the possible policies

$$\pi: S \rightarrow a$$

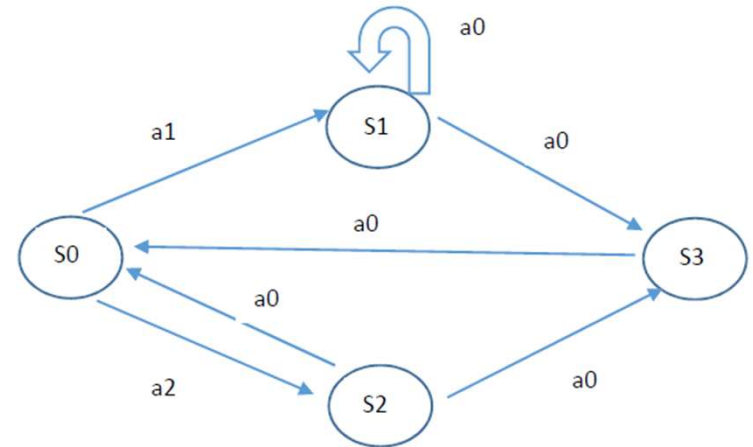
Question 2:

Write the equation for each optimal value function for each state

$$(V^*(s_0), V^*(s_1), V^*(s_2), V^*(s_3))$$

Reminder:

$$V^*(S) = R(s) + \max_a \gamma \sum_{S'} T(S, a, S') V^*(S')$$

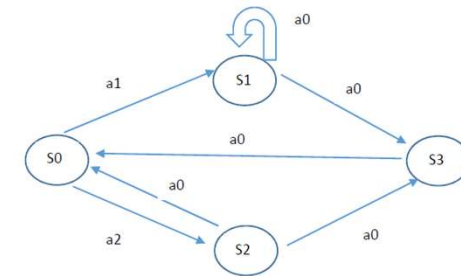


4. Task 1

Question 3:

Is there exist a value for x , that for all $\gamma \in [0,1)$, and $y \in [0,1]$, $\pi^*(s_0) = a_2$. Justify your answer.

Reminder:



$R(s) \rightarrow$ not function of " a "

$$\pi^*(s) = \arg \max_a \sum_{S'} T(S, a, S') V^*(S') \quad \leftarrow \quad V^*(S) = R(s) + \max_a \gamma \sum_{S'} T(S, a, S') V^*(S')$$

Question 4:

Is there exist a value for y , that for all $x > 0$, and $y \in [0,1]$, $\pi^*(s_0) = a_1$. Justify your answer.

Question 5: **\rightarrow python code**

Using $x=y=0.25$ and $\gamma = 0.9$, calculate the π^* and V^* for all states.

Implement value iteration.

$$\begin{aligned} Q_{k+1}(s, a) &= R(s, a) + \gamma \sum_{s'} P(s' | s, a) * V_k(s') \\ V_k(s) &= \max_a Q_k(s, a) \end{aligned}$$

Termination Rule: $|V_k(S) - V_{k-1}(S)| < 0.0001$

Rule

--In groups

--Deadline: before Tuesday 21/10

Submit:

-- Code + Readme file with Q1 to Q4

-- Github or ENSTA gitlab

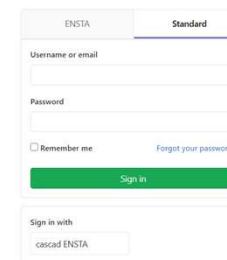


or

Bienvenue sur le serveur GitLab de DaTA, l'association d'informatique de l'ENSTA !



Hébergez vos dépôts Git simplement et en toute sécurité !

A screenshot of a GitLab login form. It has two tabs: "ENSTA" and "Standard". The "ENSTA" tab is selected. The form includes fields for "Username or email" and "Password", a "Remember me" checkbox, a "Forgot your password?" link, a green "Sign in" button, and a "Sign in with" section with a "cascad ENSTA" button.

End!
Questions?