

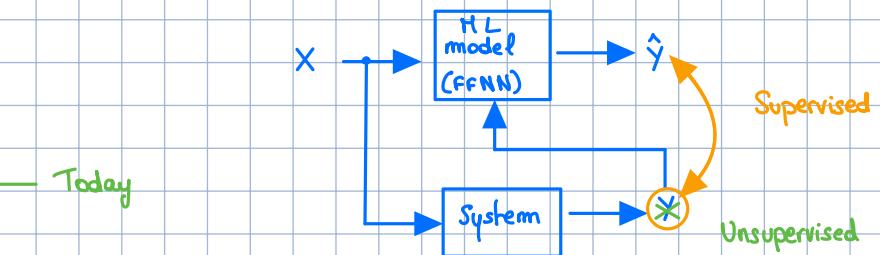
Reinforcement Learning

Summary

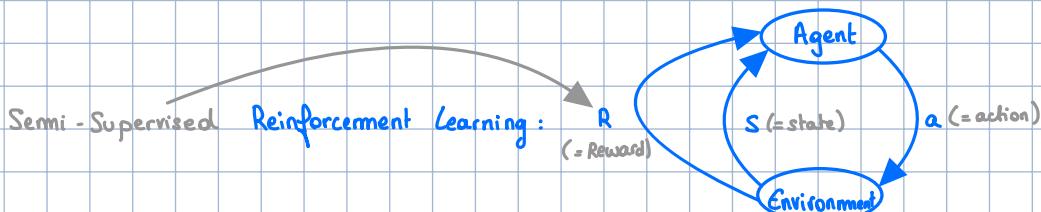
- MDP
- Policy function
- Value function
- Bellman equation
- Optimal value function
- Optimal policy function
- Value iteration
- Policy iteration
- Q - Learning

Overview of Machine Learning

HL {
 Supervised Learning
 Unsupervised Learning
 Reinforcement Learning



	Known	Unknown
Supervised Learning	X, Y	$\theta = (\nu, \omega) \rightarrow$ model parameters
Unsupervised Learning	X	Y, θ
Reinforcement Learning	X, —	a — actions



HDP

RL models the world (environment + agent) using the Markov Decision Process formalism HDP

Markov : whenever the future can be predicted knowing only the present

HDP is a 5-tuple $(S, A, \{P_{sa}\}, \gamma, R)$:

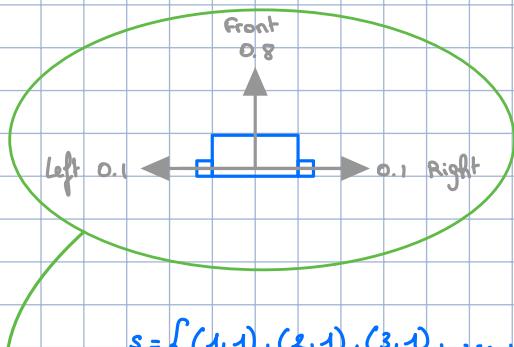
- $S \triangleq$ state space : all the values that a state can take.
- $A \triangleq$ action space : all the values that an action can take
- $\{P_{sa}\} \triangleq$ state transition distribution: probability of ending up at state s' from state s by taking action a

$$S \xrightarrow[a]{R} s', 0 \leq P_{sa}(s') \leq 1$$

$\sum_{s'} P_{sa}(s') = 1$

- $\gamma \triangleq$ discount factor : The importance of the future with respect to the present
 $0 \leq \gamma < 1$
 - $R : S \longrightarrow \mathbb{R}$ \triangleq Reward function

Example :



$$S = \{(1,1), (2,1), (3,1), \dots, (4,3)\}$$

$|S| = 11$

$$A = \{N, S, E, W\}$$

IAI

fPsaf

$$P(3,-1)_N((3,2)) = 0.8 \quad P(3,-1)_N((2,1)) = 0.1 \quad P(3,-1)_N((4,1)) = 0.1$$

$$P(3,-1)_N((3,3)) = 0 \quad P(3,-1)_W((2,1)) = 0$$

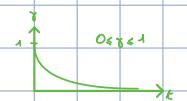
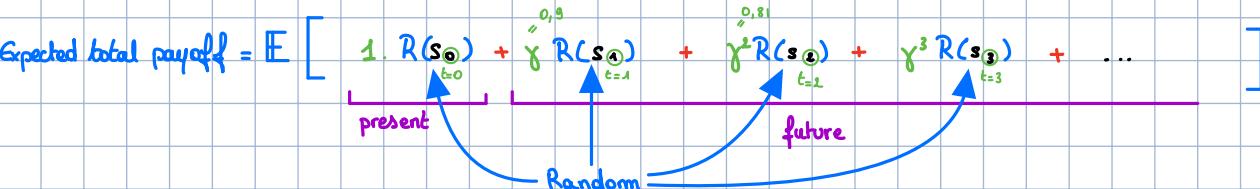
$\gamma = 0.99$ (future is very important)

$$R = \begin{cases} +1 & \text{if } s = (4, 3) \\ -1 & \text{if } s = (4, 2) \\ -0.02 & \text{otherwise (encourage robot to reach goal as quickly as possible + energy consumption)} \end{cases}$$

Stop condition: An episode finishes when the robot hits the celle $(4, 3)$ or $(4, 2)$.

Performance measure: How well the robot did by visiting the sequence of states s_0, s_1, s_2, \dots ?

- 1) Define the reward function
 - 2) Apply it to the sequence of states
 - 3) Compute the discounted cumulative reward or total payoff



Goal of RL

Choose actions over time (a_0, a_1, a_2, \dots) that maximize the expected total payoff
Now?

Define:

1) Policy function $\pi : S \rightarrow A$

It is a function that recommends which action one should take for a given state

2) State-value function $V^\pi : S \rightarrow \mathbb{R}$

For any given π , $V^\pi(s)$ is the expected total payoff starting from state s and following the policy π .

$$V^\pi(s_0) = \mathbb{E} \left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s, \pi \right]$$

Random

Knowing that

$$= R(s) + \gamma \mathbb{E} \left[R(s_1) + \gamma R(s_2) + \dots \mid \pi \right]$$

$$= R(s) + \gamma \sum_s P_{sa}(s') \mathbb{E} \left[R(s_1) + \gamma R(s_2) + \dots \mid s_1 = s', \pi \right]$$

$$\boxed{V^\pi(s) = R(s) + \gamma \sum_s P_{sa}(s') V^\pi(s')} \rightarrow \text{Bellman equation}$$

Example:

\rightarrow	\rightarrow	\rightarrow	$+1$
\downarrow		\rightarrow	-1
\rightarrow	\rightarrow	\uparrow	\uparrow

π

V^π

0.53	0.72	0.77	$+1$
-0.90		-0.82	-1
-0.88	-0.87	-0.85	-1

$(1, -1) \leftarrow 0$

$(2, 1) \leftarrow 1$

$(4, 3) \leftarrow -10$

$|s|$

\downarrow

$=$

\downarrow

Optimal value function

$$V^*(s) = \max_{\pi} V^\pi(s)$$

$$V^*(s) = R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V^*(s')$$

Bellman form
present future

Optimal policy function

$$\pi^*(s) = \arg\max_a \sum_{s'} P_{sa}(s') V^*(s')$$

Note: $\arg\max$ is the index of max.

$$z = \begin{bmatrix} 0 & 2 \\ 1 & 3 \\ 2 & 1 \end{bmatrix} \quad \max z = 3 \quad \arg\max z = 1$$

Goal of Reinforcement Learning!

Value iteration

Initialize $V(s) = 0 \quad \forall s$

Repeat until convergence { \rightarrow episodes

for every $s \{$

$V(s) = R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V(s')$ \rightarrow Bellman

}

}

$$\Rightarrow V \longrightarrow V^* \Rightarrow \pi^*(s) = \arg\max_a \sum_{s'} P_{sa}(s') V(s')$$

Policy iteration

Initialize π randomly $\forall s$

Repeat until convergence { \rightarrow episodes

$V = V^\pi = (I - \gamma A)^{-1} R$

for every $s \{$

$\pi(s) = \arg\max_a \sum_{s'} P_{sa}(s') V(s')$ \rightarrow Bellman

}

}

$$\Rightarrow V \longrightarrow V^* \\ \pi \longrightarrow \pi^*$$

Q-Learning \triangleq state-action value function

Initialize $Q(s, a) = 0 \quad \forall s, \forall a$

Repeat until convergence { \rightarrow episodes

Choose s_0 randomly

Repeat until stop condition is satisfied {

$m = \text{rand}()$

if ($m < \epsilon$):

choose a_t randomly

else:

$a_t = \arg\max_a Q(s_t, a)$

Obtain $(s_{t+1}, R(s_t))$ from (s_t, a_t)

$$Q(s_t, a_t) = (1-\alpha) Q(s_t, a_t) + \alpha [R(s_t) + \gamma \max_a Q(s_{t+1}, a)]$$

$$s_t = s_{t+1}$$

} choose action

$$\Rightarrow Q \longrightarrow Q^* \Rightarrow \pi^*(s) = \arg\max_a Q^*(s, a)$$

Exploration

Exploitation

s	N	S	E	W
O				
10				

En espérant que ce cours vous sera utile. Bon courage pour le prelab et le lab, plus que deux semaines avant la fin des cours !

GK