

Reinforcement Learning

Sequential Decisions, MDP and Policies

Erwan Le Pennec

`Erwan.Le-Pennec@polytechnique.edu`



M2DS - Reinforcement Learning – Fall 2024

Outline

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions
- 3 Prediction and Planning
- 4 Operations Research and Reinforcement Learning
- 5 Control
- 6 Survey
- 7 References

Decision or Decisions





Sequential Decision Setting

- In many (most?) settings, not a single decision but a sequence of decisions.
- Need to take into account the (not necessarily immediate) consequences of the sequence of decisions/actions rather than of each decisions.
- Different framework than supervised learning (no immediate feedback here) and unsupervised learning (well defined goal here).



Sequential Decision

Sequential Decision

- Sequence of action A_t as a response of an environment defined by a state S_t
- Feedback through a reward R_t

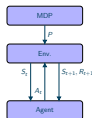
Actions?

- Is my current way of choosing actions good?
- How to make it better?

From Sequential Decision to Reinforcement Learning



Sequential Decision



MDP Modeling

Markov Decision Process Modeling

- Specific modeling of the environment.
- Goal as as a (weighted) sum of a scalar reward.

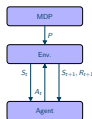
Actions?

- Is my current way of choosing actions good?
- How to make it better?

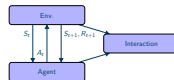
From Sequential Decision to Reinforcement Learning



Sequential Decision



MDP Modeling



Reinforcement Learning

Reinforcement Learning

- Same modeling. . .
- But no direct knowledge of the MDP.

Actions?

- Is my current way of choosing actions good?
- How to make it better?

Sequential Decisions

- MDP / Reinforcement Learning:

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t R_t \right]$$

- Optimal Control:

$$\min_u \mathbb{E} \left[\sum_t C(x_t, u_t) \right]$$

Related settings. . .

- (Stochastic) Search:

$$\max_{\theta} \mathbb{E}[F(\theta, W)]$$

- Online Regret:

$$\max \sum_k \mathbb{E}[F(\theta_k, W)]$$

References



R. Sutton and A. Barto.
Reinforcement Learning, an Introduction
(2nd ed.)

MIT Press, 2018



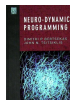
O. Sigaud and O. Buffet.
Markov Decision Processes in Artificial Intelligence.

Wiley, 2010



M. Puterman.
Markov Decision Processes. Discrete Stochastic Dynamic Programming.

Wiley, 2005



D. Bertsekas and J. Tsitsiklis.
Neuro-Dynamic Programming.
Athena Scientific, 1996



W. Powell.
Reinforcement Learning and Stochastic Optimization: A Unified Framework for Sequential Decisions.

Wiley, 2022



S. Meyn.
Control Systems and Reinforcement Learning.

Cambridge University Press, 2022



V. Borkar.
Stochastic Approximation: A Dynamical Systems Viewpoint.

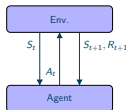
Springer, 2008



T. Lattimore and Cs. Szepesvári.
Bandit Algorithms.

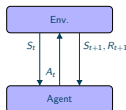
Cambridge University Press, 2020

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions
- 3 Prediction and Planning
- 4 Operations Research and Reinforcement Learning
- 5 Control
- 6 Survey
- 7 References



Decision Process and Environment

- At time step $t \in \mathbb{N}$:
 - State $S_t \in \mathcal{S}$: representation of the environment
 - Action $A_t \in \mathcal{A}(S_t)$: action chosen
 - Reward $R_{t+1} \in \mathcal{R}$: instantaneous reward
 - New state S_{t+1}
- Focus on the discrete setting, i.e. \mathcal{S} finite, $\mathcal{A}(s)$ finite and \mathcal{R} finite.
- Extension: Non finite bounded \mathcal{R} : easy / Non finite \mathcal{S} : hard / Non finite \mathcal{A} : harder.



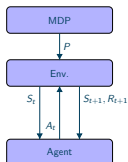
Stochastic Model

- Dynamic defined by:

$$\begin{aligned} \mathbb{P}(S_{t+1} = s', R_{t+1} = r | (S_{t'}, A_{t'}, R_{t'}), t' \leq t) \\ = \mathbb{P}(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a, H_t) \end{aligned}$$

where $H_t = (R_t, S_{t-1}, A_{t-1}, R_{t-1}, S_{t-2}, \dots)$ is the past and (S_t, A_t) the present.

$$(s', r) = F(s_t, a_t, H_t, w) \quad \text{random part}$$



Markovian Environment

- Markovian Dynamic Assumption: S_{t+1} and R_{t+1} are independent of the past $H_t = (R_t, S_{t-1}, A_{t-1}, R_{t-1}, S_{t-2}, \dots)$ conditionally to the present (S_t, A_t) .

- Dynamic entirely defined by state-reward transition probabilities

$$\begin{aligned}\mathbb{P}(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a, H_t) &= \mathbb{P}(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a) \\ &= p(s', r | s, a)\end{aligned}$$

in the discrete setting.

- Informally, this means that S_t encodes all the information related to the past.

- State-Reward transition probabilities for a given state-action:

$$\begin{aligned}\mathbb{P}(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a, H_t) &= \mathbb{P}(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a) \\ &= p(s', r | s, a)\end{aligned}$$

Induced State-action laws

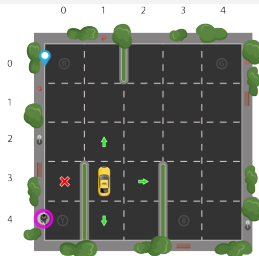
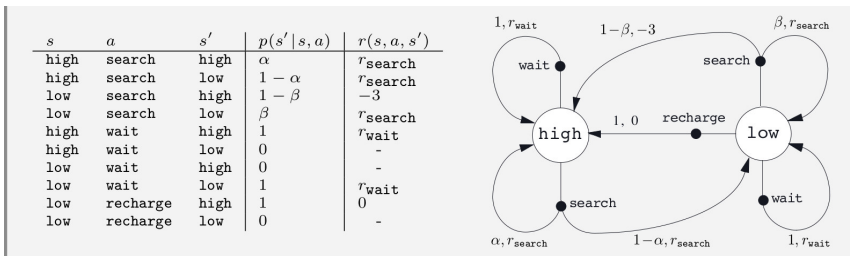
- State transition probabilities for a given state-action:

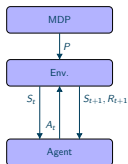
$$\begin{aligned}\mathbb{P}(S_{t+1} = s' | S_t = s, A_t = a, H_t) &= \mathbb{P}(S_{t+1} = s' | S_t = s, A_t = a) \\ &= p(s' | s, a) = \sum_r p(s', r | s, a)\end{aligned}$$

- Expected reward for a given state-action:

$$\begin{aligned}\mathbb{E}[R_{t+1} | S_t = s, A_t = a, H_t] &= \mathbb{E}[R_{t+1} | S_t = s, A_t = a] \\ &= r(s, a) = \sum_r r \sum_{s'} p(s', r | s, a)\end{aligned}$$

- From now on, we will always assume that the Markovian property holds for the environment.





Agent

- Interact with the environment by choose the action given the past.

Policy Π : specification of how to choose the action

- General stochastic policy $\Pi = (\pi_0, \pi_1, \dots, \pi_t, \dots)$:

$$\Pi_t(A_t = a) = \pi_t(A_t = a | S_t = s, A_t = a, H_t)$$

- General deterministic policy $\Pi = (\pi_0, \pi_1, \dots, \pi_t, \dots)$ (with as slight abuse of notation):

$$\Pi_t(A_t = a) = \mathbf{1}_{A_t = \pi_t(S_t = s, A_t = a, H_t)}$$

Agent

- Interact with the environment by choose the action given the past.

Policy Π : specification of how to choose the action

- History dependent stochastic policy $\Pi = (\pi_0, \pi_1, \dots, \pi_t, \dots)$:

$$\Pi_t(A_t = a) = \pi_t(A_t = a | S_t = s, H_t)$$

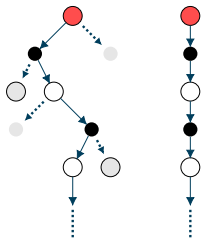
- Markovian stochastic policy $\Pi = (\pi_0, \pi_1, \dots, \pi_t, \dots)$:

$$\Pi_t(A_t = a) = \pi_t(A_t = a | S_t = s) = \pi_t(a | s)$$

- Stationary Markovian stochastic policy $\Pi = (\pi, \pi, \dots, \pi, \dots)$:

$$\Pi_t(A_t = a) = \pi(A_t = a | S_t = s) = \pi(a | s)$$

- Similar deterministic policy definition.
- Partially Observed Markov Decision Process extension: the Agent has only access to a partial observation O_t at each time step... (not the focus of the lectures)



Trajectories

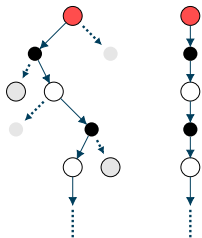
- Trajectory $(S_0, A_0, R_1, S_1, A_1, \dots)$ defined by

- an initial distribution \mathbb{P}_0 for S_0 ,
- a policy $\Pi = (\pi_0, \pi_1, \dots, \pi_t, \dots)$ specifying

$$\Pi_t(A_t = a) = \pi_t(A_t = a | S_t, H_t)$$

- an environment specifying

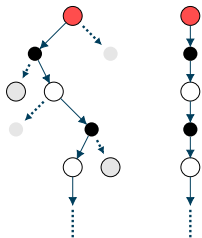
$$\mathbb{P}(S_{t+1}, R_{t+1} | S_t, A_t, H_t)$$



Trajectories

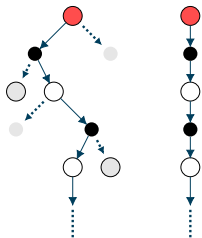
- Induced probability:

$$\begin{aligned} & \mathbb{P}(S_0 = s_0, A_0 = a_0, R_1 = r_1, S_1 = s_1, A_1 = a_1, \dots, S_t = s_t, R_t = r_t) \\ &= \mathbb{P}_0(S_0 = s_0) \\ & \quad \times \pi_0(A_0 = a_0 | S_0) \mathbb{P}(S_1, R_1 | S_0, A_0) \pi_1(A_1 = a_1 | S_1 = s_1, H_1) \\ & \quad \times \dots \times \mathbb{P}(S_t = s_t, R_t = r_t | S_{t-1} = s_{t-1}, A_{t-1} = a_{t-1}, H_{t-1}) \end{aligned}$$



Trajectories

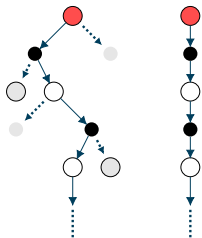
- Trajectory $(S_0, A_0, R_1, S_1, A_1, \dots)$ defined by
 - an initial distribution \mathbb{P}_0 for S_0 ,
 - a policy $\Pi = (\pi_0, \pi_1, \dots, \pi_t, \dots)$ specifying
$$\pi_t(A_t = a) = \pi_t(A_t = a | S_t, H_t)$$
 - a Markovian environment specifying
$$\mathbb{P}(S_{t+1}, R_{t+1} | S_t, A_t)$$



Trajectories

- Induced probability:

$$\begin{aligned} & \mathbb{P}(S_0 = s_0, A_0 = a_0, R_1 = r_1, S_1 = s_1, A_1 = a_1, \dots, S_t = s_t, R_t = r_t) \\ &= \mathbb{P}_0(S_0 = s_0) \\ & \quad \times \pi_0(A_0 = a_0 | S_0) \mathbb{P}(S_1, R_1 | S_0, A_0) \pi_1(A_1 = a_1 | S_1 = s_1, H_1) \\ & \quad \times \dots \times \mathbb{P}(S_t = s_t, R_t = r_t | S_{t-1} = s_{t-1}, A_{t-1} = a_{t-1}) \end{aligned}$$

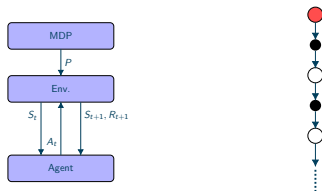


Markovian Trajectories only if the policy is Markovian

- $$\begin{aligned} & \mathbb{P}(R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \dots, R_{t+k}, S_{t+k} | S_t, A_t, H_t) \\ &= \mathbb{P}(R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \dots, R_{t+k}, S_{t+k} | S_t, A_t) \\ &= \mathbb{P}(S_{t+1}, R_{t+1} | S_t, A_t) \pi_{t+1}(A_{t+1} | S_{t+1}) \\ & \quad \times \dots \times \mathbb{P}(S_{t+k}, R_{t+k} | S_{t+k-1}, A_{t+k-1}) \end{aligned}$$

- Stationary if the policy is stationary.

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions**
- 3 Prediction and Planning
- 4 Operations Research and Reinforcement Learning
- 5 Control
- 6 Survey
- 7 References



Rewards and Total Returns

- MDP: Rewards R_t encode all the feedbacks!
- Quality of a policy Π measured from the remaining total return:

$$G_t = \sum_{t'=t+1}^{\infty} R_{t'}$$

- Expected total return following Π starting from s :

$$\mathbb{E}_{\Pi}[G_t | S_t = s] = \sum_{t'=t+1}^{\infty} \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

Issues

- G_t is a limiting process and thus may not be defined!
- Can diverge to $\pm\infty$ or not converge at all.

Fixes?

- Finite horizon: $G_t^T = \sum_{t'=t+1}^T R_{t'}$
- Episodic setting: it exists a random T such that $\forall t' \geq T, R_{t'} = 0$ and $\mathbb{E}[T] < \infty$
so that $G_t = \sum_{t'=t+1}^{\infty} R_{t'}$ is well defined.
- Discounted setting: for $0 < \gamma < 1$, $G_t^\gamma = \sum_{t'=t+1}^{\infty} \gamma^{t'-(t+1)} R_{t'}$
- Average return: $\bar{G}_t = \lim \frac{1}{T} \sum_{t'=t+1}^{t+T} R_{t'}$

$$G_t^T = \sum_{t'=t+1}^T R_{t'}$$

Finite Horizon Setting

- Always well defined and easy to interpret.
- Loss of Markovianity as we need to know the time step. . .
- Can be put in a classical Markov framework!
 - Define an absorbing state s_{abs} (a state that cannot be escaped and from which the reward is always 0).
 - Extend the state space \mathcal{S} to $(\mathcal{S} \times \{0, \dots, T\}) \cup \{s_{\text{abs}}\}$.
 - Define an state reward transition probability:

$$p(\tilde{s}', r | \tilde{s}, a) = \begin{cases} p(s', t | s, a) & \text{if } \tilde{s} = (s, t), t < T \text{ and } \tilde{s}' = (s', t+1) \\ 1 & \text{if } \tilde{s} = (s, t), t = T, \tilde{s}' = s_{\text{abs}} \text{ and } r = 0 \\ 1 & \text{if } \tilde{s} = s_{\text{abs}}, \tilde{s}' = s_{\text{abs}} \text{ and } r = 0 \\ 0 & \text{otherwise} \end{cases}$$

$$G_t = \sum_{t'=t+1}^{\infty} R_{t'}$$

Episodic Setting

- Assumption: for any policy Π , the average duration before R_t remains equal to 0 is smaller than a finite horizon H :
$$\mathbb{E}_{\Pi} \left[\min_{t, R_{t'}=0, \forall t' \geq t} t \right] \leq H < +\infty$$
- Strong assumption. . .
- Easy to interpret.
- Slightly stronger (but more convenient) def.:
 - Replace all the states from which R_t remains equal to 0 whatever the policy by a single absorbing state s_{abs} ,
 - Assumption: for any policy Π and any initial state, the average duration to reach this state is smaller than a finite horizon H :
$$\forall s, \mathbb{E}_{\Pi} \left[\min_{t, S_t=s_{\text{abs}}} t \middle| S_0 = s \right] \leq H < +\infty$$

$$G_t^\gamma = \sum_{t'=t+1}^T \gamma^{t'-(t+1)} R_{t'}$$

Discounted

- Always defined but not that easy to interpret.
- Easiest theoretical setting!
- Equivalent to an episodic setting if one adds an absorbing state s_{abs} and changes all state-reward transition probabilities to:

$$p(s', r|s, a) = \begin{cases} \gamma p(s', r|s, a) & \text{if } s' \neq s_{\text{abs}}, s \neq s_{\text{abs}} \\ (1 - \gamma) & \text{if } s' = s_{\text{abs}}, r = 0, s \neq s_{\text{abs}} \\ 1 & \text{if } s' = s_{\text{abs}}, r = 0, s = s_{\text{abs}} \\ 0 & \text{otherwise} \end{cases}$$

- Horizon $H = 1/(1 - \gamma)$.

$$\bar{G}_t = \lim \frac{1}{T} \sum_{t'=t+1}^{t+T} R_{t'}$$

Average Return

- Not always defined. (Cesaro Average)
 - Always equal to 0 in the episodic setting!
 - Natural definition in a *stationary* setting. . .
 - Complex theoretical analysis!
-
- Under a strict stationarity assumption ($R_t \sim R_{t'}$), link with discounted setting as

$$\mathbb{E}_{\Pi}[G_t^{\gamma}] = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\Pi}[R_{t+1}] = \frac{1}{1-\gamma} \mathbb{E}_{\Pi}[R_t] = \frac{1}{1-\gamma} \mathbb{E}_{\Pi}[\bar{G}_t]$$

State Value Functions

- Return expectation for a policy Π starting from s at time t

- Finite horizon setting:

$$v_{t,\Pi}^T(s) = \mathbb{E}_{\Pi}[G_t^T | S_t = s] = \sum_{t'=t+1}^T \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Episodic setting:

$$v_{t,\Pi}(s) = \mathbb{E}_{\Pi}[G_t | S_t = s] = \sum_{t'=t+1}^{\infty} \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Discounted:

$$v_{t,\Pi}^{\gamma}(s) = \mathbb{E}_{\Pi}[G_t^{\gamma} | S_t = s] = \sum_{t'=t+1}^{\infty} \gamma^{t'-(t+1)} \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Average return setting:

$$\bar{v}_{t,\Pi}(s) = \mathbb{E}_{\Pi}[\bar{G}_t | S_t = s] = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t'=t+1}^{t+T} \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Depends on t for a history dependent policy!

State Value Functions

- Return expectation for a Markovian policy Π starting from s at time t .
 - Finite horizon setting (with time extended state space):

$$v_{t,\Pi}^T(s) = \mathbb{E}_{\Pi}[G_t^T | S_t = s] = \sum_{t'=t+1}^T \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Episodic setting:

$$v_{t,\Pi}(s) = \mathbb{E}_{\Pi}[G_t | S_t = s] = \sum_{t'=t+1}^{\infty} \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Discounted:

$$v_{t,\Pi}^{\gamma}(s) = \mathbb{E}_{\Pi}[G_t^{\gamma} | S_t = s] = \sum_{t'=t+1}^{\infty} \gamma^{t'-(t+1)} \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Average return setting:

$$\bar{v}_{t,\Pi}(s) = \mathbb{E}_{\Pi}[\bar{G}_t | S_t = s] = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t'=t+1}^{t+T} \mathbb{E}_{\Pi}[R_{t'} | S_t = s]$$

- Becomes independent on t if the policy is stationary and Markovian the generic case (except in the finite horizon setting).

State-Action Value Functions

- Return expectation for a policy Π starting from s and an action a at time t .

- Finite horizon setting:

$$q_{t,\Pi}^T(s, a) = \mathbb{E}_{\Pi}[G_t^T | S_t = s, A_t = a] = \sum_{t'=t+1}^T \mathbb{E}_{\Pi}[R_{t'} | S_t = s, A_t = a]$$

- Episodic setting:

$$q_{t,\Pi}(s, a) = \mathbb{E}_{\Pi}[G_t | S_t = s, A_t = a] = \sum_{t'=t+1}^{\infty} \mathbb{E}_{\Pi}[R_{t'} | S_t = s, A_t = a]$$

- Discounted:

$$q_{t,\Pi}^{\gamma}(s, a) = \mathbb{E}_{\Pi}[G_t^{\gamma} | S_t = s, A_t = a] = \sum_{t'=t+1}^{\infty} \gamma^{t'-(t+1)} \mathbb{E}_{\Pi}[R_{t'} | S_t = s, A_t = a]$$

- Average return setting:

$$\bar{q}_{t,\Pi}(s, a) = \mathbb{E}_{\Pi}[\bar{G}_t | S_t = s, A_t = a] = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t'=t+1}^{t+T} \mathbb{E}_{\Pi}[R_{t'} | S_t = s, A_t = a]$$

- Different strategy for action at time t than after. . .
- Independent of t for a Markovian policy except for the finite horizon setting!



$$v_{t,\Pi}(s) = \mathbb{E}_{\Pi}[G_t | S_t = s] \quad q_{t,\Pi}(s, a) = \mathbb{E}_{\Pi}[G_t | S_t = s, A_t = a]$$

State vs State-Action

- Performance measure of a policy Π :
 - starting from a state s for the state value function,
 - starting from a state s and an action a (not necessarily related to Π) for the state-action value function.
- State value function at time t from state-action value function:

$$v_{t,\Pi}(s) = \sum_a \Pi_t(a) q_{t,\Pi}(s, a)$$

Equivalent Markovian policy in terms of value function

- **Thm:** For any policy Π and any initial distribution $\mathbb{P}_0(S_0)$, it exists a Markovian policy $\tilde{\Pi}$ such that

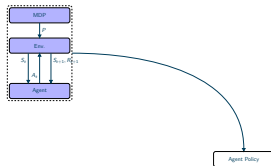
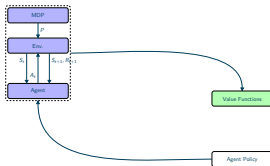
$$\forall t, \forall s, v_{t,\Pi}(s) = v_{t,\tilde{\Pi}}(s).$$

- Relies on the Markovian environment.
- Possible choice:

$$\tilde{\pi}_t \{A_t = a_t | S_t = s_t\} = \mathbb{E}_{\mathbb{P}, \mathbb{P}_0} [\pi_t(A_t = a_t | S_t = s_t, H_t) | S_t = s_t, S_0]$$

- **No need to consider non Markovian policy** if the goal is entirely defined in terms of value functions.

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions
- 3 Prediction and Planning**
- 4 Operations Research and Reinforcement Learning
- 5 Control
- 6 Survey
- 7 References

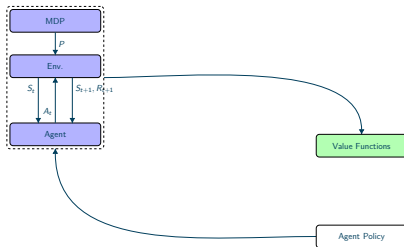


Prediction

- What is the performance of a given policy?
- Planning is harder than predicting.

Planning

- What is the *best* policy?

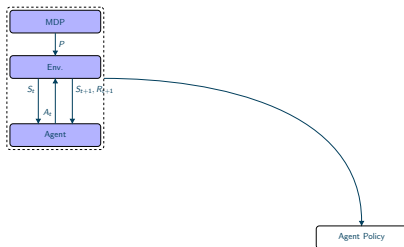


Prediction

- What is the performance of a given policy?
- Compute/Approximate/Estimate

$$v_{t,\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$$

- Well defined provided the expectation exists.

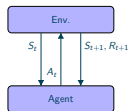
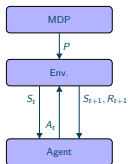


Planning

- What is the *best* policy?
- A possible definition: $\operatorname{argmax}_{\Pi} \sum_{s,t} \mu(s, t) v_{t,\Pi}(s)$
- Not necessarily well defined...
- Several choices for μ !
- More realistic goal: find a *good* policy...

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions
- 3 Prediction and Planning
- 4 Operations Research and Reinforcement Learning**
- 5 Control
- 6 Survey
- 7 References

What Do We Know?



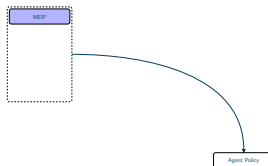
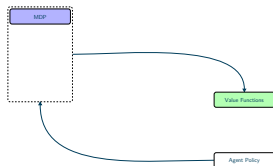
Model

- Able to use the MDP transition probabilities.
- Markov Decision Process / Operations Research.
- Probability world.

Only Observations

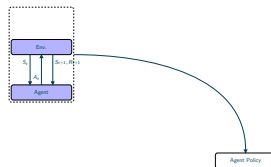
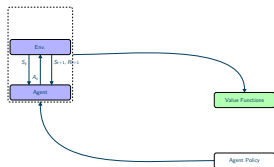
- No access to the MDP transition probabilities.
- Reinforcement Learning.
- Statistic world.

- Reinforcement Learning is harder than Markov Decision Process / Operations Research.



MDP / OR

- Stochastic setting in which the world is known.
- MDP model assumption.
- Probability world / Idealized setting. . .
- Lots of insight for the RL problem.



RL

- Stochastic setting in which the world is observed through interactions.
- Still MDP model assumption.
- More realistic setting?
- More difficult theoretical analysis.

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions
- 3 Prediction and Planning
- 4 Operations Research and Reinforcement Learning
- 5 Control**
- 6 Survey
- 7 References

MDP

- State s and action a
- Dynamic model:
$$\mathbb{P}(s'|s, a)$$
- Reward r defined by $\mathbb{P}(r|s', s, a)$.
- Policy Π : $a_t = \pi_t(S_t, H_t)$
- Goal:

$$\max \mathbb{E}_{\Pi} \left[\sum_t R_t \right]$$

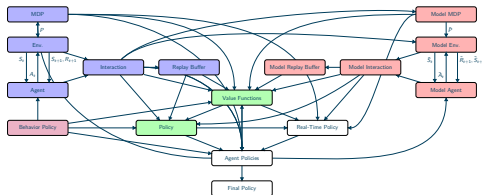
Discrete Control

- State x and control u
- Dynamic model:
$$x' = f(x, u, W)$$
with W a stochastic perturbation.
- Cost: $C(x, u, W)$.
- Control strategy U : $u_t = u(x_t, H_t)$
- Goal:

$$\min_U \mathbb{E}_U \left[\sum_t C(x_t, u_t, W_t) \right]$$

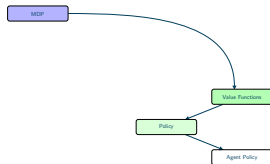
- Almost the same setting but with a different vocabulary!

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions
- 3 Prediction and Planning
- 4 Operations Research and Reinforcement Learning
- 5 Control
- 6 Survey**
- 7 References



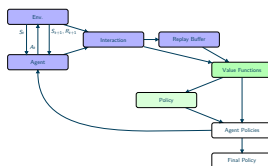
Outline

- Operations Research and MDP.
- Reinforcement learning and interactions.
- More tabular reinforcement learning.
- Reinforcement and approximation of value functions.
- Actor/Critic: a Policy Point of View



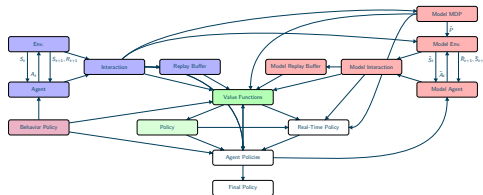
How to find the best policy knowing the MDP?

- Is there an optimal policy?
- How to estimate it numerically?
- Finite states/actions space assumption (tabular setting).
- Focus on iterative methods using value functions (dynamic programming).
- Policy deduced by a statewise optimization of the value function over the actions.
- Focus on the discounted setting.



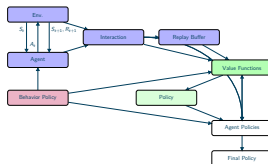
How to find the best policy not knowing the MDP?

- How to interact with the environment to learn a good policy?
 - Can we use a Monte Carlo strategy outside the episodic setting?
 - How to update value functions after each interaction?
-
- Focus on stochastic methods using tabular value functions (Q learning, SARSA...)
 - Policy deduced by a statewise optimization of the value function over the actions.



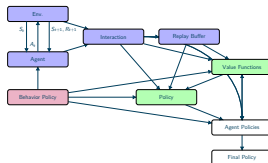
Can We Do Better?

- Is there a gain to wait more than one step before updating?
 - Can we interact with a different policy than the one we are estimating?
 - Can we use an estimated model to plan?
 - Can we plan in real-time instead of having to do it beforehand?
-
- Finite states/actions space setting (tabular setting).



How to Deal with a Large/Infinite states/action space?

- How to approximate value functions?
- How to estimate good approximation of value functions?
- Finite action space setting.
- Stochastic algorithm (Deep Q Learning. . .).
- Policy deduced by a statewise optimization of the value function over the actions.



Could We Directly Parameterized the Policy?

- How to parameterize a policy?
- How to optimize this policy?
- Can we combine parametric policy and approximated value function?
- State Of The Art Algorithms (DPG, PPO, SAC...)

- 1 Decision Process and Markov Decision Process
- 2 Returns and Value Functions
- 3 Prediction and Planning
- 4 Operations Research and Reinforcement Learning
- 5 Control
- 6 Survey
- 7 References**



R. Sutton and A. Barto.
Reinforcement Learning, an Introduction
(2nd ed.)

MIT Press, 2018



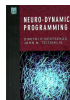
O. Sigaud and O. Buffet.
Markov Decision Processes in Artificial Intelligence.

Wiley, 2010



M. Puterman.
Markov Decision Processes. Discrete Stochastic Dynamic Programming.

Wiley, 2005



D. Bertsekas and J. Tsitsiklis.
Neuro-Dynamic Programming.

Athena Scientific, 1996



W. Powell.
Reinforcement Learning and Stochastic Optimization: A Unified Framework for Sequential Decisions.

Wiley, 2022



S. Meyn.
Control Systems and Reinforcement Learning.

Cambridge University Press, 2022



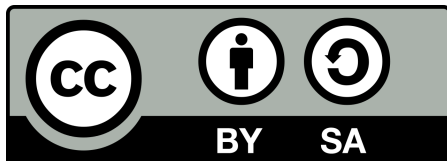
V. Borkar.
Stochastic Approximation: A Dynamical Systems Viewpoint.

Springer, 2008



T. Lattimore and Cs. Szepesvári.
Bandit Algorithms.

Cambridge University Press, 2020



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.