CS 59000_RL
Markov Decision Processes (MDPs)
  - Agenda:
        - settings.
        - Policy classes
        - Value
        - Dynamic programming

Examples: - Game of Casino
          - Game of racing
          - Continuous time Control system (plane)

MDPs can be continuous time, like the plane example,
They can be discrete time, like the casino game.
Or, time domain, where the process is defined can be
general normed spaces (it need no be time any way)

For simplicity, we study discrete time MDPs in this class.

For the game of racing, we considered four possible states,
and three possible actions. We could consider the case
that states and actions are continuous, or some what
general mensurable spaces (e.g. Polish space)
In the casino night setting, both state and action
spaces are finite spaces and categorical.

In this class we study tabular $MDP_s$ where
both state and action spaces are finite sets.
(Later we also talk about general ones)

---

MDP in a generic case is a tuple

$$M = (\mathcal{X}, \mathcal{A}, P, R, P_1)$$

Initial state distribution

$\to$ Reward Kernel

Transition Kernel

state space    action space

Protocol: At the first time step, $X_1$ is drawn

from $P_1$. Observing $X_1$, i.e, the history up to time

$t: h_1 := (X_1)$, the agent (the decision maker)

makes decision by choosing an action $A_1$,

and receives $r_1 \sim R_f(h_1, A_1)$

Given $h_1$, and $A_1$, the environment succeed to new

state $X_2$.

*Wednesday, October 7, 2020*

At each time step $t$, the environment is at

some state $X_t$, given $h_t : (X_1, A_1 - X_t)$, the agent

makes a decision $A_t$ and receives reward

$$r_t \sim R_1(h_t, A_t)$$

and the environment succeeds to state

$$\underline{X_{t+1}} \sim \underline{P_t(h_t, A_t)}$$

(The reward could also be $r_t \sim R(h_{t+1})$ )

MDP is a controlled Morkov Process. Therefore

$$P_t(h_t, A_t) = P_t(X_t, A_t)$$

In other words, $P_t(X_{t+1} | h_t, A_t) = P_t(X_{t+1} | X_t, A_t)$

similarly for the reward : $R(h_t, A_t) = R(X_t, A_t)$

$$r_t \sim R(X_t, A_t)$$

How the agent makes decision.

Policy sets:

- History dependent and randomized policy: $\pi \in \Pi^{HR}$

$$A_t \sim \pi(h_t)$$

- History dependent and deterministic: $\pi \in \Pi^{HD}$

$$A_t = \pi(h_t)$$

- Markov and randomized: $\pi \in \Pi^{MR}$

$$A_t \sim \pi(x_t, t)$$

- Markov and deterministic: $\pi \in \Pi^{MD}$

(Memory less policy is a special case of Markov policies, where $A_t \sim \pi(x_t)$ which is stationary.)

$$A = \pi(x_t)$$

(Note: Be careful about the notation since we used $\pi$ for anything.)

which policy class is the largest? $\Pi^{HR}$

*Wednesday, October 7, 2020*

Let's consider the case that we are interested in policies that provide us with large expected return.

Let's mak it concrebe.

---

− Infinite Horizon MDP

− Expected return : $\eta^{\Pi} = E^{\Pi}\left[\sum_{t=1}^{\infty} r_t\right]$
(Undiscounted reward)

Example: Collecting bib coin.

− Expected average return: $\eta^{\Pi} = \lim_{T \to \infty} \frac{1}{T} E^{\Pi}\left[\sum_{t=1}^{T} r_t\right]$

(A robot in assembling line keeps working for every

− Expected discounted return:

which one you choose ?
$\longrightarrow$ − 1 grand today
− 1 grand in a year
− 1 grand and 5 buchs in two years.

$$\eta^{\Pi} = E^{\Pi}\left[\sum_{t=1}^{\infty} \lambda^{t-1} r_t\right] \quad 0 < \lambda < 1$$

Future reward worth less to us right now

Episodic:

- Expeted return (undiscounted)

$$-\eta^{\Pi} = E\left[\sum_{t=1}^{\tau} r_t\right]$$

where $\tau$ is a termination (stopping) time, and is a random variable such that

$\tau < \infty$ a.s., i.e. the process terminates in finite time.

– Discounted reward

$$\eta^{\Pi} = E^{\Pi}\left[\sum_{t=1}^{\tau} \lambda^{t-1} r_t\right]$$

Note; there is a relationship between Episodic expected return setting and Infinite horizon discounted setting

For the latter $\eta^{\Pi} = E^{\Pi}\left[\sum_{t=1}^{\infty} \lambda^{t-1} r_t\right]$

Now, consider the case, where we terminate the process at time $t$ with probability $(1-\lambda)\lambda^{t-1}$.

$$\Rightarrow \eta^{\Pi} = E^{\Pi}\left[\sum_{t \geq 1} r_t\right]$$

$$= E^{\Pi}\left[\sum_{\tau=1}^{\infty}\left(\sum_{t=1}^{\tau} r(X_t, A_t)\right)(1-\lambda)\lambda^{\tau-1}\right]$$

if sumable and integrable
$$= E^{\Pi}\left[\sum_{t=1}^{\infty}\sum_{\tau=t}^{\infty} r(X_t, A_t)\underbrace{(1-\lambda)}\underbrace{\lambda^{\tau-1}}\right]$$

$$= E^{\Pi}\left[\sum_{t=1}^{\infty} \lambda^{t-1} r(X_t, A_t)\right]$$

which is the expected discounted reward in infinite horizon setting.

Finite horizon MDP: The set of finite horizon MDPs is a subset of episodic

MDPs, where, there exists a finite number

$H < \infty$, such that $T \leq H$

Fixed Horizon: is a subset of episodic MDPs where $T = H$ . for $1 \leq H < \infty$

Martin Putterman

MDP