# CMPUT 365: Markov Decision Processes

Rupam Mahmood

Jan 17, 2022

# Admin

Due dates C1M2:

- Practice quizzes: Tues Jan 18
- Peer-graded assignment: Thu Jan 20
- Peer-review: Sat Jan 22

Assignment 1:

- Will be released soon (This Sat?)
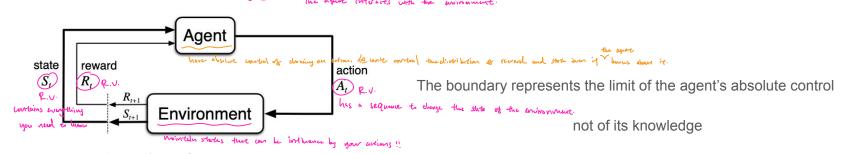- One week time
- Two worksheet(-like) questions

Midterm:

- Based on worksheet questions, book reading, and lectures

# Bandit review

- Policy  $\pi(a) \doteq \Pr(A = a)$.

  $\pi(a) \doteq \Pr(A = a \mid \pi)$

  $\pi_1(\cdot), \pi_1 = \Pr(A = a \mid \pi_1)$

  $\pi_2 = \Pr(A = a \mid \pi_2)$

  $\Pr\{\quad\}$

- Action value  $q_*(a) \doteq E[R \mid A=a] = \sum_r r \, P(R=r \mid A=a)$

  The pro of observing a reward of choosing an action.

  - Does not depend on the policy in bandits

- Value  $\Rightarrow$ goal: maximize the expect rewards $q_*(a)$.  $E_\pi[R] = \sum_r r \, P(R=r \mid \pi)$

  $= \sum_r r \sum_A P(R=r \mid A=a) \left(P(A=a \mid \pi)\right)$   $\rightarrow$ under the control of the agent.

  $\underbrace{P(R=r, A=a)}$

  $P(R=r) = \sum_A P(R=r, A=a)$.

  $= \sum_A P(A=a \mid \pi) \sum_r r \, P(R=r \mid A=a)$

  $\underbrace{\quad}_{\pi(a).}$  $\underbrace{\quad}_{q_*(a).}$

  $= \sum_a \pi(a) \cdot q_*(a)$

  ( $\therefore \exists \pi_1, \pi_2 \; E_{\pi_1}[R] < E_{\pi_2}[R]$ )

  - Depends on the policy even in bandits

- Contextual bandit

  $\pi(a \mid s) = P(A=a \mid X=x)$

  the particular context.

  (under this particular situation, what's the pro of the agent to choose action A?)

  (现实的 state).

  $E[R \mid A=a, X=x]$

  $= q_*(x, a)$

# Coursera: video 1 - Intro to MDPs

- The dynamics of an MDP

The agent interacts with the environment.



The boundary represents the limit of the agent's absolute control

not of its knowledge

have absolute control of choosing an action. No sure control the distribution of reward and state even if the agent knows about it.

$S_t$ R.V. contains everything you need to know

$R_t$ R.V.

$A_t$ R.V. has a sequence to change the state of the environment.

Maintain states that can be influenced by your actions !!

time steps, $t = 0, 1, 2, 3, \ldots$

Sequence or a trajectory
$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \ldots$

Dynamics of the MDP   The next state and reward depends on the previous state and action.
$$p(s', r \mid s, a) \doteq \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$$

dynamics function $p : \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \to [0, 1]$

$$\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r \mid s, a) = 1, \text{ for all } s \in \mathcal{S}, a \in \mathcal{A}(s)$$

David Wang to Everyone          1:48 PM

If I may. In layman's terms, what the dynamics means is the probability of moving to a specific following state and getting a specific reward, given the state you were in previously and the action you chose. Yes?   Yes!!

# Intro to MDPs (cont'd)

→ Probability of setting $s'$ as the next state given 当前状态为 $s$ 当前 action $a$.

*state-transition probabilities*

The boundary between the agent and the environment 是可变的.

a three-argument function $p : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \to [0, 1])$

$$p(s'|s,a) \doteq \Pr\{S_t = s' \mid S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathcal{R}} p(s', r \mid s, a)$$

doesn't depend on the policy.

only for the next state base on the current state and action.

Parham Golestaneh Talaei to Everyone 1:51 PM
is the first line the expected reward of choosing an action leading to state $s'$ ...

<- according to which law? ⇒ marginal probabilities
$P(R = r) = \sum_a P(R = r, A = a)$

The likelihood of observing different form $P(S'|S,a)$, 当前是 next state $s'$ and reward given the current state and the action.

*expected rewards for state–action pairs*

define a new object

$$r(s,a) \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r \mid s, a)$$

$\doteq$: mathematically equality

previous state     previous action

$\doteq$: this two expression equality

$A \cap B$: A 与 B 的 intersection

$= \sum_r r \, P(R_t = r \mid S_{t-1} = s, A_{t-1} = a)$  ~对s

$= \sum_r r \sum_{s'} p(s', r \mid s, a)$  对s

<- according to which law?

*expected rewards for state–action–next-state triples*

inputs (are given)

$$r(s,a,s') \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathcal{R}} r \frac{p(s', r \mid s, a)}{p(s' \mid s, a)}$$

$= \sum_r r \, P(R_t = r \mid S_{t-1} = s, A_{t-1} = a, S_t = s') =$

<- according to which law?

$Pr(A|B) = \frac{Pr(A \cap B)}{Pr(B)}$

Probability of A given B.

The joint probability of A and B (A和B同时发生)

# Coursera: video 2 - Examples of an MDP

- **Task**: the goal of the robot is to pick and place object

- **State**: latest readings of joint angles and velocities

- **Action**: the amount of voltage applied to each vector

- **Reward**: -c every time step, where c > 0 => every timestep you get a negative reward.
  when the reward is all negative, ie force the robot to terminate the process as soon as possible.

- **Termination**: when an object is placed successfully

# Coursera: video 3 - The goal of RL
# Video 5: Continuing tasks

Return in an episodic task, where an episode ends in a terminal state and T is the terminal step at that episode

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T$$

*=) goal: maximize the sum of future rewards. (everything matters).*

*reward*

Return in a continuing task *在 discount rate γ*

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$$

*goes to infinite.*

where $\gamma$ is a parameter, $0 \le \gamma \le 1$, called the *discount rate*.

Return expressed recursively

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots$$
$$= R_{t+1} + \gamma \left( R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \cdots \right)$$
$$= R_{t+1} + \gamma G_{t+1}$$

- Short-sighted agent vs. Far-sighted agent

*γ 接近 0.*      *γ 接近 1.*

Cameron Jen to Everyone    1:49 PM

CJ   So it's best for gamma to be 1 for continuing tasks to ensure we are optimizing all steps ?

David Wang to Everyone    1:49 PM

DW   No. If it is one, then the thing doesn't converge.

# Coursera: video 4 - The reward hypothesis

That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward).

# Coursera: video 6 - Examples of Episodic and Continuing Tasks

- Examples?

## Notations

r.v.  $X, Y, A, R, S$

outcomes  $x, y, a, r, s$

大写 $P$ → probability distribution
$P(event)$     output $\in [0,1]$
take an event as input     →output一个0~1之间的数字。(概率∈0,1)

大写 $E$ → expectation
input: $r.v.$   output $\in \mathbb{R}$
take an random variable as input     output一个real number

$P(A=a \mid S=s)$   conditional probability: take 2 events as input，输出 output $a$ given $s$ 的概率是多少概率。

$MDP \rightarrow$ 大写 $P\left(s' \mid s,a\right) = \underline{P(S_{t+1}=s' \mid S_t=s, A_t=a)}$  the probability of next state given the current state and current action.
take three outcomes as inputs.

$P是$ 的变量初值
~后面考虑 仪式 MDP，但$MDP_2 \rightarrow r(s' \mid s,a)$
是，$s,a \rightarrow r$ 决定的MDP

for action policy:
$\left.\begin{array}{r} \pi(a \mid s) \\ \mu(a \mid s) \end{array}\right\}$ 我们要一种 policy.
我们要找到 best policy.

# Notations (2)

$p \left( r \mid s, a \right)$ :> probability of reward given current state and action [this, current]

$= \sum_{s'} p(s', r \mid s, a)$ by law of total probabilities.

$p \left( s', r \mid s, a \right)$ :> probability of next state and this current reward given current state and action

Expected reward for $s, a$ [state - action pair]

$r(s, a) = \sum_r r \; p(r \mid s, a)$ :> by law of total probabilities.

$s, a$ are not implicated. :> 이곱은 방정식과 $r(s, a)$이를 input.

[left margin notes, partly in Korean]
확률은 어떻게 policy 기, 요가 behaviour [pink]
저거가 7개.

Probability of next state given the current state.

$p_\pi \left( s' \mid s \right) = \sum_a p_\pi \left( s', a \mid s \right)$

[pink] 저거가 현재의 어떤 확률을 구한 것도.

$= \sum_a p(s' \mid s, a) \; \pi(a \mid s)$ an action probability, which comes up to be policy.

currently denotes by $\pi$ (policy) [green]

:> depend on our policy!! if the policy is different, then this value will be different!!

$p(s' \mid s) = \sum_a p_\pi(s', a \mid s)$

$= \sum_a p(s' \mid s, a) \pi(a \mid s)$

$p \left( r \mid s \right)$ is this independent of the policy $\pi$? :> Yes.

probability of taking a particular reward r given current state.

# Bellman equation for $v_\pi$

Bellman function: express the value function of the given state in terms of the value of more states.

需要 more states reward
特殊 condition.
∴会 bring in new R.V.

LOTE
LOTUS
MP — Markov property: everything happen just depends on the state given.
LE

$$v_\pi(s) = E_\pi \left[ G_t \mid S_t = s \right]$$

depend on policy $\pi$.
use law of total expectation.

$$= E_\pi \left[ E_\pi \left[ G_t \mid S_t = s, A_t \right] \mid S_t = s \right] \Rightarrow \text{by law } E[E[X \mid Y]] = E[X].$$

the probability of this R.V. taken value.

$$= \sum_a E_\pi \left[ G_t \mid S_t = s, A_t = a \right] \pi(a \mid s)$$

def of action value.
∴所有动作 $q_\pi(s,a)$
$P(A_t = a \mid S_t = s)$ — the policy probability. $\pi(a \mid s)$

$$= \sum_a \pi(a \mid s) \, q_\pi(s, a)$$

the relation between state value and action value.

引入新变量

$$q_\pi(s, a) = E_\pi \left[ G_t \mid S_t = s, A_t = a \right]$$

LOTE

$$= E_\pi \left[ E_\pi \left[ G_t \mid S_t = s, A_t = a, R_t, S_{t+1} \right] \mid S_t = s, A_t = a \right]$$

by law $E[E[X \mid Y]] = E[X]$.

LOTUS

$$= \sum_{s', r} E_\pi \left[ G_t \mid S_t = s, A_t = a, R_{t+1}=r, S_{t+1}=s' \right] \, P(s', r \mid S_t = s, A_t = a)$$

state-action pair.

write the probability between of this R.V. taken value.

$P(s', r \mid s, a)$  Probability of next state $S_{t+1}=s'$ and current rewards $R_{t+1}=r$

$$= \sum_{s', r} P(s', r \mid s, a) \, E_\pi \left[ R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a, R_{t+1}=r, S_{t+1}=s' \right]$$

∵ $R_{t+1}$ already given.
可忽略之无关.

use linearity of expectation

$$= \sum_{s', r} P(s', r \mid s, a)\left[ r + \gamma \; \underset{\pi}{E}\left[ G_{t+1} \mid S_{t+1} = s' \right] \right]$$

$\rightarrow$ state value of the next state. 即 $V_{\pi}(s')$

by Markov Property. 依赖当前状态与它对应的 state 有关!!
∴ 只与 state 对应.

$$= \sum_{s', r} P(s', r \mid s, a)\left[ r + \gamma \; U_{\pi}(s') \right]$$

注!! don't write policy here!! 因为 the transitions dynamics of the environment chanis depend on policy!!
即 不出现 $\pi$!!

$$\boxed{V_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{s', r} P(s', r \mid s, a)\left[ r + \gamma \, U_{\pi}(s') \right]}$$
Bellman eq.

注: 因为 value functions 's depend on policy $\pi_i$, 即所以 policy 下都 should appear 出现 at the R.H.S. as well!!

S 的所有都需考虑!!
(即 $s$ is given!)

store value also depends on the policy.

有出现了 $a$ 与 $s'$ 这两种 extra 的额外的 variables 都 should be bounded by the summation!! 对应 $\sum_a$ 与 $\sum_{s', r}$ !!

- 因为 $\sum_{a} \pi(a \mid s)$ 中, 因该有新的为 input 在该式里(即式) L.H.S 里面, 所以 $V_{\pi}(s)$ 式 bounded by $\sum_{a}$ !!

- 在 $\sum_{s', r} P(s', r \mid c, a)$ 中, $c$, $r$ 没有对应input在该式里(即L.H.S)
所以, $V_{\pi}(s)$ 式 bounded by $\sum_{s', r}$ !!