

# CMPUT 365: Value Functions and Bellman Equations

Rupam Mahmood

Jan 24, 2022

## Admin

### Due dates C1M3:

- Practice quiz: Tues Jan 25
- Graded quiz: Sat Jan 29

### Assignment 1:

- was released on last Sun Jan 23
- Two worksheet-like questions
- Due on this Sun Jan 30

### Midterm:

- Based on worksheet questions, book reading, and lectures

# Coursera video: Specifying policies

- Definition

Formally, a *policy* is a mapping from states to probabilities of selecting each possible action. If the agent is following policy  $\pi$  at time  $t$ , then  $\pi(a|s)$  is the probability that  $A_t = a$  if  $S_t = s$ . Like  $p$ ,  $\pi$  is an ordinary function; the “|” in the middle of  $\pi(a|s)$  merely reminds us that it defines a probability distribution over  $a \in \mathcal{A}(s)$  for each  $s \in \mathcal{S}$ . Reinforcement learning methods specify how the agent’s policy is changed as a result of its experience.

*π belongs to the agent.  
agent has to learn how to apply actions.*

- Notations for deterministic and stochastic policies
- Invalid actions

All environments need to do is provide the next reward. The environment doesn't need to know  $v_{\pi}(s)$  and  $q_{\pi}(s, a)$ , agent B不知道, 加基它們需要知道  $v_{\pi}(s) \leq q_{\pi}(s)$ , 這點沒有 computation 來得快。

## Coursera video: Value Functions

- If the agent is following particular policy  $\pi$  from state  $s$ , what is the expected return the agent will get?
- **Definition of state value function**

The **value function** of a state  $s$  under a policy  $\pi$ , denoted  $v_{\pi}(s)$ , is the expected return when starting in  $s$  and following  $\pi$  thereafter. For MDPs, we can define  $v_{\pi}$  formally by

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right], \text{ for all } s \in \mathcal{S}, \quad (3.12)$$

*you can't change  $G_t$  because it's fixed*

*following a particular policy*

*( $G_t$  depends on policy  $\pi$ , not on the  $\pi$  that gave  $s$ )*

Saif Qasem to Everyone 1:30 PM
 

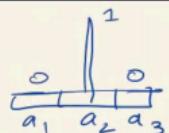
pi is the policy we follow so it determines the states we'll be in and the rewards we receive ...

- **Definition of action-value function**

$$q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]. \quad (3.13)$$

*given: doesn't matter how you arrive there.*

$$\begin{aligned}\pi(a_1|s) \\ \pi(a_2|s) = 1 \\ a_2 = \pi(s)\end{aligned}$$



$$\begin{aligned}
 \text{Reward} & \quad \left\{ E \left[ R_t \mid S_{t-1} = s \right] \right. \\
 \text{depends on} & \quad \left. = E \left[ R_{t+1} \mid S_t = s \right] \right. \\
 \text{the previous} & \\
 \text{state and} & \\
 \text{action} & 
 \end{aligned}$$

action. i. 这里比  $\pi$  more depend on  $\pi$  (policy) 但  $\pi$  是 given 的只有  $s$ . ii. To a. not depend on policy. 但  $\pi$  不知道 action (策略) 是什, 就谈不出 reward



## Coursera video: History of RL

- Rich S. Sutton
- Andy Barto
- A. Harry Klopf
- Philosophy of science: empiricism, logical positivism, Karl Popper
- Phenomenology: Hubert Dreyfus

# Coursera video: Bellman equations

- Bellman equations

$v_\pi(s) \stackrel{\text{def}}{=} \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_\pi(s')], \quad \text{for all } s \in \mathcal{S},$

the value function  $v_\pi(s)$  is the sum of the expected values of the immediate reward  $r$  and the discounted value of the next state  $s'$  for all possible actions  $a$ .

We can write the state value recursively.

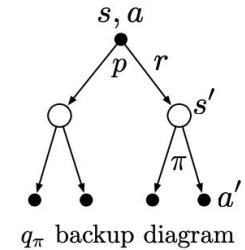
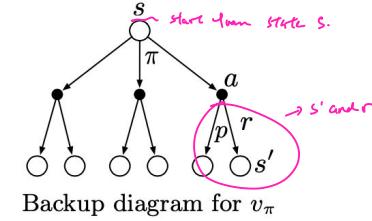
$\stackrel{=} : \text{Equality (equal-to) (logical equality)}$   
 $\stackrel{\text{def}}{=} : \text{by def } s\text{th is equal. (you can't redefine } s\text{th that has been defined, i.e. } v_\pi(s) \text{ is not } \stackrel{\text{def}}{=} \mathbb{E}_\pi[G_t | S_t=s].$

$v_\pi(s) \stackrel{\text{def}}{=} \mathbb{E}_\pi[G_t | S_t=s] \stackrel{\text{def}}{=} \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t=s \right], \quad \text{for all } s \in \mathcal{S}, \quad (3.12)$

is already defined.  
 logically equal to.

- Let's derive it next time

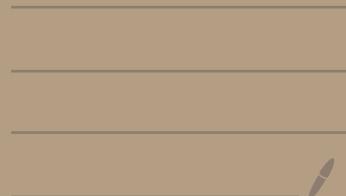
$$\begin{aligned}
 v_\pi(s) &\stackrel{\text{def}}{=} \mathbb{E}_\pi[G_t \mid S_t=s] \\
 &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t=s] \quad (\text{from (3.9)}) \\
 &= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t=s] \quad (4.3) \\
 &= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_\pi(s')], \quad (4.4)
 \end{aligned}$$



# Deriving Bellman Equations

---

$V_f$



# Bellman equation for $v_\pi$

$$v_\pi(s) = E_\pi \left[ a_t \mid s_t = s \right]$$

LOTE  
 LOTUS  
 MP  
 LB  
 $E[F(X|Y)]$   
 $E[X]$

$$\begin{aligned}
 &= E_\pi \left[ E_\pi \left[ a_t \mid s_t = s, A_t = a \right] \mid s_t = s \right] \\
 &= \sum_a E_\pi \left[ a_t \mid s_t = s, A_t = a \right] P(A_t = a \mid s_t = s) \\
 &= \sum_a \pi(a \mid s) q_\pi(s, a)
 \end{aligned}$$

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

$$\begin{aligned}
 &= E_\pi \left[ E_\pi \left[ a_t \mid s_t = s, A_t = a, R_{t+1}, S_{t+1} \right] \mid s_t = s, A_t = a \right] \\
 &= \sum_{s', r} E_\pi \left[ a_t \mid s_t = s, A_t = a, R_{t+1} = r, S_{t+1} = s' \right] p(s', r \mid s, a)
 \end{aligned}$$

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

$$\begin{aligned}
 & \downarrow \text{LE + MP} \\
 & = \sum_{s', r} p(s', r | s, a) \left[ r + \underbrace{\gamma \mathbb{E}_{\pi} \left[ \mathcal{G}_{41} \mid S_{t+1} = s' \right]}_{\text{By M.P. (G}_{t+1} \text{ 由前状态 } S_{t+1} \text{ 决定)}} \right] \\
 & \quad \xrightarrow{\text{用 LE 把这个 constant 提出来.}}
 \end{aligned}$$

$$\mathcal{V}_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \mathcal{V}_{\pi}(s') \right]$$

Show that

$$v_{\pi}(s) = r_{\pi}(s) + \gamma \sum_{a, s'} \pi(a|s) p(s'|s, a) v_{\pi}(s')$$

where

$$r_{\pi}(s) = \underset{\pi}{E} \left[ R_t \mid S_{t-1} = s \right]$$

$$= \sum_r r \frac{P(R_t = r \mid S_{t-1} = s)}{\pi} \quad \downarrow \text{LOT P}$$

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

$$v_{\pi}(s) = r(s) + \gamma \sum_{a, s'} \pi(a|s) p(s'|s, a) v_{\pi}(s')$$

depend on policy  $\pi$  but action  
 is not given.

$$\underbrace{r_{\pi}(s)}_{\text{expected reward of state } s.} = \mathbb{E}_{\pi} \left[ R_t \mid S_{t-1} = s \right]$$

↳ follow the belief

$$= \sum_r r \frac{P(R_t = r \mid S_{t-1} = s)}{\pi}$$

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

# Properties of expectations



$$\checkmark \text{ Linearity: } E[X + Y] = E[X] + E[Y]$$

Sum of two different R.V.

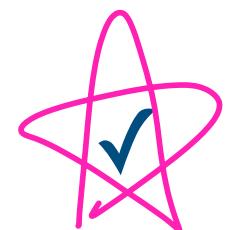


$$\checkmark \text{ Linearity: } E[aX] = aE[X]$$

→  $E[X]$  is R.V.,  $a$  is constant



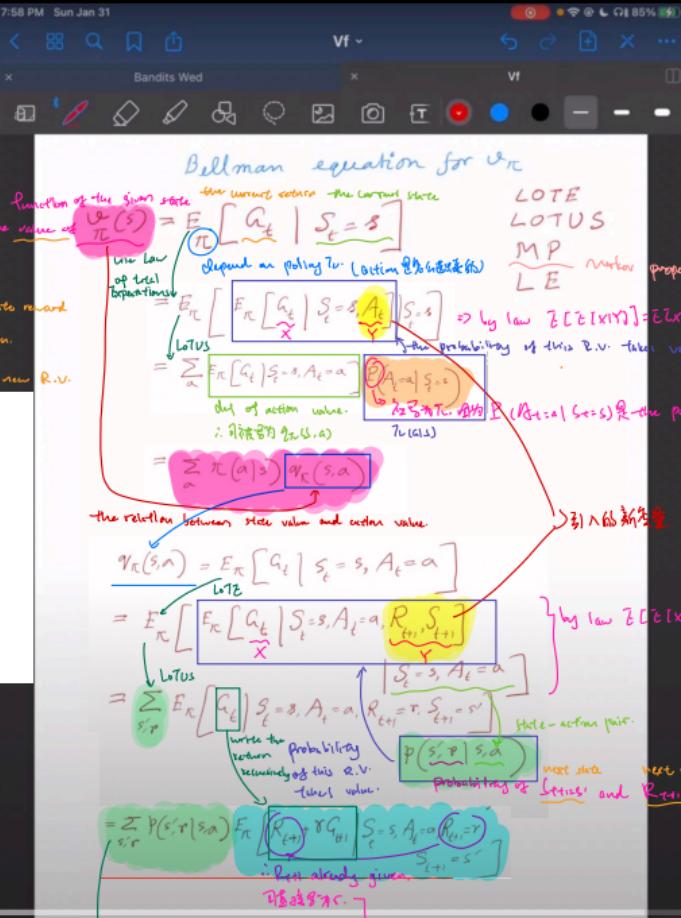
$$\checkmark \text{ Non-multiplicativity: } E[XY] \neq E[X] E[Y]$$



$$\text{Law of the unconscious statistician: } E[g(X)] = \sum_{x \in \mathcal{X}} g(x) P(X=x)$$

A R.V. which itself is a function of another R.V.

$$\begin{aligned} & \text{if } E[Y], \quad Y = g(X). \\ &= \sum_{y \in Y} y P(Y=y) \quad \text{probability of the taken value of the R.V. } Y \text{ is } y. \\ & \quad \text{a bigger sample space from which the R.V. } Y \text{ can take values.} \\ &= \sum g(x) P(g(X)=g(x)) \quad \text{probability of } g(X) \text{ taken value of } g(x). \text{ It is } g(x) \text{ is weighting.} \\ &= \sum_{x \in X} g(x) P(X=x) \end{aligned}$$



8:13 PM Sun Jan 31

Bandits Wed

VI

use linearity of expectation

$$= \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \mathbb{E}_{\pi} [q_{\pi} | s_{t+1} = s'] \right]$$

→ Note value of the next state at  $V_{\pi}(s')$

by Markov Property. 从该状态的奖励和状态转移概率决定下一状态的值

$$= \sum_{s', r} p(s', r | s, a) \left[ r + \gamma V_{\pi}(s') \right]$$

!! don't write policy here!! Because the transition dynamics of the environment doesn't depend on policy!!

take  $s$  as an input,  $\pi$  is the only function depend on  $s$  (because  $\pi$  is a function of  $s$ !!)

pink box:  $V_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \mathbb{E}_{\pi} [V_{\pi} | s'] \right]$

!! The action depends on the policy

!! The next state value also depends on the policy.

!!  $s$  is given !!  $\rightarrow$   $s$  is the only variable that should be bounded by the summation!!  $\pi$  is  $\sum_a$  !!

pink box:  $\sum_a \pi(a | s) \neq 1$ ,  $\pi$  is the only input in the L.H.S. 出现

pink box: bounded by  $\sum_a$  !!

pink box:  $\sum_{s', r} p(s', r | s, a) \neq 1$ ,  $s'$  有多个输入值!!

pink box: bounded by  $\sum_{s', r}$  !!

## Expectations & conditional expectations

- ✓ An expected value of a random variable is a weighted average of possible outcomes, where the weights are the probabilities of those outcomes

$$\mathbf{E}[X] = \sum_{x \in \mathcal{X}} x P(X=x)$$

- ✓ An expected value of a random variable conditional on another event is a weighted average of possible outcomes, where the weights are the conditional probabilities of those outcomes given the event

$$\mathbf{E}[X | Y=y] = \sum_{x \in \mathcal{X}} x P(X=x | Y=y)$$

- ✓ Expectation conditional on a random variable  $\mathbf{E}[X | Y]$  itself is a random variable, which is a function of another random variable  $Y$

$$v_{k+1}(s) \doteq \mathbb{E}_{\pi}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s]$$

this is not value function!

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_k(s')],$$

# Coursera video: Optimal policies and Optimal Value Functions

- Def'n:  $\pi \geq \pi'$  if and only if  $v_\pi(s) \geq v_{\pi'}(s)$  for all  $s \in \mathcal{S}$

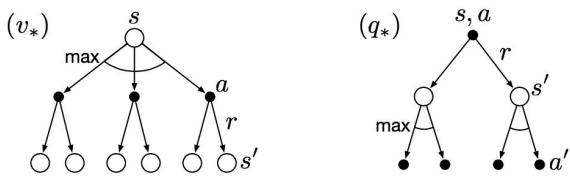
$$v_*(s) \doteq \max_{\pi} v_{\pi}(s) \Rightarrow \text{optimization of the optimal value.}$$

When you optimize, you define the optimal policy.  
Goal: maximize the  $v_*$  (s).

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

- Bellman optimality equations:

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]$$



$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]. \end{aligned}$$

Figure 3.4: Backup diagrams for  $v_*$  and  $q_*$

## Additional videos

- Notations for MDPs
- Deriving the Bellman equation with all steps shown